

TECHNICAL REPORT

COMPUTER SCIENCE

**The Right Stuff: Appropriate Mathematics for Evolutionary and
Development Biology**

**Chrystopher L Nehaniv
Gunter P Wagner**

June 1998 Report No 315

The Right Stuff:
Appropriate Mathematics for Evolutionary and
Developmental Biology

CHRISTOPHER L. NEHANIV
University of Aizu, Japan
University of Hertfordshire, U.K.

GÜNTER P. WAGNER
Yale University, U.S.A.

Full-Day Workshop at
the Sixth International Conference on Artificial Life
Life and Computation: the Boundaries are Changing
University of California, Los Angeles
26 June 1998



CONTENTS

Chrystopher L. Nehaniv and Günter P. Wagner,
The Right Stuff: Appropriate Mathematics for Evolutionary and Developmental Biology 1

INVITED LECTURES:

Giuseppe Pirillo, <i>Automatically Finding Genes in Genomes</i>	3
James P. Crutchfield, <i>The Evolutionary Unfolding of Complexity</i>	5
Erik van Nimwegen, <i>The Statistical Dynamics of Epochal Evolution</i>	6

ACCEPTED CONTRIBUTIONS:

Paulien Hogeweg, <i>Tangled Hierarchies and Tangled Spaces in Development and Evolution: Computational methods and biological insights</i>	7
Stefan Reimann, <i>Structure, Symmetry, and Decomposition</i>	11
Lionel Barnett, <i>Applying Markov Process Analysis to Evolutionary Systems</i>	15
Hamid Bolouri, Rod Adams, Stella George, and Alistair Rust, <i>Molecular self-organisation in a developmental model for the evolution of large-scale Artificial Neural Networks</i>	19
Misha Kapushesky, <i>Pattern formation by lateral inhibition</i>	22
Tim Taylor, <i>Using Bottom-Up Models to Investigate the Evolution of Life: Steps Towards an Improved Methodology</i>	23
Peter Dittrich, <i>Real Evolution in Artificial Chemistries</i>	27
C. L. Nehaniv and J. L. Rhodes, <i>Algebra, Evolution and Complexity of Living Systems</i> ...	32
D. Repsilber and F. Scholz, <i>Genetic Networks as a Model for the Regulatory Domain applied in Ecological Genetics</i>	33
Ehud Shapiro, Doron Lancet, Daniel Segre, <i>Molecular Transition Systems: A Computational embodiment for the Graded Autocatalysis Replication Domain (GARD) model</i>	36



The Right Stuff: Appropriate Mathematics for Evolutionary and Developmental Biology

CHRISTOPHER L. NEHANIV
University of Aizu, Japan
University of Hertfordshire, U.K.

GÜNTER P. WAGNER
Yale University, U.S.A.

We are looking for ‘the right stuff’, i.e. appropriate mathematical and computational tools/models for describing, studying, building or understanding fundamental aspects of natural living systems or living systems as-they-could-be (whether carbon-based, digital or otherwise) as opposed to inanimate systems.

Classical mathematical methods of population genetics tend to set out a fixed space of possibilities for the evolution of gene frequencies within a population. Unfortunately, by circumscribing the state-space at the outset, such an approach excludes the possibility of expressing change in developmental mechanisms or new evolutionary innovations such as body plans. While differential equation descriptions have proved crucial for understanding physics and chemistry and aspects of evolution, they seem to have largely failed as an appropriate language for some key aspects of biological systems. Living systems present special difficulties for such a mathematical treatment to particular problems of

- (1) death, damage, and development,
- (2) replication, inheritance and maintenance,
- (3) the relationship between genetic information and its realization via expression, and
- (4) the origin and evolution of biological complexity in populations of developing individuals.

The search for the right stuff strives to identify aspects special to living systems outside the scope of classical formal and conceptual tools, that can be treated formally with mathematical tools or computational models appropriate for natural (and artificial) biology.

Candidate areas where new, appropriate mathematical and computational approaches are needed include:

- Origins of Life
- Constructive Dynamical Systems
- Genetic Systems
- Algebraic Aspects of Evolutionary Change
- Self-Replicating / Self-Maintaining Systems
- Developmental Models and Evolutionary Change
- Evolution of Individuality
- Units of Evolution
- Body Plans
- Symbiogenesis
- Epigenetic Inheritance
- Modularity in Development
- Evolution and Maintenance of Sex
- Irreversibility in Biosystems and Development
- Algebraic Structure of Landscapes
- Symmetry and Decomposability
- Scaling Laws
- Community Construction

Acknowledgments

We are deeply grateful to our program committee members, who helped make the first right stuff workshop a success. The program committee consisted of Kurt Fleischer (Pixar Animation Studios, USA), Richard Michod (University of Arizona, USA), Melanie Mitchell (Santa Fe Institute, USA), Chrystopher L. Nehaniv (University of Aizu, Japan & University of Hertfordshire, UK), Joel R. Peck (University of Sussex, UK), Thomas S. Ray (ATR Human Information Processing Research Labs, Japan), Karl Sigmund (University of Vienna, Austria), and Günter Wagner (Yale University, USA)

Results about automatically finding genes in genome

GIUSEPPE PIRILLO
IAMI CNR
Viale Morgagni 67/A
50134 FIRENZE, Italy
and
Institut Gaspard Monge
Bâtiment IFI
Université de Marne-la-Vallée
2 rue de la Butte Verte
93160 NOISY-LE-GRAND
e-mail: pirillo@udini.math.unifi.it

The terminology concerning sequences (words) is that of [4]. Given an alphabet B , the *free monoid* (resp. *free semigroup*) over B is denoted by B^* (resp. B^+). An element of B is a *letter* (*nucleotide* or *base*). We use a four letter alphabet $B = \{A, C, G, T\}$ where A, C, G, T are the bases *Adenine*, *Cytosine*, *Guanine*, *Thymine*, respectively. An element of B^* is a *sequence* (*word*); a sequence of length 3 is a *trinucleotide*. The *empty sequence* is denoted by 1. A subset of B^* is a *language* (or a *gene population*).

Definition [3]. A language X in B^+ is a *code* if for $x_1, \dots, x_n, x'_1, \dots, x'_m$ in X an equality

$$x_1 \cdots x_n = x'_1 \cdots x'_m$$

implies $n = m$ and $x_i = x'_i$, $i = 1, 2, \dots, n$.

Definition [3]. A language X in B^+ is a *circular code* if for $x_1, \dots, x_n, x'_1, \dots, x'_m$ in X , $p \in B^*$ and $s \in B^+$, the equalities

$$sx_2 \cdots x_n p = x'_1 \cdots x'_m, \quad x_1 = ps$$

imply $n = m$, $p = 1$ and $x_i = x'_i$, $i = 1, 2, \dots, n$.

The theoretical biology team of Didier Arquès of University of Marne-la-Vallée has recently discovered an interesting partition of the 64 trinucleotides in three classes T_0 , T_1 and T_2 :

$T_0 = \{AAA, AAC, AAT, ACC, ATC, ATT, CAG, CTC, CTG, GAA, GAC, GAG, GAT, GCC, GGC, GGT, GTA, GTC, GTT, TAC, TTC, TTT\}$,

$T_1 = \{AAG, ACA, ACG, ACT, AGC, AGG, ATA, ATG, CCA, CCC, CCG, GCG, GTG, TAG, TCA, TCC, TCG, TCT, TGC, TTA, TTG\}$,

$T_2 = \{AGA, AGT, CAA, CAC, CAT, CCT, CGA, CGC, CGG, CGT, CTA, CTT, GCA, GCT, GGA, GGG, TAA, TAT, TGA, TGG, TGT\}$.

Let $X_0 = T_0 - \{AAA, TTT\}$, $X_1 = T_1 - \{CCC\}$, $X_2 = T_2 - \{GGG\}$. Then each of the sets X_0, X_1, X_2 contains 20 trinucleotides and has the remarkable properties of *complementarity* and *circularity* (see [1, 2]).

Moreover,

Theorem [1, 2]. *The sets X_0, X_1, X_2 are maximal circular codes.*

The purpose of this talk is to provide some information about those maximal circular codes that can be of interest in the study of the DNA and to present the progress that has been made in collaboration with D. Arquès on the preparation of the program of the automatic research of the genes in the genome.

REFERENCES

[1]. Arquès D. G. and Michel C. J., *A possible code in the genetic code*, STACS 95, (1995) 640-651.

[2]. Arquès D. G. and Michel C. J., *A complementary circular code in the protein coding genes*, J. theor. Biol., **182** (1996) 45-58.

[3]. Berstel J. and Perrin D., *Theory of codes*, Academic Press, London, 1985.

[4]. M. Lothaire, *Combinatorics on words*, Addison-Wesley, London, 1983.

[5]. Watson J. D. and Crick F. H. C., *A structure for deoxyribose nucleic acid*, Nature **171**, (1953) 737-738.

**THE EVOLUTIONARY UNFOLDING OF
COMPLEXITY**

JAMES P. CRUTCHFIELD
SANTA FE INSTITUTE & UNIVERSITY OF CALIFORNIA, BERKELEY,
U.S.A

**THE STATISTICAL DYNAMICS OF EPOCHAL
EVOLUTION**

ERIK VAN NIMWEGEN
SANTA FE INSTITUTE, U.S.A

Tangled Hierarchies and Tangled Spaces in Development and Evolution: Computational methods and biological insights.

P. Hogeweg

Theoretical Biology and Bioinformatics Group

Utrecht University, Padualaan 8, 3584CH Utrecht, the Netherlands.

email: ph@binf.bio.uu.nl

Evolution and development are preeminently multilevel processes, in which various entities can be recognized at different space as well as time scales. It is conceptually, mathematically and computationally convenient to treat such multilevel processes as hierarchical processes, which by separating time scales and/or space scales can be studied one hierarchical level at the time. Other levels then define the prerequisites and/or constraints on the behavior of the level under consideration, or the interactions (e.g. conflicts) between various independently defined levels is studied.

In this way ecological and evolutionary dynamics is mostly studied separately, and separately from e.g. genetic coding and from development.

We will argue that such approaches may 'miss the right stuff'. Instead we need to view these processes as 'tangled hierarchies' with overlapping and mutually defining space-time dynamics at several scales. Failing to do so may result in pseudo problems and/or pseudo solutions.

Two lines of research will be reviewed which demonstrate this position.

The first line is concerned with eco-evolutionary processes. Even in the most simple case of Predator-Prey interactions, we have shown that ecological and evolutionary time-scales interlock (van der Laan and Hogeweg 1995). In fact the ecological time-scale would be an order of magnitude larger when the evolutionary processes would be ignored (parameters held constant). Moreover, such an interlocked eco-evolutionary process throws new light on long studied ecological as well as evolutionary problems. With respect to ecology we show that stability of a relatively diverse ecosystem is maintained by the evolutionary process. With respect to evolution we show that the much debated issue of sympatric speciation occurs easily in such an eco-evolutionary system. Further entanglements occur when such eco-evolutionary processes are studied in space. Savill and Hogeweg (1997) showed that the dynamics of spatial patterns and patterns in phenotype space mutually determine each other.

In these studies genetic coding was fixed, and genotype phenotype mapping were left out of consideration. If we do take genetic coding into consideration, and allow the system to 'choose' its coding scheme (as is e.g. the case in Genetic programming) - the choice of coding scheme (i.e. the definition of genotype space) depends strongly on ecological and spatial processes (Pagie and Hogeweg 1998). Evolvability in locally interacting systems is improved and an interesting trade-off is found between generalisability (i.e. robustness to

environmental changes) and mutational stability (i.e. phenotypic robustness to genotypic change): local interactions lead robustness to environmental change, but sensitivity of genotypic change, while for global interactions it is the other way around. (see also Huynen and Hogeweg 1994) Moreover, we have demonstrated long term information integration in this system. This renders many 'what is this good for' questions examined on a single time scale meaningless.

The second line of research tries to exploit the entangled hierarchies as a research tool. We show that instead of being a nuisance, it can in fact help is to focus on interesting systems. In particular we propose to use evolutionary processes and the type of 'solution' it chooses when confronted with a very general problems for which many solutions exist, to map entanglements between levels. (see also Hogeweg 1998, for further discussion of this approach and some examples). Here we examine work in progress on development as an entangled multilevel process involving intra-cellular processes (gene-regulation leading to differential expression patterns) inter-cellular processes (cell sorting and morphogenesis through differential adhesion and/or chemotaxis) and the evolution thereof. We show an amazing morphogenetic versatility in such systems when as a fitness criterion simply 'number of cell types' is used.

Figure 1 gives an example of the development of an evolved creature. Striking is the pseudo isomorphic outgrowth of the creature. Cell differentiation is initiated by one 'maternal' signal at the first cell division, and remains fully reversible. This reversibility appears to stabilise the morphogenesis. The shape changes are triggered by cell death which occurs due to differential adhesion (cells are squeezed to death). The first 7 cell divisions are preprogrammed and occur simultaneously for all cells. Later cell divisions are triggered by stretching due to differential adhesion and cell death. Both cell division and cell death occur throughout the development (see upper curve of number of cells plot; the lower curve shows the number of cells of the same creature when later cell divisions are blocked: in that case morphogenesis is reversed due to cell loss, see lowest figure). Only through evolution, the feasibility of such an orchestrated development based on the simple process of differential adhesion can be demonstrated.

Finally, we will note that only in the light of entangled hierarchies, which arise automatically when no priori separations are imposed, we can pose the interesting evolutionary question if and how disentangling can take place in evolution.

References

- Hogeweg, P. (1998) On searching generic properties of non-generic phenomena: an approach to bioinformatic theory formation. Proceedings ALIFE6 in press
- Huynen, M.A. and Hogeweg, P. (1994) Pattern generation in molecular evolution: ex-

ploitation of the variation in RNA landscapes. *J. Mol.Evol* 39:71-79

Laan, J.D. van der, Hogeweg, P. (1995) Predator-prey coevolution: interactions among different time scales. *Proc. R. Soc. Lond. B.* 259: 35-42

Pagie, L. and Hogeweg, P. (1998) Evolving adaptability due to coevolving targets. *Evolutionary computation* (in press).

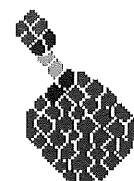
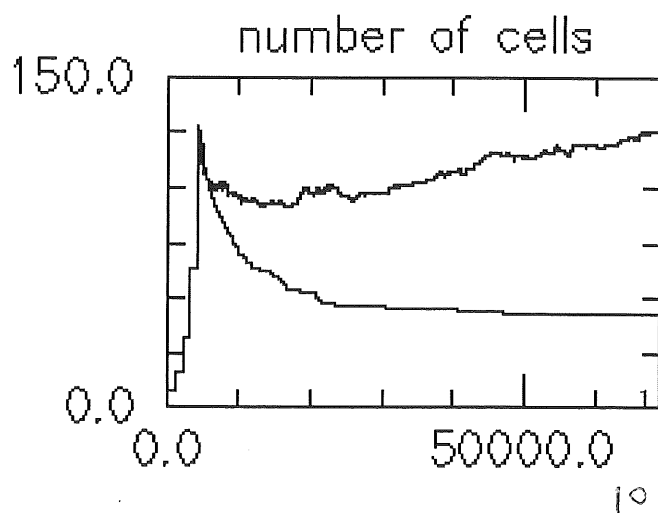
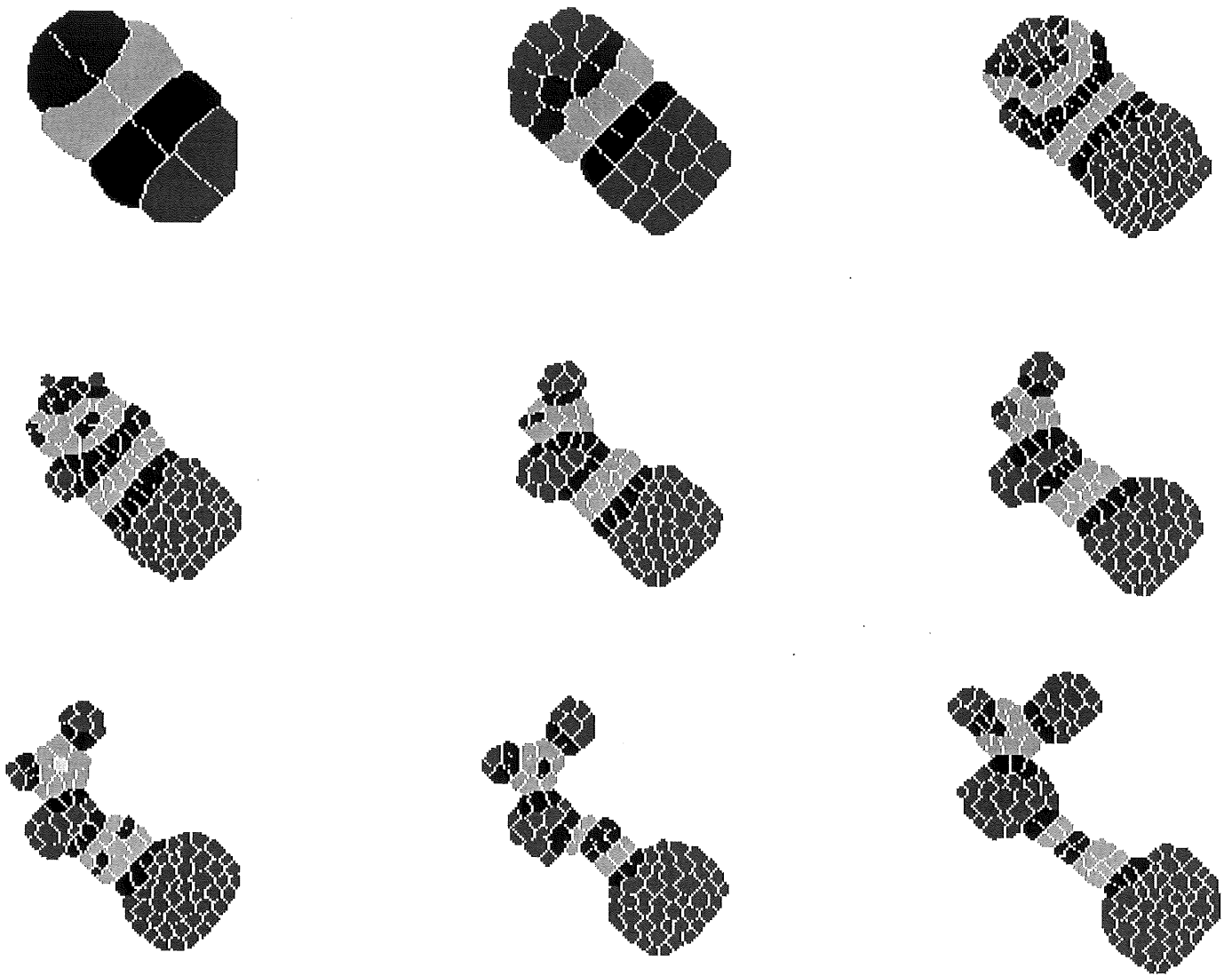
Savill, N. J. and Hogeweg, P. (1997) Evolutionary stagnation due to pattern-pattern interactions in a co-evolutionary predator-prey model. *Artificial Life* 3:81-100

Figure 1: Morphogenesis of an evolved creature.

differential gene expression is triggered by a maternal signal at the first cell division.

differential gene expression leads to differential adhesion.

differential adhesion leads to differential cell death (apoptosis) and differential cell growth



Structure, Symmetry, and Decomposition

STEFAN REIMANN

GMD

D-53754 Sankt Augustin

phone: (49) 2241 14 2068, fax: (49) 2241 14 2146,

e-mail: stefan.reimann@gmd.de

Biological systems certainly belong to the most successful complex systems known in nature. While each of their parts has a finite life duration and functions under quite restrictive constraints, the whole system exhibits an outstanding stability concerning a large amount of qualitatively as well as quantitatively different and disperse environmental conditions. This stability has at most two interrelated aspects. One is that, obviously, during their development the conditions have been maintained being necessary for the survival of the system. The second one concerns stability in terms of functional reliability: Biological systems appear to have the capability to systematically construct sufficiently good solutions not to only a single, but to a set of problems under a wide range of constraints. Today the question about the functional structures providing this kind of stability remains unanswered within a mathematical as well as in a technical framework.

Thus, being concerned with the structural features of biological systems and their related functional capabilities, the question arises about an appropriate mathematical framework. Obviously, biological systems are finite sets of interacting elements. Thereby, the elaboration of the structure and the function of a biological system are closely linked by the various modes of interactions between its components: While a particular structure determines the functional capabilities of a system, it evolves due to the embedding in its environment. Thus, interaction between its various subsystems plays a fundamental role concerning the structural and functional development of a biological system interrelating its structural organization with its functional capabilities. Accordingly, finiteness, discreteness, and interaction constitute very fundamental and general features of biological systems. These features provide the basis for the approach to be proposed in the following.

Biological systems as structured sets

The aim of the following is to show that the algebraic approach to be proposed is very natural as a fundamental concept and allows for investigating the interdependence of structural features and functional capabilities of biological systems.

Concerning mathematical formalization, it is convenient to start by considering a most simple but non-trivial setting of the problem which may provide a clear basis for further investigation. Thus, for simplicity, let us consider a (biological) system M and regard it as a finite set $M = \{m_1, \dots, m_N\}$, $N < \infty$, of identical elements m_i , which may be molecules, cells, organismic parts, or individuals. The structure of the system then is represented by the set of

all permutations that leave this set invariant. Together with the usual composition law, this set constitutes a group which synonymously is called the *symmetry group* of M or the *structure* of this set. Thus, we represent a biological system as the *structured set* (M, G) with G acting on M in a natural way by $g(m_1, \dots, m_N) := (m_{g^{-1}(1)}, \dots, m_{g^{-1}(N)})$. G can be thought of representing the morphology of the system. In fact, G can be represented as the adjacency matrix of a graph displaying the various mechanistical realizations of the interactions within the system. The spirit of this approach can be regarded as referring back to the work of H. Weyl [10]

In order to represent the functional properties of the system, consider the set of all the local states X_i of m_i , thereby specifying this element in terms of all of its possible states, and let $X = \prod_{i=1}^N X_i$ be the set of state of the entire system M . Accordingly, interaction within M is viewed as a mapping $\phi : X \rightarrow X$ where $X_i := \phi_i(X)$ is the state of the element m_i subject to the systemic interactions within M . We may therefore identify ϕ with the set of systemic interactions in the system. Since we restricted ourselves to identical elements only, we have $\phi_i = \phi_j$ for all i, j . We may further divide ϕ_i into two parts, i.e. $\phi_i = f_i + C_i$ where f_i is the restriction of ϕ_i to X_i , i.e. $f_i : X_i \rightarrow X_i$ thus representing the autonomous behaviour of the system m_i , while $C_i : X \rightarrow X_i$ describes the effect of the whole system to the state of m_i . Often, C_i displays a pairwise coupling of elements in the system. In this particular case, $C_i(x) := \sum_{j=1}^n c_{ij}(x_i, x_j)$, where the term $c_{ij}(x_i, x_j)$ means the influence of m_j to m_i . Consequently, a biological system can be represented as a tuple (X, ϕ) where ϕ is an operation on the global state space X representing the interaction within the system. If we consider pairwise interaction only, ϕ has to understood as the set of binary operations $\{C_i\}$. This case often is met in a system of pairwise interacting chemical elements and seems to be natural when regarding interaction concerning population dynamics, for example.

Note, that until now, we have obtained two different aspects: One is to regard the system as a tuple (M, G) , where M is the set of its elements and G is its global symmetry group acting by $gm := (m_{g^{-1}(1)}, \dots, m_{g^{-1}(N)})$. The other one is to represent the biological system by its state set X on which the function ϕ operates displaying the interaction between the elements m_i . A very natural way for joining these two aspects together is to consider a biological system as a tuple (X, G) consisting of state set X with the group G acting on it by $gx := (x_{g^{-1}(1)}, \dots, x_{g^{-1}(N)})$, thus regarding X as a G -space. The interaction represented by the mapping ϕ must be due to the global structure of the system, thus fulfilling $g\phi(x) = \phi(gx)$ for all $x \in X$ and $g \in G$. Such a mapping is called G -equivariant or a G -morphism. Accordingly, we describe a biological system an element of the category of G -sets (X, G) and G -morphisms $\phi : X \rightarrow X$ on it. This approach can be extended to additionally include symmetries of the sets of local states X_i by considering the *wreath product* of the local and the global symmetries of the system involved in the interaction of the system. For details see [3] and for further group theoretical material, for example, [6].

Decomposition and factorization

Considering a biological systems as a structured set (X, G) , we can easily define a (semi-)dynamical system φ on X by defining $\varphi^{t+s}(x) = \varphi^s \circ \varphi^t(x)$, $t, s \in \mathbf{N}, \mathbf{Z}$, or \mathbf{R} , subject to the initial condition $\varphi^0(x) = x$ for all $x \in X$ [2]. The resulting structure of the system depends on the properties of φ , of course. In particular, if φ is equivariant only according to a proper

subgroup H of G , $h\varphi(x) = \varphi(hx)$ for all $h \in H$, the effect of the mapping is due to break the symmetry G of the system down φ . This can be regarded as one aspect of the decomposition of a given system into various functionally distinguishable subsystems. A large amount of phenomena concerning pattern formation, in the temporal as well as in the spatial domain, seem to be due to such kind of symmetry breaking. Analysis of continuous dynamical systems was carried out by M. Golubitski, D. Schaeffer, I. Stewart, and others, for example [4], and has been related to dynamical phenomena in biological systems, including the coupling of cells, pattern formation, and the gaits of animals [1].

In the following, we consider a different aspect of decomposition, whose spirit is substantially algebraic and appears to be natural within the framework sketched above. Suppose that we consider a system (X, G) with an interaction given by $\phi : X \rightarrow X$ transforming a set of signals $Z \subseteq X$ into a set $Z' \subseteq X$, Z and Z' both having the symmetry group $H \leq G$. According to above, the system must be equivariant with respect to H fulfilling $h\phi(z) = \phi(hz)$ for all $z \in Z$ and all $h \in H$. Thus, we regard the system as the G -set (X, G) with ϕ , restricted to Z , being an H -morphism with range Z' . [8]. The question is: Can we find a decomposition of the system M into two parts, M_1 and M_2 , such that the composition of these two parts realize the same operation as the whole ensemble? One suggests that it may be possible to find smaller ensembles M_1 and M_2 such that the whole problem is reduced according to the symmetry of the set Z . It will turn out that one can, in fact, find a finite number of equivalent solutions of this problem. The argument is standard within universal algebra [7]. According to our setting, we know that the system operates as a H -morphism φ on the set Z . We further identify all signals z that are identical with respect to the functioning of the system, i.e. $z \equiv z' \pmod{\varphi}$ if and only if $\varphi(z) = \varphi(z')$. This defines a congruence relation on Z which may be denoted by ρ_φ . It is known that the image of the set Z under the action of φ is an H -invariant subspace of X' that is H -isomorphic to the set of these congruence classes, Z/ρ_φ . Accordingly, φ may be factorized into a projection $\varphi_1 : Z \rightarrow Z/\rho_\varphi$ and an injective mapping $\varphi_2 : Z/\rho_\varphi \rightarrow Z'$, defined by $\rho_\varphi \mapsto \varphi(z)$ such that $\varphi(z) = \varphi_2 \circ \varphi_1(z)$ for all signals $z \in Z$. Thus we can regard (M_1, φ_1) and (M_2, φ_2) as a decomposition of the system M . Note that this decomposition is unique only up to an isomorphism. For a related treatment of artificial networks and the question about the minimal size of the partial ensembles concerning the auto-associator problem see [9].

In summary:

Due to its inherent discreteness (and finiteness), we started by regarding each biological system as a finite set M of elements m_i . Moreover, due to the systemic interactions in M we proceeded to regard a biological system as an algebraic object (X, G) where X is the set of its global states and G is its symmetry group together with a G -equivariant mapping ϕ representing the interaction in M according to the global symmetry of the system G . Thus, the category of G -sets and G -morphisms naturally arises from basic features of biological systems. The time-evolution of a biological system is represented as the its trajectory of a dynamical system (X, φ) defined on the (discrete) structured set (X, G) and to study its dynamics with respect to this algebraic structure. Thereby, the trajectory of the system will have a symmetry reflecting the properties of the φ corresponding to the global symmetry of the system G . The symmetry breaking effect of φ thus leads to a decomposition of the whole system into parts which are defined according to the resulting symmetry of the system. An other aspect of decomposition

concerns the question of whether one can find subsystems of an ensemble such that their composition preserves the function of the entire system. Within the algebraic framework proposed, the answer to this question is closely related to the factorization of the mapping representing the functioning of the entire systems into two mappings each representing the functioning of its subsystems. Factorization is a common tool within universal algebra and provides a large amount of deep insights into the structure of the algebraic situation considered. Thus, roughly speaking, an algebraic approach seems to be natural with respect to very fundamental features of biological system as considered above and, in particular, serves as a powerful mathematical conception for describing and discussing the structure of biological systems. Additionally, even from a conceptual point of view, algebra seems to provide a very natural tool for investigating biological systems in that algebra is constructive in nature in that it studies the formation of larger structures (groups, rings, . . . , vectorspaces) and their corresponding properties as being established by composition rules from more elementary objects. Thus, by being locally vague but globally rigid [5], algebra may provide a conceptual and methodological frame to investigate the properties of large systems constituted by a number of elementary interaction modes of its components.

References

- [1] J.J. Collins, I.N. Stewart, "Coupled nonlinear oscillators and the symmetries of animal gaits", *J. Nonlinear Science*, **3**, 349-392, 1993;
- [2] M.J. Field, "Equivariant dynamics", in: *Contemporary mathematics* **56**, 69-96, 1986;
- [3] M. Golubitski, I. Stewart, B. Dionne, "Coupled cells: wreath products and direct products", in: *Dynamics, Bifurcation and Symmetry*, P. Chossat (ed.), 127-138, 1994;
- [4] M. Golubitski, D.G. Schaeffer, "Singularities and groups in Bifurcation Theory", *Applied Mathematical Science* **51**, Vol. 1 + 2, 1985;
- [5] R.E. Kalman, "Remarks on mathematical brain models", in: *Biogenesis, Evolution, Homeostasis*, A. Locker (ed.), 173-179, 1973;
- [6] J.D.P. Meldrum, *Wreath products of groups and semigroups, Pitman monographs and surveys in pure and applied mathematics* **74**, Harlow, Essex: Longman Group LTD, 1995;
- [7] P. M. Neumann, G. A. Stoy, E.C. Thompson, *Groups and Geometry*, Oxford University Press, 1995.
- [8] S. Reimann, "On the design of artificial auto-associative networks", *Neural Networks* (in print), 1998;
- [9] S. Reimann, W. Terhalle, "On the functioning of an encoder", (submitted), 1998;
- [10] H. Weyl, *Symmetrie*, Birkhäuser, Basel, 1955.

Collapsing the State Space

Applying Markov Analysis to Evolutionary Systems

Lionel Barnett

Centre for Computational Neuroscience and Robotics
Centre for the Study of Evolution
Department of Cognitive and Computing Sciences
University of Sussex, Brighton BN1 9QH, UK
lionelb@cogs.susx.ac.uk

Abstract

Evolutionary systems may often be accurately modelled by Markov processes. However the state space invariably turns out to be vast and multi-dimensional, thus limiting the application of Markov theory to broad abstraction rather than specific problems. One notable exception is the analysis of error thresholds in finite populations by (Nowak & Schuster 1989), where a (seemingly unjustifiable) approximation is made to “collapse” the state space, reducing the problem to an analytically tractable form. In this paper we outline Nowak and Schuster’s analysis and discuss the methodology of their approach.

Error thresholds for Finite Populations

(Nowak & Schuster 1989) investigated the extension of established results from “quasispecies” theory (Eigen *et al.* 1989) on error thresholds for infinite populations to finite populations. The basic problem is as follows: we are given a “single spike” fitness landscape of binary genotypes of sequence length v . All genotypes have fitness 1 except for the genotype consisting of all zero’s (the *master* genotype or *optimum*), which has fitness $\sigma > 1$. Genotypes Hamming distance α from the optimum are said to belong to the *error class* Γ_α - the Γ_α for $\alpha > 1$ constitute the *error tail*.

Consider a fixed-size population of N genotypes evolving via fitness-proportional selection¹ and mutation at a per-locus rate of μ ($0 \leq \mu \leq \frac{1}{2}$). There is no recombination. The observed long-term behaviour of such a system is as follows: at low mutation rates the population clusters around the optimum (Fig 1a). At higher mutation rates more genotypes are to be found at a small (Hamming) distance from the optimum (Fig 1b). Beyond a critical mutation rate, the *error threshold*, the population “loses” the optimum altogether and drifts randomly

around the landscape² (Fig 1c). In the infinite population limit the error threshold may be calculated from quasispecies theory using perturbation methods (Eigen *et al.* 1989). For finite populations the error threshold is less easy to define, let alone calculate. Nevertheless, there is still a sharp transition between long-term behaviours in the sense that the transition (for reasonably long sequence length v) occurs within a very small range of mutation rates.

To analyse the transition we must examine the distribution (over time) of the number of optimum genotypes, $\pi_i \equiv \mathbf{P}(\hat{X} = i)$, where the random variable \hat{X} represents the number of optimum genotypes in the long term. (Nowak & Schuster 1989) found that at low mutation rates the distribution peaks at some characteristic value of i (Fig 2a). At high mutation rates the distribution decreases monotonically from $i = 0$ (Fig 2c). At intermediate mutation rates the distribution develops a second peak at $i = 0$ (Fig 2b). The authors then *define* the error threshold to be that value of μ at which the distribution changes from monotone decreasing to one with a peak at $i > 0$.

Now it is apparent that we could calculate the distribution π_i if it were true that the random variables $X(t)$ representing the number of optimum genotypes at (discrete or continuous) time t constituted a Markov process. The distribution π_i would then be simply the stationary distribution of the process (Karlin & Taylor 1975). However it is clear that the Markov property does not hold, for the following reason: while the probability that a genotype be selected for replication depends only on whether it is of the optimum type or in the error tail, the probability that a genotype in the error tail “back-mutates” to the optimum type depends, in

¹ The exact selection algorithm in effect does not alter the qualitative phenomena; thus selection may be roulette-wheel, tournament, etc. as long as the *expected* number of offspring of a genotype is proportional to its fitness. The algorithm may, in addition, be discrete or continuous time.

² This implies that the distribution of number of genotypes among the error classes is *binomial*, as the α ’th error class occupies a fraction $2^{-v} \binom{v}{\alpha}$ of the landscape; cf. Fig 1c.

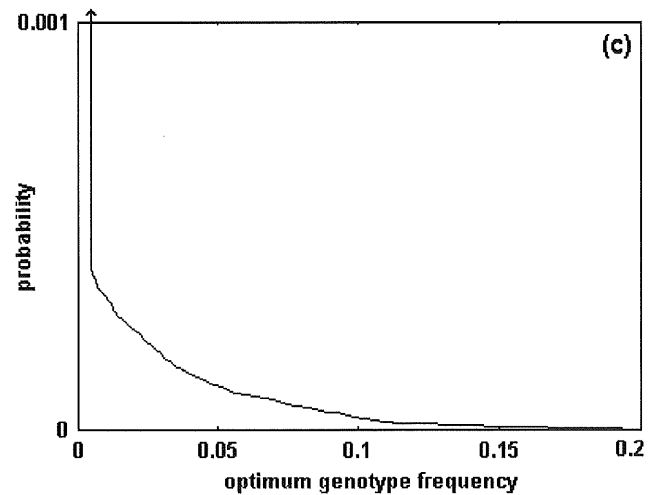
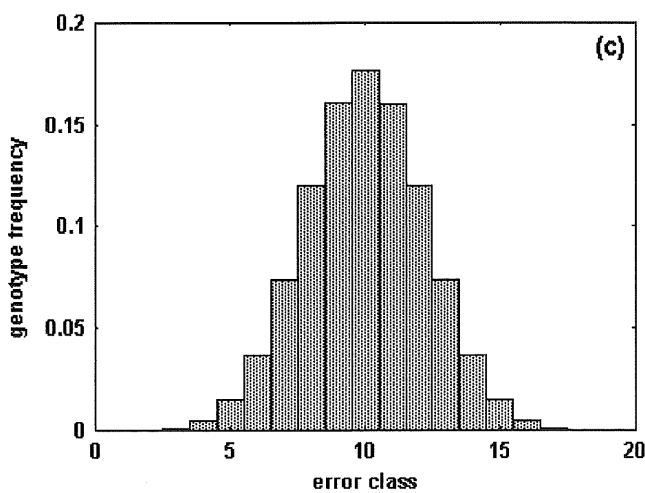
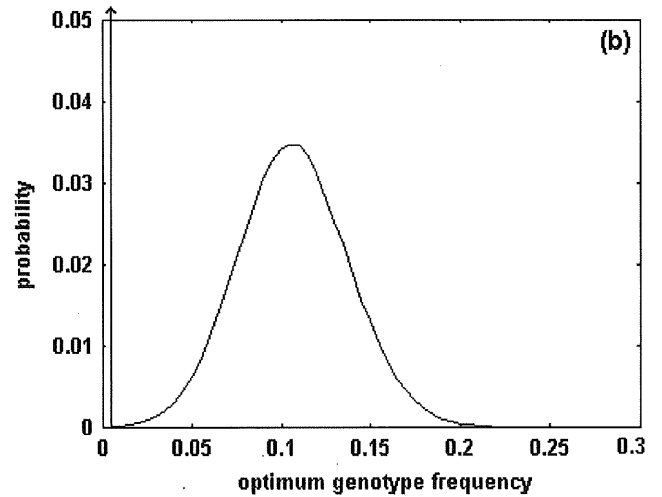
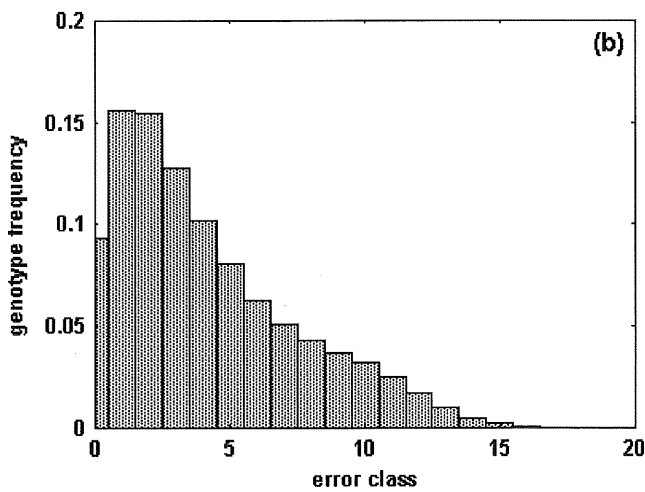
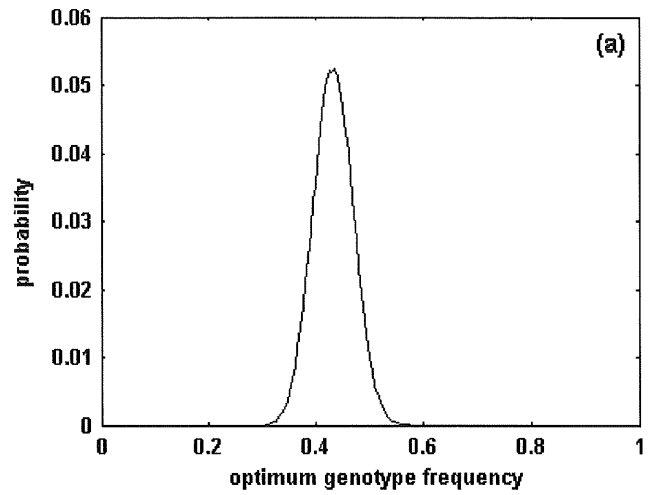
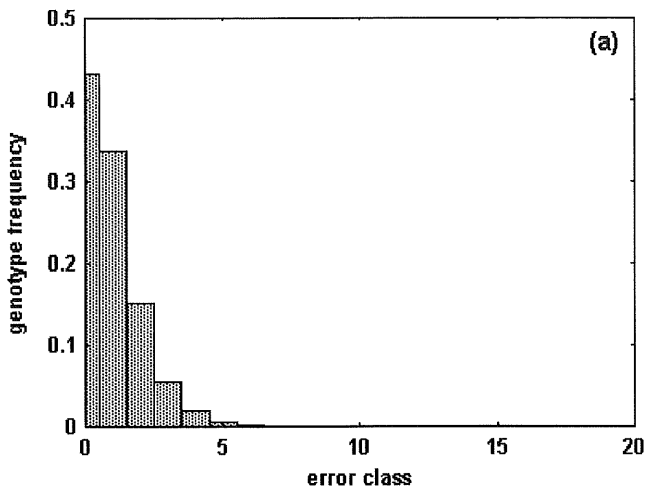


Fig 1. Distribution of genotype frequency over error classes for $\sigma = 5$, $v = 20$, $N = 200$ and per-locus mutation rates (a) $\mu = 0.03$, (b) $\mu = 0.061$ and (c) $\mu = 0.07$. Simulation was over 100,000 generations of a roulette-wheel fitness-proportional selection algorithm.

Fig 2. Stationary probability distribution π_i of optimum (error class 0) genotype frequency for the same simulations as in Fig 1. The error threshold is approximately $\mu = 0.063$.

addition, on how many 1's it has; i.e. its error class. Nowak and Schuster address this issue by making what on the face of it is an unjustifiable assumption: that the distribution of genotypes is *uniform* in the error tail. Under this assumption both selection and mutation probabilities depend only on the number of optimum genotypes and the Markov property holds. The authors then carry through the Markov analysis for a particular evolutionary algorithm [a continuous birth-death model from population genetics - see (Moran 1958)] for which the stationary distribution is explicitly solvable, thus arriving at what turns out to be a very accurate estimate for the error threshold.

So why should this scheme work at all? It is clear from Figs. 1a and 1b that the assumption of a uniform distribution of genotypes in the error tail is manifestly false. A point apparently missed by the authors, though, is that at *high* mutation rates, when the optimum is "lost", the uniform distribution assumption is actually quite sound, as is evidenced by Fig. 1c. This may explain the accuracy of their result to some degree; as long as the error threshold is approached "from above" the assumption holds good. Another point worth noting is that for reasonably long sequence length v the probability of back-mutation from the error tail to the optimum becomes small (of the order of 2^{-v}), even for genotypes a small Hamming distance from the optimum. Then ignoring back mutation entirely is a reasonable approximation and the Markov property holds without any further assumptions.

Another possible approach might be as follows: it is possible to calculate (numerically at least) the distribution of genotypes in the infinite population limit (van Nimwegen *et. al.* 1997). This limiting distribution is quite a good approximation to the finite population case for reasonably large populations, although it is not clear what "reasonably" large might mean. Thus, rather than assuming a uniform distribution of genotypes in the error tail we could assume instead the infinite population limit. Preliminary tests by this author suggest that this can give a significantly more accurate approximation to the stationary distribution of the optimum than the cruder uniform distribution assumption, particularly near the error threshold.

Methodological Issues

It is worth examining how the procedure outlined above tackles the issue of state space size and dimension. The full state space for the problem of a fixed-size population of size N evolving on a fitness landscape is the set of all possible populations. A population is naturally identified with an integer vector $\mathbf{n} = (n_g)$ indexed by all possible

genotypes g , where n_g represents the number of copies of genotype g in the population. The n_g must satisfy $n_g \geq 0 \forall g$ and $\sum_g n_g = N$. The state space is thus vast and multi-dimensional; if sequence length is v then the cardinality of the state space is $\binom{v+N}{v}$ which is of the

order of N^v for $N \gg v$. The crucial point is that if all we are interested in is the error threshold, the only quantity we need to know is the stationary probability distribution of the frequency of optimum genotypes. Now since the single-spike landscape is "isotropic" with respect to the optimum genotype we can immediately "collapse" the state space into the frequencies of genotypes in the error classes without losing either the Markov property or the quantity we wish to measure. This is possible because mutation and selection probabilities (and hence the transition probabilities of the Markov process) depend only on error class. Thus our state space may be immediately reduced to the set of vectors $\mathbf{n} = (n_\alpha)$ indexed now by the error classes α . [If recombination were present this would no longer be true - see below.] Note that thus lose all information as to the distribution of genotypes *within* error classes - but we do not need this information for the problem at hand! The state space is then reduced still further by (cautious) approximation to an analytically tractable 1-dimensional space.

Another point that may be overlooked in the quest for quantitative results is that even if various approximations introduce quantitative inaccuracies (as they do to some degree in Nowak and Schuster's analysis), the *qualitative* picture may still hold up. Thus valuable insights may be gained into the dynamical behaviour of an evolutionary system by the introduction of simplifying assumptions; this is certainly the case for Nowak and Schuster's analysis of error thresholds.

As a further case in point this author [in preparation] has extended Nowak and Schuster's analysis to include recombination, revealing a rich and often surprising range of dynamics. This necessitated the introduction of further (quantitatively unjustifiable) assumptions, specifically because the Markov property does not hold even for recombination *within* error classes. Comparing analytical results with simulations, however, reveals that the approximation retains almost all qualitative features of interest.

Finally, it would seem to be feasible to extend these principles to the analysis of evolution on more complex landscapes, particularly if they feature analogues of the error classes. A comparable approach (although not for

the purposes of Markov analysis) can be found in (van Nimwegen *et. al.* 1997).

Thus we might define a partition $\{\Gamma_\alpha | \alpha \in A\}$ of a fitness landscape with fitness function $f(g)$ and (stochastic) mutation operator M_μ to be a *Markov Partition* if it satisfies:

$$\forall \alpha \in A, \forall g, g' \in \Gamma_\alpha \forall \beta \in A \text{ and } \forall \mu (0 \leq \mu \leq 0.5):$$

$$\text{MP1} \quad f(g') = f(g)$$

$$\text{MP2} \quad P(M_\mu(g') \in \Gamma_\beta) = P(M_\mu(g) \in \Gamma_\beta)$$

The quantities:

$$f_\alpha \equiv f(g) \text{ for some } g \in \Gamma_\alpha$$

and:

$$M_{\alpha\beta}(\mu) \equiv P(M_\mu(g) \in \Gamma_\beta | g \in \Gamma_\alpha)$$

are then well-defined. An evolutionary process may then be considered as a Markov process on the state space of all integer vectors $\mathbf{n} = (n_\alpha)$ with $n_\alpha \geq 0 \forall \alpha$ (and $\sum_\alpha n_\alpha = N$ for fixed population size N), the transition probabilities being determined (for the particular evolutionary algorithm employed) by the f_α and $M_{\alpha\beta}(\mu)$.

The *coarsest* such partition yields the smallest and most manageable state space. If possible, depending on specific features of the fitness landscape, a coarser partition might be found by relaxing conditions **MP1** and/or **MP2** to hold approximately. A further condition to cover recombination could also be defined, although it seems doubtful that useful partitions could be found which would respect such a condition exactly.

MP1 corresponds to the statement that the Γ_α are *neutral subsets* (Barnett 1998) of the fitness landscape. A particular case of interest is where they constitute the *neutral networks* (Forst *et. al.* 1995, Huynen *et. al.* 1996, Barnett 1989) of the landscape i.e. maximal connected neutral subsets. In particular, when the neutral networks "percolate" the landscape (Forst *et. al.* 1995) it seems reasonable that **MP2** might be expected to hold to some approximation.

Acknowledgements

This paper arose out of discussions of the Neutral Networks group at the Sussex University CCNR. The author would like to thank all those who took part and Inman Harvey in particular.

References

- Barnett, L. 1998. *Ruggedness and Neutrality - the NKp Family of Fitness Landscapes*. Proc. 6th Intl. Conf. on Artificial Life (in print).
- Eigen, M., McCaskill, J. & Schuster, P. 1989. *The Molecular Quasispecies*. Adv. Chem. Phys. 75:149-263.
- Forst, C.V., Reidys, C. & Weber, J. 1995. *Neutral Networks as Model-Landscapes for RNA Secondary-Structure Folding-Landscapes*. Lecture Notes in Artificial Intelligence, vol. 929: Advances in Artificial Life (Morán, F., Moreno, A., Merelo, J.J. & Chacón eds.), Springer-Verlag, Berlin.
- Huynen, M.A., Stadler, P.F. & Fontana, W. 1996. *Smoothness Within Ruggedness: The Role of Neutrality in Adaptation*. Proc. Natl. Acad. Sci. (USA) 93:397-401.
- Karlin, S. & Taylor, H.M. 1975 *A First Course in Stochastic Processes (2nd ed.)*. Academic Press, New York.
- Moran, P.A.P. 1958 *The Effect of Selection in a Haploid Genetic population*. Proc. Camb. Phil. Soc. 54: 463-465.
- van Nimwegen, E., Crutchfield, J.P. & Mitchell, M. 1997. *Statistical Dynamics of the Royal Road Genetic Algorithm*. Santa Fe Institute Pre-print 97-04-035, Santa Fe, NM, USA.
- Nowak, M. & Schuster, P. 1989. *Error Thresholds of Replication in Finite Populations - Mutation Frequencies and the Onset of Muller's Ratchet*. J. Theor. Biol. 137:375-395.

Evolving large-scale Artificial Neural Networks

Hamid Bolouri

Engineering R & D Centre, University of Hertfordshire, UK
Division of Biology, California Institute of Technology, USA

Rod Adams, Stella George, Alistair Rust
Department of Computer Science, University of Hertfordshire

Introduction

One of the attractions of Artificial Neural Networks (ANNs) has been the possibility of designing intelligent systems capable of optimising their functionality according to application requirements.

Adapting the architecture of an ANN (number of neurons, neuron function and connectivity pattern) to a given application can be viewed as a combinatorial optimisation problem. For sophisticated applications, the problem tends to be high-dimensional and highly nonlinear. Direct applications of current combinatorial optimisation methods to such problems tend to be unacceptably inefficient.

To date, research in self-adapting ANNs has been limited to the design of unsupervised learning algorithms, and the use of algorithmic (generative) techniques to 'grow' or 'prune' neurons and their connections. Unsupervised learning, and generative algorithms both make strong assumptions about the ANN architecture and the characteristics of the target application domain.

To alleviate these restrictions, a number of researchers have exploited evolutionary methods to design adaptive generative algorithms. The generative algorithm is encoded in a genome. Its execution (the mapping from genotype to phenotype) mimics embryonic development. Current implementations of this approach limit the architectural search space through apriori assumptions and algorithmic restrictions. However, restricting the range of genotype-phenotype mappings can hinder rather than help the optimisation process: non optimal mappings result in more nonlinear search spaces which are more difficult to search. We argue that it is possible to evolve unconstrained optimal mappings through the use of two strategies:

- 1) Cellular (neuron) characteristics should be defined in terms of interacting molecular processes. These molecular interactions result in nonlinear genotype-phenotype mappings. Their evolutionary optimisation reduces the degree of nonlinearity in the genotype search space and simplifies the optimisation process.
- 2) Evolution should be staged, mimicking the process of speciation. Specifically, evolution should start with the simplest set of generative rules, search for the best achievable, then add new (more sophisticated) rules (mimicking the emergence of a new species), and search again; repeating this procedure until satisfactory phenotypes emerge.

We are currently investigating ways of exploiting these strategies to robustly evolve large scale ANNs, and present simulation results demonstrating molecular, staged evolution of a simple edge detecting retina.

Methodology

Our strategy is based on a molecular model of the evolution of embryonic development (see figures 1 & 2 for a schematic overview) with the following characteristics:

- 1) Neurons are modelled as cells. Each distinct cell type is defined by the interactions of a specific subset of genes within the genome.
- 2) The genome, and the evolutionary operators mutation and cross-over, are defined at the molecular level (i.e. a much finer level of detail than cell characteristics such as receptive field size).

- 3) The translation from genetic encoding to neural network (genotype to phenotype mapping) is performed using sequences of operations defined by individual genes. These operations mimic gene regulation (the sequential activation and interaction of genes) and embryonic development in biological systems.
- 4) Although cells of the same type (ensembles of neurons) share the same genetic description, individual neurons within an ensemble may differ from each other if differences in cell-cell or cell-environment interactions lead to different developmental histories.
- 5) Hierarchic network structures can be defined with nested loops of gene interactions. Thus, similar cells in different parts of a network may utilise common segments of the genome to define the parts they have in common.
- 6) A gene (more accurately a gene product) has two functional aspects: the definition of what it will interact with, and the definition of the nature of its interactions¹ (affinity and type respectively).
- 7) The genome is defined as a (variable) number of distinct chromosomes each comprising a set of genes which may, in principle, interact with each other. During cross-over only alike chromosomes can exchange parts.
- 8) The neural system is evolved gradually and in stages corresponding to evolution within species and the emergence of new species.
- 9) Speciation is achieved by allowing the size of the genome to vary through gene duplications and deletions.
- 10) Evolution starts by using only the simplest of developmental programs. For instance, neuron connectivity is specified using only intrinsic growth rules. When this type of genome is deemed to have been adequately optimised, additional genes describing potentially more powerful developmental processes are added to the genome and a new cycle of evolutionary optimisation starts.
- 11) No information other than the performance of the evolved networks is used to guide the evolutionary process.

¹ In biological systems, a single gene product (a polypeptide) also has two aspects. Firstly, it will only interact with specific other gene products. Affinity is defined by molecular characteristics such as geometric shape, distribution and type of electrostatic and covalent bonds, and the presence of particular metallic ions such as copper or zinc at particular locations. Secondly, these characteristics also combine to determine *how* a polypeptide interacts with other molecules. For instance, a gene product may act exclusively as a catalyst (an enzyme), or it may form a chemical product with another molecule (e.g. in signal transduction), or it may be involved in a physical/mechanical interaction with other molecules and form a structural part of a cell (e.g. an ion channel, or part of a lipid bilayer).

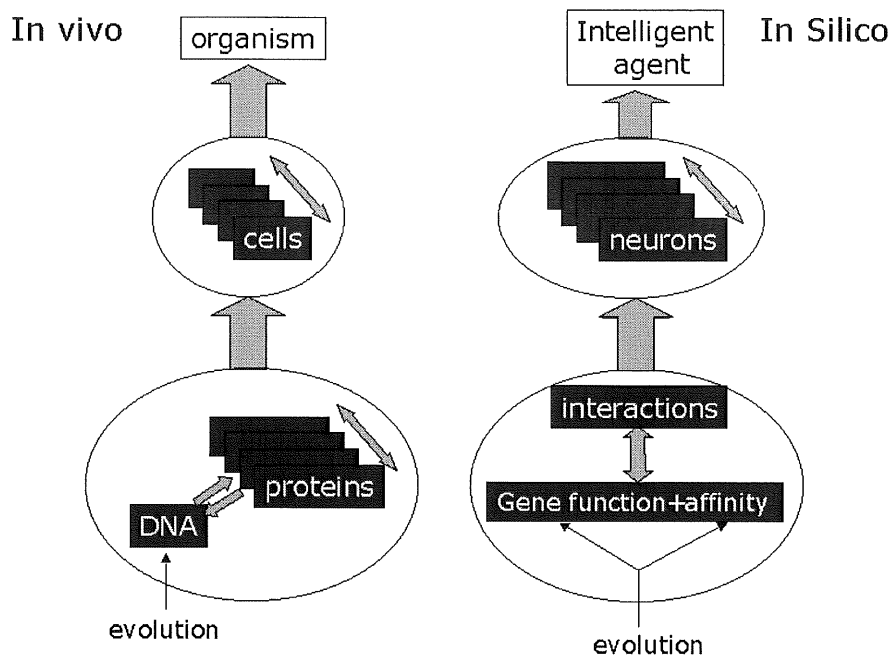


Figure 1, overview of proposed evolutionary process

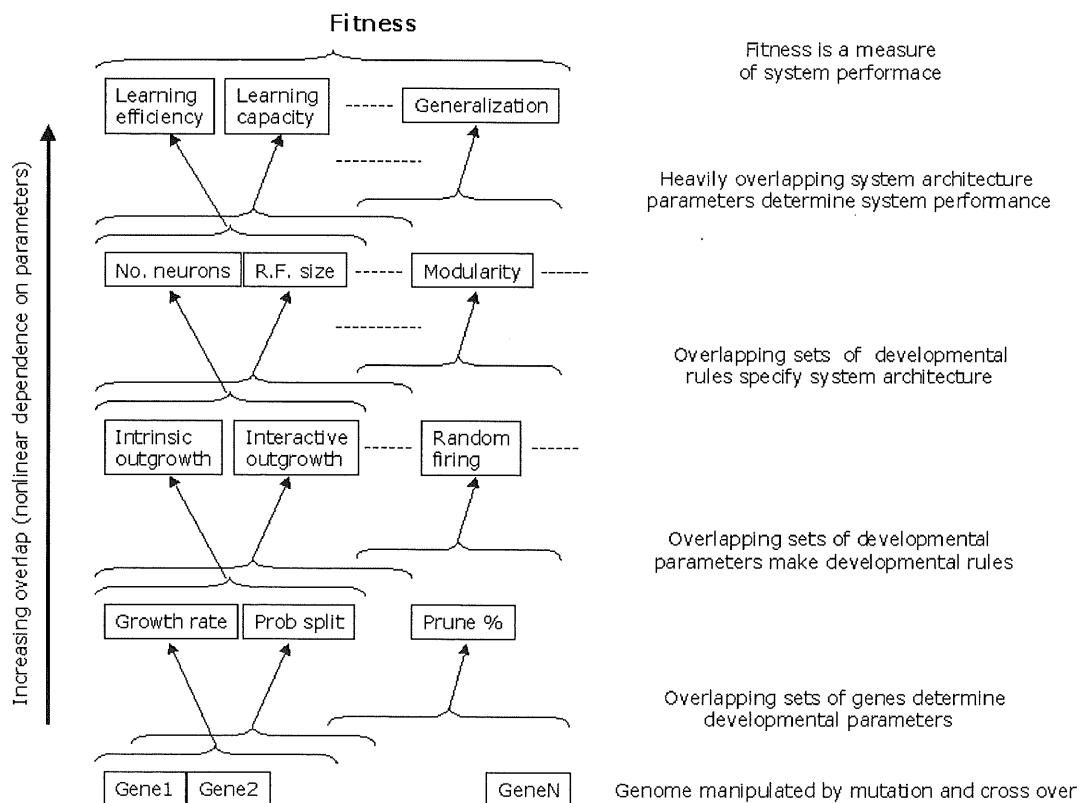


Figure 2, Proposed nonlinear mapping between the genetic code and the objective function

PATTERN FORMATION BY LATERAL INHIBITION

MISHA KAPUSHESKY
CORNELL UNIVERSITY, USA

Abstract

In many developing tissues, adjacent cells diverge in character so as to create a fine-grained pattern of cells in contrasting states of differentiation. It has been proposed that such patterns can be generated through lateral inhibition - a type of cell-cell interaction whereby a cell that adopts a particular fate inhibits its immediate neighbours from doing likewise. Lateral inhibition is well documented in flies, worms and vertebrates. In all of these organisms, the transmembrane proteins Notch and Delta (or their homologues) have been identified as mediators of the interaction - Notch as receptor, Delta as its ligand on adjacent cells. However, it is not clear under precisely what conditions the Delta-Notch mechanism of lateral inhibition can generate the observed types of pattern, or indeed whether this mechanism is capable of generating such patterns by itself. Here we construct and analyse a simple and general mathematical model of such contact-mediated lateral inhibition. In accordance with experimental data, the model postulates that receipt of inhibition (i.e. activation of Notch) diminishes the ability to deliver inhibition (i.e. to produce active Delta). This gives rise to a feedback loop that can amplify differences between adjacent cells. We investigate the pattern-forming potential and temporal behaviour of this model both analytically and through numerical simulation. Inhomogeneities are self-amplifying and develop without need of any other machinery, provided the feedback is sufficiently strong. For a wide range of initial and boundary conditions, the model generates fine-grained patterns similar to those observed in living systems.

Using Bottom-Up Models to Investigate the Evolution of Life: Steps Towards an Improved Methodology

Tim Taylor

Department of Artificial Intelligence, University of Edinburgh
5 Forrest Hill, Edinburgh EH1 2QL, U.K.

timt@dai.ed.ac.uk

1 Introduction

Perhaps one of the few features shared by most artificial life approaches is that a phenomenon observed in biological life is studied by constructing a *bottom-up* model, in which a number of low-level components and interactions are explicitly encoded, and one or more higher-level phenomena are expected to emerge. While this is a perfectly valid approach, one has to be careful about how the model is constructed if it is to bring any scientific insight to bear on the phenomenon in question. Too many (but not all) studies of artificial life (my own included) have adopted a sloppy approach in the past, and this has meant that the field of artificial life has not contributed as much as it might have done to broader areas of scientific knowledge.

This paper highlights some areas of general methodology which should be carefully considered when designing a bottom-up simulation for scientific experimentation, and also suggests some considerations that are specifically relevant to A-life models designed to investigate the evolution of life. This is not, of course, the first time that concerns have been raised about the methodology of A-life, and very little is said here that has not been said before (see, for example, [4], [1], [5]). However, I believe that much current A-life work still suffers from poor methodology, and that it is therefore important to stress such issues at every available opportunity.

2 General Considerations

2.1 Explicit Assumptions and Predictions

The bottom-up approach to studying high-level phenomena is of scientific value only to the extent that the investigator has (a) made explicit exactly *what* high-level phenomenon he/she is trying to explain or investigate, and (b) explicitly enumerated a list of low-level phenomena (components and interactions) that he/she believes are necessary and sufficient to explain the high-level phenomenon.

If the assumptions and predictions have not been made explicit, then the output of the model will be able to tell us little of scientific value, no matter how surprising, interesting, or 'life-like' it may be. Although this point is fairly fundamental to scientific methodology in general and may seem so obvious that it is unnecessary to point it out, a quick skim through any A-life conference proceedings should be enough to demonstrate that these basic considerations are (very) often overlooked.

The number of assumptions (the low-level phenomena) that go into the model does not have to be large (e.g. they may be, say, (1) inert entities capable of (2) reproduction and (3) heritable

variation), and the high-level phenomenon under investigation does not have to be small (e.g. it could be, say, the evolution of life [but see Section 3]). However, the more explicit assumptions there are, and the more restricted the phenomenon to be explained, the more likely the model is to produce the desired results.

2.2 Minimal Models

Having devised an explicit list of low-level phenomena as a tentative reductive explanation for a specific high-level phenomenon, a model should be constructed that encapsulates these low-level phenomena *and nothing else*. In other words, it should be a minimal model. The model can then be run to see if it produces the expected results.

In practice, one generally has a choice of representations and algorithms that could be used to capture the low-level phenomena, and it may prove hard to be sure that no extra assumptions have crept into the model in the course of implementing it as a computer program (or as any other physical realization). However, as the list of low-level phenomena is explicit, the final implementation is open to testing, criticism and possible revision by others. David Marr essentially made the same point in his discussion of the three levels at which information-processing systems should be understood; he suggested that fields such as Artificial Intelligence were for too long hampered by a failure to recognize the theoretical distinction between *what* a system does (the ‘computational theory’ level), and *how* it does it (the ‘representation and algorithm’ and ‘hardware implementation’ levels) [3] (pp.19–29).

With the above in mind, once the model has been implemented, then if the expected results *are* observed, the model has demonstrated that the given assumptions are *sufficient* to explain the high-level phenomenon. To test whether all of the assumptions are *necessary*, further tests may be carried out in which assumptions are removed or relaxed one by one.

On the other hand, if the expected results are *not* observed, then the model has demonstrated that the assumptions are not sufficient. The model can then be revised by changing existing assumptions, or adding new ones.

Both cases can tell us something about the subject we are investigating, as we are always clear exactly *what* it is that we are trying to explain, and *how* we are trying to explain it. (It is much harder to conclude that assumptions are necessary to explain a given behaviour than it is to prove they are sufficient—indeed, we can *never* know that the behaviour may not also be achievable by completely different means. However, this problem is not specific to A-life, but is true of all science. All we can do is put forward our explanation as a possible model of the real world, and choose to accept the model that performs better (by some metric) than its competitors as our current ‘best guess’ on the matter [5].)

3 Specific Considerations for Models of the Evolution Of Life

3.1 The Low-Level Phenomena That Must Be Made Explicit

Darwinian (or, indeed, Lamarckian) evolution is a process of *change*. It tells us something about the *trajectory* of reproducing entities through their space of possible forms, and explains how reproducing entities become adapted to their environment. However, it *assumes* the existence of reproducing entities to begin with, and does not specify what sort of entities they should be, other than that they must be able to reproduce. Similarly, it does not specify that any particular *sort* of environment is necessary—evolution is a very general phenomenon.

A model in which a population of integers reproduce with occasional mutation, and differential survival based, perhaps, upon how large the integer is, will exhibit evolution, but it will never produce anything more than just integers. To take a more familiar example, genetic algorithms (see, e.g., [2]) satisfy the basic requirements for the evolution of the individual 'chromosomes', but all that is generally evolving is the encoded solution to some predetermined problem. Thus it is clear that if we are interested in modelling the evolution of *life*, we must (a) have a clear idea of what sorts of functions or roles a reproducing entity must fulfil if we are to consider it alive, i.e. a definition of life (this does not, of course, have to be universally agreed upon, but it *does* have to be explicitly stated), and (b) include in our model explicit components and interactions not only to allow for an (open-ended) evolutionary process to emerge, but also to allow for the existence of entities that fulfil any other functions or roles that we have specified as necessary for life.

In other words, evolution is not sufficient to explain life; we also require a theory of living organisation, and of the sorts of worlds which are able to support the emergence and evolution of such organisations. We must incorporate all of these into any A-life model designed to investigate the emergence and evolution of life.

3.2 A Definition of Life

If we are to build models to investigate or explain the emergence and evolution of life, we therefore need an explicit definition of life. That is, we need to be clear about exactly *what* we are trying to explain (as pointed out in Section 2). A number of definitions may be found in the literature, e.g. Maturana and Varela's notion of autopoiesis (see, e.g., [6]). It is emphasised that any definition adopted does not have to be universally agreed upon (although it would obviously be desirable if it were widely accepted), but it does have to be explicitly stated if we are hoping that the model will be able to tell us anything of scientific value about the evolution of life (rather than just evolution in general). Many existing A-life models that claim to have been designed to investigate the evolution of life are accompanied with no explicit statement of exactly what they are trying to demonstrate (and often also have no explicit list of the assumptions and theory involved in the construction of the model), so it is impossible to judge whether they have succeeded or failed and they can therefore tell us little of any interest.

3.3 Ecological Considerations

An aspect of biological life that seems to be particularly overlooked in many A-life models is that biological organisms are dissipative organisations that participate in exchanges of energy and matter with their biotic and abiotic environment. Perhaps more accurately, most A-life models tend to focus upon *either* evolutionary *or* ecological aspects of life, but few consider both equally. Any acceptable definition of living organisation is likely to concentrate on an organism's capability of self-maintenance in the face of environmental perturbations (caused by biotic or abiotic factors). It therefore seems probable that any A-life model of the sort we are considering will have to make explicit assumptions about the sorts of ecological interactions that are necessary and sufficient, as well as what sorts of organisations should be classified as living, and by what mechanisms they may evolve. A model that contains all of these things, and is capable of supporting a large population of organisms, may turn out to be prohibitively large for most computers at present (but maybe not). However, these are the design criteria we should move towards if models of this sort are to make significant contributions to the more general study of living systems.

4 Summary

It has been suggested in this paper that too much of the current research being done in A-life still suffers from a poor methodological approach. Specific recommendations are given to improve the situation; these basically boil down to having an explicit high-level natural phenomenon to be explained, and proposing an explicit list of low-level phenomena as a tentative reductive explanation. In Section 3 specific weaknesses are identified in the particular area of current A-life research into the evolution of life. It is suggested that such studies require a definition of living organisation, together with consideration for the sorts of environment which can support such organisation, as well as a mechanism for open-ended evolution.

Acknowledgements

Thanks to John Hallam, John Demiris and the other members of the Mobile Robots Group in the Department of Artificial Intelligence at Edinburgh University for their interesting discussions and suggestions on this topic.

References

- [1] Ezequiel A. Di Paolo. Some false starts in the construction of a research methodology for artificial life. In J. Noble and S. Parsowith, editors, *The 9th White House Papers: Graduate Research in the Cognitive and Computing Sciences at Sussex*. School of Cognitive and Computing Sciences, University of Sussex, 1996.
- [2] David E. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley, Reading, MA, 1989.
- [3] David Marr. *Vision*. W.H. Freeman, New York, 1982.
- [4] Geoffrey F. Miller. Artificial life as theoretical biology: How to do real science with computer simulation. Cognitive Science Research Paper 378, School of Cognitive and Computing Sciences, University of Sussex, 1995.
- [5] Jason Noble. The scientific status of artificial life. Poster presented at the Fourth European Conference on Artificial Life (ECAL97), Brighton, U.K., 1997.
- [6] Francisco J. Varela. *Principles of Biological Autonomy*. North Holland, Amsterdam, 1979.

Real Evolution in Artificial Chemistries

Peter Dittrich

University of Dortmund

Dept. of Computer Science, Systems Analysis

D-44221 Dortmund, Germany

dittrich@LS11.cs.uni-dortmund.de

<http://ls11-www.cs.uni-dortmund.de>

1 Overview

The intention of this contribution is to stimulate discussion concerning the following question: Are artificial chemistries a useful tool to study pre-biotic or chemical evolution? Or more general: Are artificial chemistries a useful tool to study the evolution of organizations? The latter question implies that the insights gained from the investigation of chemical-like systems can be transferred to other systems which are also forming organizations through local interactions of many components.

The first part gives an introduction to artificial chemistries. In the second part an example for “real evolution” in an artificial chemistry is shown. The term “real evolution” refers to the phenomena of self-evolution, where every variation and implicit selection is performed by the individuals (molecules) themselves and not by explicit external selection, mutation or recombination operators.

2 Introduction to Artificial Chemistries

An **artificial chemistry** is an artificial system, which is similar to a chemical system. Usually, an artificial chemistry consists of:

1. **a set of objects S** : These objects may be abstract symbols [16], character sequences [1, 12, 14], lambda-expressions [8], binary strings [3, 6, 15], numbers [4], or proofs [10].
2. **a set of rules R , describing the interaction among objects**: The rules can be defined explicitly [16, 7] or implicitly by using string matching/string concatenation [2, 13, 14], lambda-calculus [8, 9], Turing machines [15], finite state machines or machine code language [6], proof theory [9], matrix multiplication [3], or simple arithmetic operations [4].
3. **an algorithm A driving the system**: The algorithm describes how the rules are applied to a collection of objects (soup/population). The algorithm may simulate a well-stirred reaction vessel with no topology [1, 6, 8], an Euclidean discrete CA-like (fixed) topology [13, 16], a continuous 3-D space [17], or a self-organizing topology [5].

Both, the set of object and the interaction rules, can be defined explicitly or implicitly (e.g. by an algorithm or mathematical expression). An example for an implicit definition is the **number-division chemistry** [4]: In the number-division chemistry the set of objects are natural numbers s $S = \{2, 3, 4, \dots\}$. Two objects can interact, if one object can be divided by the other. The result of the interaction is the divisor and the division of the two objects. Thus, $R = \{s_1 + s_2 \implies s_3 | s_1 \bmod s_2 = 0 \wedge s_3 = s_1 / s_2\}$.

A typical algorithm simulates a well-stirred tank reactor which contains a population P of molecules out of S . Here, the population is implemented as an array of fixed size M :

```
while not terminate() do
  s1 := P[randomInteger(1, M)]
  s2 := P[randomInteger(1, M)]
  if a rule (s1 + s2 => s3) exists in R
    P[randomInteger(1, M)] := s3
  fi
od
```

This simple algorithm instantiates a mass-action kinetics equivalent to second order catalytic reactions of the form $s_1 + s_2 + X \longrightarrow s_1 + s_2 + s_3$, where the concentration of the substrate X is kept constant. The population P can be initialized by randomly selecting elements out of S .

An important aspect of this framework is, that it allows to setup *constructive* dynamical systems. i.e., the collision of molecules can generate new molecules [8, 9]. A constructive dynamical system can be treated as a dynamical system where dimension and components are changing through interactions of the components [1].

A second important aspect is, that the dynamics is *not* abstracted from the structure, as it is the case for example in the works by Eigen and Schuster on hypercycles [7]. The definition of the reaction mechanism relies on the structure of the interacting substances. The mapping from structure to function plays a key role in the process of self-organization. It allows a structure (or an organization) to operate on itself.

Finally, the inherent parallelism should be noted which allows very efficient implementation on massively parallel hardware [14].

3 Example for Self-Evolution in an Artificial Chemistry

In Fig. 1 an example for evolution without any explicit mutation, recombination, or fitness operator is shown. Every variation and implicit selection are only performed by the interacting species.

In the example, molecules are represented by binary strings with constant length of 32 bits, $S = \{0, 1\}^{32}$. An interaction among two strings $s_1, s_2 \in S$ is performed in two steps: (1) s_1 is mapped to a finite state automaton A_{s_1} by interpreting s_1 as 4-bit machine code¹. (2) The automaton A_{s_1} is applied to s_2 to generate the output s_3 . In addition elastic collisions are introduced by not allowing exact replications [8]. Thus, $R = \{s_1 + s_2 \implies s_3 \mid s_3 = A_{s_1}(s_2) \wedge s_1 \neq s_3 \wedge s_2 \neq s_3\}$. The population is initialized with $M = 10^6$ strings out of S and the algorithm mentioned above applied.

During the run displayed in Fig. 1 a lot of completely new strings are produced. Some strings are generated with rapidly increasing concentrations. A few generations later they are replaced by "better" ones. It is also interesting to note that the structure of the strings is evolving (a string is similar to its "predecessor").

¹The automaton is described in more detail in [6] and available as C++ source from <ftp://lumpi.informatik.uni-dortmund.de/pub/biocomp/src/>

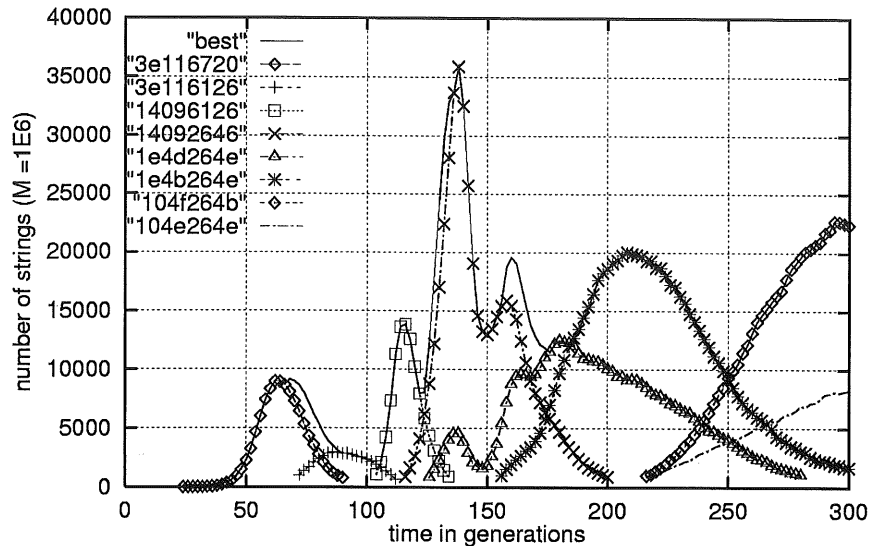


Figure 1: *Example for self-evolution in an artificial chemistry. Concentration of some representative strings are shown. Parameters: population size $M = 10^6$, automata reaction with table 2, no exact replication, population seeded with M random strings out of $\{0, 1\}^{32}$.*

4 Questions for Discussion

(1) Is an artificial chemistry suitable to study molecular/prebiotic evolution ? (1a) What is missing ? (2) What level of abstraction is reasonable ? (2a) Do we have to incorporate detailed physical/chemical knowledge ? (3) How many meta-levels of evolution are there ? (4) How can "evolution" be measured ? (5) Does "information processing" emerges in the context of evolution ? And how ? (6) How can an artificial chemistry be investigated and analyzed ? (7) Is an artificial chemistry able to create information ?

It is interesting to note that the same models used to describe chemical systems can also be found in other domains, such as population dynamics, immune system, social dynamics, economy, memetics etc. [11]. Therefore we may suspect that artificial chemistries can also serve as a tool for understanding the formation and decay of organizations in other domains. Here, a key question is: When an artificial chemistry should be used as a social dynamics model, what can be kept and what should be added ?

Acknowledgments

The author is supported by the DFG (Deutsche Forschungsgemeinschaft), grant Ba 1042/2-1 and Ba 1042/2-2.

References

- [1] R. J. Bagley and J. D. Farmer. Spontaneous emergence of a metabolism. In C. G. Langton et al. (editors), *Proceedings of the Workshop on Artificial Life (ALIFE '90)*, Redwood City, 1992, Addison-Wesley.

- [2] R. J. Bagley and J. D. Farmer. Spontaneous emergence of a metabolism. In C. G. Langton et al. (editors), *Proceedings of the Workshop on Artificial Life (ALIFE '90)*, pages 93–140, Redwood City, 1992. Addison-Wesley.
- [3] W. Banzhaf. Self-replicating sequences of binary numbers – foundations i and ii: General and strings of length $n = 4$. *Biological Cybernetics*, 69:269–281, 1993.
- [4] W. Banzhaf, P. Dittrich, and H. Rauhe. Emergent computation by catalytic reactions. *Nanotechnology*, 7(1):307–314, 1996.
- [5] P. Dittrich and W. Banzhaf. A topological structure based on hashing - emergence of a "spatial" organization. *Fourth European Conference on Artificial Life (ECAL97)*, Brighton, <http://www.cogs.susx.ac.uk/ecal97/>, 1997.
- [6] P. Dittrich and W. Banzhaf. Self-evolution in a constructive binary string system. *submitted to Artificial Life, revision in preparation*, 1998.
- [7] M. Eigen and P. Schuster. The hypercycle: a principle of natural self-organisation, part a. *Naturwissenschaften*, 64(11):541–565, November 1977.
- [8] W. Fontana. Algorithmic chemistry. In C. G. Langton et al. (editors), *Proceedings of the Workshop on Artificial Life (ALIFE '90)*, pages 159–210, Redwood City, 1992. Addison-Wesley.
- [9] W. Fontana and L. W. Buss. What would be conserved if 'the tape were played twice'? *Proc. Natl. Acad. Sci.*, 91:757–761, 1994.
- [10] W. Fontana and L. W. Buss. The barrier of objects: From dynamical systems to bounded organizations. In J. Casti and A. Karlqvist, editors, *Boundaries and Barriers*, pages 56–116. Addison-Wesley, 1996.
- [11] J. Hofbauer and K. Sigmund. *Dynamical Systems and the Theory of Evolution*. University Press, Cambridge UK, 1988.
- [12] S. A. Kauffman. *The Origins of Order*. Oxford University Press, 1993.
- [13] M. W. Lugowski. Computational metabolism: Towards biological geometries for computing. In C. G. Langton, editor, *Artificial Life*, pages 341–368. Addison-Wesley, Redwood City, CA, 1989.
- [14] J. McCaskill, H. Chorngiewski, D. Meikelburg, U. Tangen, and U. Gem. Configurable computer hardware to simulate long-time self-organization of biopolymers. *Ber. Bunsenges. Phys. Chem.*, pages 1114–1115, 1994.
- [15] M. Thürk. *Ein Modell zur Selbstorganisation von Automatenalgorithmen zum Studium molekularer Evolution*. PhD thesis, Universität at Jena, naturwissenschaftliche Fakultät, 1993.
- [16] F. J. Varela, H. R. Maturana, and R. Uribe. Autopoiesis: The organization of living systems. *BioSystems*, 5(4):187–196, 1974.
One of the first presentations of the concept of autopoiesis.
- [17] K.-P. Zauner and M. Conrad. Simulating the interplay of structure, kinetics, and dynamics in complex biochemical networks. In R. Hofestädt et al. (editors), *Computer Science and Biology—Proceedings of the German Conference on Bioinformatics (GCB'96)*, IMISE Report, pages 336–338, Leipzig, 1996. Univ. Leipzig.

- Hattemer, H.H.; Bergmann, F. and Ziehe, M. (1993): Einführung in die Genetik. Sauerlnders Verlag (publisher), Frankfurt / Main 1993
- Holland, J.J.(1975): Adaptation in natural and artificial systems. Ann Arbor, Michigan, University of Michigan Press.
- Ipsen, A.; Kasten, B.; Scholz, F. and Ziegenhagen, B.(accepted 1997): Studying allelic diversity and stress response of PEPC (Phosphoenolpyruvate carboxylase) in Norway Spruce (*Picea abies*). *Chemosphere* in press.
- Kang, M. S. and H. G. Gauch (1996): Genotype-by-Environment-Interaction. CRC Press, Inc 1996
- Kauffman, S.A. 1967: Sequential DNA replication and the control of differences in gene activity between sister chromatids - a possible factor in cell differentiation. *J. Theoret. Biol.* 17, 483 ff
- Kauffman, S.A. (1969): Metabolic Stability and Epigenesis in Randomly Constructed Genetic Nets. *J. Theoret. Biol.* (1969) 22: 437-467
- Kauffman, S.A. (1974): The Large Scale Structure and Dynamics of Gene Control Circuits: An Ensemble Approach. *J Theoret. Biol.* (1974) 44: 167-190
- Karhu, A.; Hurme, P.; Karjalainen, M. and Karvonen, P. (1996): Do molecular markers reflect patterns of differentiation in adaptive traits of conifers?.*Theor. Appl. Genet.* (1996) 93: 215-221.
- Langton, C.G.(1992): Artificial Life. In: Lectures in Complex Systems.eds: Nadel, L and Stein, D.(1991) Reading, Mass., Addison-Wesley
- Mitton, J. B. (1995): Genetics and the Physiological Ecology of Conifers. In: Smith, W. K.; Hinkley, T. M.: *Ecophysiology of Coniferous Forests*. Academic Press, pp 1-35
- Nolfi, S. and Parisi, D. (1997): 'Genotypes' for neural networks. Seth, A.K. (1997): Interaction, Uncertainty, and the Evolution of Complexity. In: Fourth European Conference of Artificial Life (ECAL97), Husbands, P. and Harvey, I. (eds.) Repsilber,D.; Kasten, B.;Ziegenhagen, B.;Gregorius, H.-R.; Scholz, F.: Investigations on the regulation of isozyme- gene-analysis using PEPCase of Norway spruce as a model. In: Abstracts of the Annual Meeting of the Genetics Society 1997 in Gieen, Germany; Sept. 23/24, 1997. Ed. R. Renkawitz. *Genes, Chromosomes, Genomes; Heidelberg* 5 (1997), p. 78
- Rothe, G. M. amd Bergmann, F. (1995): Increased Efficiency of Norway Spruce Heterozygous Phosphoenolpyruvate Carboxylase Phenotype in Response to Heavy Air Pollution. *Angew. Bot.* 69, 27 - 30 (1995)
- Wagner, G. P. and Altenberg, L. (1996): Complex adaptations and the evolution of adaptability. *Evolution*, 50 (3), 1996, pp 967-976

ALGEBRA, EVOLUTION AND COMPLEXITY OF LIVING SYSTEMS

CHRISTOPHER L. NEHANIV
UNIVERSITY OF HERTFORDSHIRE
HATFIELD HERTS AL10 9AB U.K.

JOHN L. RHODES
DEPARTMENT OF MATHEMATICS
UNIVERSITY OF CALIFORNIA AT BERKELEY
BERKELEY CA 94720 U.S.A.

Abstract

We give a natural axiomatization for the notion of hierarchical complexity measures for biological systems modelled by finite-state automata. The algebraic theory of automata is applied to show the existence of a unique maximal complexity measure satisfying these axioms, and relates hierarchical complexity to global semigroup theory. We then study the rate at which hierarchical complexity can evolve in biological systems assuming evolution is "as slow as possible" from the perspective of computational power of organisms.

Explicit bounds on the evolution of complexity are derived showing that, while the evolutionary changes in hierarchical complexity are bounded, in some circumstances complexity may more than double in certain 'genius jumps' of evolution. In fact, examples show that our bounds are sharp. We sketch the structure where such complexity jumps are known to occur and note some similarities to previously identified mechanisms in biological evolutionary transitions.

Furthermore, constructions show that in principle evolution of complexity may proceed at a surprisingly fast rate: doubling every two generations.

Genetic Networks as a Model for the Regulatory Domain applied in Ecological Genetics

D. REPSILBER AND F. SCHOLZ

Federal Research Centre for Forestry and Forest Products,
Institute for Forest Genetics, Department of Ecological Genetics
Sieker Landstr.2, D - 22927 Grosshansdorf, Germany

In population genetics stability of ecosystems is defined via the persistence of certain key species in time (GREGORIUS 1996). The persistence of a given species is determined by the adaptability of the populations within the species concerning changing environment in time and space. Within populations as the units of adaptation the genetic system determines the adaptive potential and is object to evolution (DARLINGTON 1939, GREGORIUS 1995). Population genetics analyses the genetic structure of populations by assessing allele and genotype frequencies at certain gene loci. Here, the adaptive potential is estimated in terms of allelic diversity (HATTEMER et al. 1993). On the other hand the analysis of phenotypes of quantitative adaptive traits is governed by genetic and environmental components. In quantitative genetics it is shown that adaptation includes both levels of variation, the structural level of allele and genotype frequencies and the regulatory level by which phenotypic traits are varied due to environment influences. The regulatory level is also governed genetically, i.e. by regulatory genes. These genes can vary, causing varying regulatory effects on the expressed phenotype. So far, models of population genetics and of quantitative genetics do not include the regulatory process as caused by the regulatory background of structural genes. This, however is necessary if we want to understand adaptation at the structural and at the regulatory level. Therefore so far there is a lack of explanation while trying to estimate the adaptive potential in real populations (MITTON 1995). Furthermore traits closely related to fitness tend to have rather low heritabilities (HARTL and CLARK 1989) and molecular gene markers often are poor predictors of the population differentiation of quantitative adaptive traits (KARHU et al 1996). In contrast to the allelic adaptive potential which is easily determined in the case of special assessable loci, e.g. isozyme loci, by counting alleles, the regulative adaptive potential can be measured for a special trait, mostly in provenance trials, but there is no general model yet how to integrate the understanding of the dynamics of the genetic system and the role of the regulatory domain.

The present work tries to close this gap to get a more comprehensive model which includes a simulation of the regulative domain to enable the estimation of the adaptive potential of populations on the phenotype level. In a first approach the genetic system of a population is represented by a genetic algorithm (HOLLAND 1975) where the gene-expression-system is modelled as a matrix like a neuronal net, similar to the genetic nets used for cell cycle models (KAUFFMAN 1967). For this general model population dynamics in response to changing environments are analysed and on the other hand parameters are evaluated using the results

of a specially designed genetic clone experiment with *Picea abies* as organism and PEPC as isozyme-gene-system. Predictions of the dynamics of ecosystems in response to changing environment then can take into account the regulatory domain on the population level and predict the phenotype adaptive potential for natural populations.

Approach:

The aim of a model for the regulatory domain has to be the extension of population genetic models, e.g. ECO-GENE (DEGEN and SCHOLZ 1995), considering the adaptive potential taking into account the regulative domain. On the basis of the genetic system such a model has to explain phenotypic adaptability concerning the information flow within the gene-expression system. This problem in theoretical biology is known as mapping-problem (WAGNER and ALTENBERG 1996), in population and quantitative genetics as genotype-environment interaction (FALCONER 1996, KANG 1996) and in physiology as regulation of gene activity (MITTON 1995). To model the gene-expression system in a first order approach the representation as neuronal net was chosen: Similar to the so-called genetic net (KAUFFMAN 1967, 1969 and 1974) the phenotype is not necessary a bijectivism of the genotype, but a linear image under a matrix given by the connectivity table and weight factors of the regulating system (NOLFI and PARISI 1997). Of course non-linear interactions will form an essential part of the physiological regulating system (LANGTON 1991), but on the other hand neuronal nets are well known in mathematical aspects so that first steps of analysis are enhanced. Taking neuronal nets as an image for the regulatory system can possibly explain system properties related to the phenotype of adaptive traits and contribute to an integrated understanding of the adaptive potential determined by the genetic system, but formed out to act as phenotype by the system of trait formation.

The experimental foundation to get the heuristics and to evaluate the models parameters was a specially designed physiological genetic clone experiment. The experimental design was chosen to evaluate the genetic determined part of a trait's variation caused by the regulatory domain. The type of organism looked on requires an immense extent of adaptability, because being sessile and long-lived. A gene locus has been chosen, which plays a role in adaptation and whose modules of the expression-system are object to actual research projects: *Picea abies* as organism and PEPC as isozyme-gene-system (ROTHER and BERGMANN 1995, IPSEN et al. 1996, REPSILBER et al. 1997) meet the requirements for studying the system of trait formation and its role in determining the adaptive potential. A hierarchy of clones with different relationships and from different provenances were exposed to different temperature regimes to evaluate the regulation system in dependence of different genetic backgrounds, whereas the PEPC-genotype is known for each clone. As physiological trait the temperature dependence of the specific enzyme-activities and amounts of enzyme were measured.

The results of the experimental part will give an estimation of the natural "connectivity" for this enzyme-system. Differences between homozygote and heterozygote individuals could indicate the system properties leading to the observation of the so-called heterozygous advantage, which was observed in spruce populations adapting to heavy air pollution (BERGMANN and SCHOLZ 1989). The impact of the regulatory domain on the realisation of the genetic system can be measured for this example.

Processing status:

The experimental basis has been completed in 1997. First results concerning different adaptive potentials of PEPC homozygote and heterozygote individuals (publication in preparation) show that heterozygote organisms make full use of their higher variability in enzyme composition to regulate their enzyme activity. Variance analysis to exploit the data set is being carried out to estimate the regulation system's parameters as there are: "Connectivity" of the regulation system's network, dimension of the regulatory matrix, degree of hierarchical organisation of the regulation net. So far SAS-analyses of experimental data is going on and showing interesting results concerning variance components in the trait expressing system consisting of enzyme-genotype genetical background and regulation and the environmental conditions.

The model approach employs a genetic algorithm (HOLLAND 1975) to simulate an evolution of neural networks similar to the approach of Nolfi and Parisi (NOLFI and PARISI 1997). The neural network is identified with the regulating network. Fitness is calculated as the similarity of the networks output in comparison to the environments input. A rough draft of the simulation program has been tested for basic new population dynamics compared to the variant using direct genotype phenotype mapping: Direct mapping populations show faster adaptation, but are less well adapted to changing environments. This behaviour is nothing new for the comparison of complex and simple systems (SETH, A.K. 1997). Further analyses in comparison with models that simulate the dynamics of the genetic system will feature differences in prognosis tendencies due to the integration of the genotype-phenotype-mapping module. It is planned to investigate if on the population level the adaptive potential is systematically underestimated if only the allelic diversity is taken into account. At the moment only phylogenetic adaptation is modelled - the next step is taken in modelling the ontogenetic part of regulative adaptation.

Literature:

- Bergmann, F. and Scholz, F. (1989): Selection effects of air pollution in Norway spruce populations. In: Scholz, F., Gregorius, H.-R. and Rudin, D. (eds.), Genetic effects of Air Pollutants in Forest Tree Populations, 143 - 160. Springer-Verlag, Berlin, Heidelberg.
- Darlington, C. D. (1939): The evolution of genetic systems. Cambridge: University Press, pp 1-254
- Degen, B.; Gregorius, H.-R. and Scholz, F.(1996): ECO-GENE, a model for simulation studies on the spatial and temporal dynamics of genetic structures of tree populations. *Silvae Genetica* 45: 323-329.
- Falconer, D. S. (1996): Introduction to Quantitative Genetics. Longman Group Ltd (publishers) 1996, pp131 ff
- Gregorius, H.-R. (1995): Measurement of genetic diversity with special reference to the adaptive potential of populations. In: Boyle, C. E. B.: Measuring and Monitoring Biodiversity in Tropical and temperate Forests.
- Gregorius, H.-R. (1996): The Contribution of the Genetics of Populations to Ecosystem Stability, *Silvae Genetica* 45 (1996), pp 5-6
- Hartl, D. L. and Clark, A. G.(1989): Principles of population genetics. Sinauer Associates, Inc (publishers) pp 480 ff

Molecular Transition Systems: A Computational embodiment for the Graded Autocatalysis Replication Domain (GARD) model

Preliminary Draft, Extended Abstract

Ehud Shapiro(1), Doron Lancet(2), Daniel Segr?(2),
(1) Department of Applied Math and Computer Science
(2) Department of Molecular Genetics and the Genome Center
Weizmann Institute of Science
Rehovot 76100
Israel

[For correspondence: Ehud Shapiro, 314 W. 78 St., New York, NY 10024
udi@cs.weizmann.ac.il]

Abstract

We describe a computational model for molecular transition systems and discuss its significance to mutual catalysis scenarios geared towards the study of the origins of life.

1. Introduction

Research indicates that all existing life on earth originated some 3.7 billion years ago from a single cell type, or progenote. While there is ample research explaining how complicated life forms can arise through evolution once the basic mechanisms of the cell are in place, the question of how these mechanisms came into being, or even how could have they come into being, is a profound mystery.

From a computer science perspective, the progenote is an extremely complicated and sophisticated entity, with multiple components, each of which can function only in conjunction with some or all of the others. Several proposals were made regarding a plausible explanation for the emergence of the progenote, but the gap between inanimate matter and the complexity of the structure and the processes of the progenote remains to be explained.

The seminal work regarding abstract models of living systems is von Neumann's work on self-reproducing automata, dated back to 1949. Several aspects of his work are fascinating, especially with a 50-years or so hindsight. First, while his work predates the discovery of the DNA, the basic element of his self-replicating automata, namely a description of the automaton that is both interpreted during the "life" of the automaton and replicated into the automaton "progeny" is identical in concept and in function to the DNA. Second, while his work on self-reproducing automata was done at the same period in which he developed the stored-program computer, follow-on work to his that

utilizes developments in computer science that occurred since the first stored-program computer was built is scarce (notable exceptions include the Chemical Abstract Machine by Gerrard Berri and the work of Fontana and colleagues).

In parallel to the standard organic chemistry approaches to the origin of life, scientists from various disciplines have been trying to investigate biogenesis by theoretical modeling. Prominent examples are Dyson's statistical physics model for homeostatic catalytic networks, Morowitz's studies on the thermodynamics of protocellular aggregates, Kauffman's graph theory analysis of catalytic networks, Bagley's kinetic approach to mutual catalysis and Fontana's lambda-calculus-based algorithmic chemistry. Such models describe complex chemical interactions among organic molecules and their assemblies, prior to the emergence of DNA, RNA and proteins. They purport to demonstrate how a transition could occur from early random mixtures of organic molecules to the first protocell.

In the framework of this approach, two of us (DL and DS) have recently developed the Graded Autocatalysis Replication Domain (GARD) model (2,3), which provides a thermodynamic and kinetic analysis of mutually catalytic assemblies combined with statistical tools. It assumes a finite micellar enclosure, containing the catalytic set members, and utilizing energy-rich chemical precursors from the external environment. Through the solution of differential equations and by Monte Carlo simulations, GARD predicts the spontaneous emergence of assemblies with idiosyncratic molecular compositions, capable of carrying information, as well as of undergoing rudimentary self-replication and chemical evolution.

The present degree of formalization of GARD allows only a limited category of reaction topologies (isomerization, dimerization and non-covalent recruitment). A broader scenario, en-route to a protocell, should involve a much larger variety of more complex chemical species and reactions, including the potential emergence of templating and primitive genetic codes. We believe that this may be achieved through the implementation of Theoretical Computer Science concepts, including molecular transition systems.

2. The GARD model

In the Graded Autocatalysis Replication Domain (GARD), mutual catalysis within a set of N types of molecules (A_i) derived from a common precursor (A_0), can sustain self-replication of the entire ensemble (2,3). The components of GARD may be any organic molecules, endowed with sufficient complexity to allow for structural diversity and mutual complementarity. The mutual catalytic rate enhancement exerted on the species A_i by the species A_j is denoted by a matrix element β_{ij} . For GARD simulations we use a formalism that allows one to assign likelihood values for any degree of catalysis between two randomly chosen species A_i and A_j . This is described in the form of a probability distribution $\phi(\beta_{ij})$, analogous to our previously developed Receptor Affinity Distribution (RAD) model (1). We further assume that the system is subjected to a constant dilution effect, (e.g. due to expansion of its vesicle enclosed volume). The time-dependent concentrations of the species A_i then obey the differential equations ($i=1,N$)

$$\frac{dA_i}{dt} = k_i A_0 - k_{-i} A_i + \sum_{j=1}^N k_i \beta_{ij} A_0 A_j - \sum_{j=1}^N k_{-i} \beta_{ij} A_i A_j - \lambda A_i$$

where k_i and k_{-i} are the uncatylsed rate constants of the reaction of formation and degradation of A_i and l describes the system's exponential expansion rate, according to

$$V(t)=V(0) \cdot \text{Exp}(\lambda t).$$

The main question addressed by the GARD model is whether a chemical system, connected through random catalytic interactions, and governed by a statistical catalysis formalism (1-3), can propagate its own chemical composition with no absolute requirement for autocatalysis for any of its individual components. For this, we define a graded quantitative measure for GARD's self replication through a critical rate of dilution (λ_c), that is shown to increase with the extent and connectivity of mutual catalysis. We then envisage an evolutionary process, where the content of GARD is subjected to random compositional fluctuations, which affect the prevailing network of catalytic interactions. We analyze this process by stochastic computer simulation, where compositional "mutants" with an augmented capacity of self replication may spontaneously appear, which may take over the prevailing GARD.

The GARD model demonstrates quantitatively how self replication may be a property of a molecular ensemble, without any specific constraints on the structure of the components. No individual molecule needs to be endowed with the specialized chemical properties currently associated with replicating macromolecules such as DNA and RNA. In the above, the simplest form of GARD is described. Higher level simulations, with higher degree of polymers may be considered, and their analysis is expected to be made possible by the molecular transition system described below. Our analysis also allows to compute the probability for a primordial spontaneous emergence of a GARD-like entity based on first chemical principles. GARD may thus be considered as a feasible paradigm for understanding the early emergence of chemical self-replication and chemical selection.

3. Requirements from a computational model.

We believe a more "advanced" answer for the problem von Neumann was interested in, namely a computational model suitable for the study of the origins of life, should satisfy the following requirements:

1. The model should bear higher structural resemblance to biochemical environments. As in biochemistry, the building blocks should be (abstractions of) monomers and polymers, and the basic transition rules should be the rules that govern the interactions of (abstract) monomers and polymers.
2. The model should be concurrent, to model both parallelism within a cell as well as interaction among multiple organisms (We believe cellular automata are inadequate modeling both internal parallelism and interaction among organisms).
3. The model should not distinguish between "program" and "data" and, ideally, the notion of "program" should be emergent, not built-in. Specifically, the "meaning" of DNA or its abstract equivalent should not be "built in", but derived through interpretative mechanisms as in the living cell.

The purpose of our work is to devise and investigate such models.

4. Molecular Transition Systems

We describe here Molecular Transition Systems, a preliminary step in the direction outlined above.

Transition systems are one of the standard tools for defining and studying abstract computational models, and are especially geared towards the study of concurrent computing. A transition system consists of the following components:

1. A set of states, and within this set a subset of allowed initial states.
2. A set of transitions $S \rightarrow S'$, where S and S' are states, defining how one state can change to another.

A *computation* of a transition system is a sequence of states S_1, S_2, \dots such that S_1 is an initial state and $S_i \rightarrow S_{i+1}$ is a transition, for every pair S_i, S_{i+1} in the sequence. Infinite computations typically have to satisfy additional constraints which are beyond the scope of this discussion.

Molecular Transition Systems are transition systems with the following characteristics:

- ? States are multisets of polymers, which in turn are orientation-free strings of monomers (we do not distinguish between a string and its inverse),
- ? Reactions consist of ligation and cleavage of polymers, possibly with the aid of a third polymer serving as a catalyst.

Definition: A *Molecular Transition System* is a transition system with the following:

- ? A set of monomers $M = \{M_1, M_2, \dots, M_n\}$.
- ? *Polymers* are orientation-free strings over M .
- ? *States* are multisets of polymers.
- ? *Reactions* among polymers have the form:
 - ? Ligation: $A, B, C \rightarrow AB, C$
 - ? Cleavage: $AB, C \rightarrow A, B, C$where A, B , and C are polymers, with C possibly being absent. AB is the string resulting from the concatenation of the strings A and B .
- ? *Transitions* are pairs of states $S \rightarrow S'$ where S' is obtained from S by replacing the polymers on the left-hand side of a reaction by the polymers on its right-hand side.

We ignore for now what are the initial states and what are the constraints on the application of a transition.

5. Conclusion

The introduction of a general chemical reaction programming language in the form of a molecular transition system, will allow to reformulate the GARD model in a more formal and rigorous fashion. Most importantly, the unlimited number of chemical species and

reactions that could be present in such a generalized system has the potential to analyze the transition from a compositional information based system (as GARD is), to a protoliving unit in which a primordial coding mechanism resembling the modern genetic code could arise.

6. References

1. Lancet, D., Sadovsky, E. and Seidemann, E. Probability Model for Molecular Recognition in Biological Receptor Repertoires: Significance to the Olfactory System. Proc. Natl. Acad. Sci., USA 90: 3715-3719 (1993).
2. Lancet, D., Segr?, D., Kedem, O. and Pilpel, Y. Graded Autocatalysis Replication Domain (GARD): Kinetic Analysis of Self-Replication in Mutually Catalytic Sets. Origins of Life and Evolution of the Biosphere 28:000-000 (1998).
3. Segr?, D., Pilpel, Y. and Lancet, D. Mutual Catalysis in Sets of Prebiotic Organic Molecules: Evolution Through Computer Simulated Chemical Kinetics. Physica A 249(1-4): 558-564 (1998).