

Design and Characterisation of a Novel Artificial Life System Incorporating Hierarchical Selection

Ciarán Kelly BSc. MSc.

A dissertation submitted in fulfilment of the requirements for the award of

Doctor of Philosophy (Ph.D.)

to the



Dublin City University

Faculty of Engineering and Computing

Supervised by Prof. Barry McMullin & Dr. Darragh O'Brien

Examined by Dr. Peter J. Bentley & Dr. Noel Murphy

September 2010

I hereby certify that this material, which I now submit for assessment on the programme of study leading to the award of *Doctor of Philosophy* is entirely my own work, that I have exercised reasonable care to ensure that the work is original, and does not to the best of my knowledge breach any law of copyright, and has not been taken from the work of others save and to the extent that such work has been cited and acknowledged within the text of my work.

Signed: _____ (Ciarán Kelly) ID No.: 99290146

Date: _____

To Nora, Maureen, Ted and Billy

Acknowledgements

I extend my sincerest gratitude to my parents, Des and Teresa, and especially to Victoria. Your love and patience were the fuel that kept the fire burning. To Prof. Barry McMullin and Dr. Darragh O'Brien: you gave me the best guidance and supervision I could have hoped for—I couldn't have gotten here without you. Your advice and encouragement were always there whenever I needed a push in the right direction. To Joan, for providing a "home from home", and for being such a good sport during our many heated "debates". And to Conor, Niamh, Orla and Molly, for making Tonlegee Road such a fun place to live. To Esmond, James, George, Elaine, Johanna, Pietro and all the interns that passed through the DCU aLife Lab: thanks for everything.

And finally, to all my friends and family: it's been a long journey and your support helped me see the light at the end of the tunnel to get where I am today.

Thank You!

Contents

Abstract	i
1 Introduction	1
1.1 Context	1
1.2 Research Problem and Thesis Contribution	5
1.3 Thesis Overview	6
1.4 Associated Publications	9
1.4.1 Conference Proceedings	10
1.4.2 Posters	10
I Literature Review	11
2 Life	12
2.1 What is Life?	13
2.1.1 The “definition” problem	13
2.1.2 Definitions can be restrictive	14
2.1.3 A satisfactory definition	15
2.2 Darwin’s Insight	15
2.2.1 Evolution by Natural Selection	16
2.2.2 Heritability	17
2.3 Information Chemistry	18
2.3.1 DNA Replication	19
2.3.2 Transcription	20
2.3.3 Translation	21
2.3.4 DNA—A Ubiquitous Language?	21
2.4 The Origin of Life	22
2.4.1 Overview of timescale of Origin of Life	22
2.4.2 Major Evolutionary Transitions	23
2.4.3 The Transition to Cells	24
2.5 Challenges to Life	28
2.5.1 Environmental Challenges	28

2.5.2	Organisational Challenges	28
2.5.3	Proposed Solutions	29
2.6	Conclusion	31
3	Artificial Life	32
3.1	Introduction	34
3.2	Soft Artificial Life	35
3.2.1	Von Neumann’s <i>Theory of Automata</i>	36
3.2.2	Genetic Algorithms, Classifier Systems & α -universes	38
3.2.3	CoreWar, Tierra, Avida	39
3.2.4	Algorithmic Chemistry	41
3.3	Wet-Lab Artificial Life	41
3.3.1	Protocells	42
3.3.2	Bottom-Up Approach	43
3.3.3	Top-Down Approach	43
3.4	Conclusion	44
II	Experimental System Design	45
4	Building MCS	46
4.1	Methodology	46
4.2	Underlying Theory	48
4.2.1	Replicator Theory	48
4.2.2	Catalysed Replication	49
4.3	Idealisation	49
4.4	System Design	51
4.4.1	Modelling Platform	52
4.4.2	MCS-“world”	52
4.4.3	Molecular Sub-System	55
4.4.4	Protocell Sub-System	56
4.4.5	Reaction Mechanism	57
4.5	Theoretical Analysis of Molecular Chemistry	60
4.5.1	General Dynamic Equation	60
4.5.2	Terminology	61
4.5.3	“Self”-Systems	61
4.5.4	Binary Replicase Interaction Systems	62
4.6	Conclusion	74

III	Experimental Results	75
5	MCS-0	76
5.1	MCS-0 Specification	76
5.1.1	Molecular Selection	77
5.1.2	Parasitic Mutation Network	81
5.1.3	Further Characteristics of Molecular Mutation	82
5.1.4	Cellular Selection	84
5.2	Experiments: MCS-0	90
5.2.1	Survival of the Common	90
5.2.2	Takeover by Facultative Parasite	94
5.2.3	Elongation Ratcheting	97
5.2.4	Protocell Evolution in MCS-0 Systems	100
5.3	Conclusion	105
6	MCS-1	106
6.1	Introduction	106
6.2	MCS-1: The Modifications	107
6.2.1	Molecular Folding	109
6.2.2	Molecular Binding	110
6.2.3	Expanded Matching Capabilities	112
6.3	MCS-1: The Experiments	112
6.3.1	Collapse of Seed Self-Replicase Species	113
6.3.2	MCS-1 Molecular Evolution	121
6.3.3	Protocell Evolution in MCS-1 Systems	124
6.4	Summary of Experimental Results	134
6.5	Conclusion and Wider Impact of Results	134
IV	Concluding Remarks	137
7	Discussion	138
7.1	Recap of Key Achievements	139
7.2	Potential for Further Work	140
7.2.1	Further Investigation of Evolutionary Phenomena	140
7.2.2	Extension of Computational Capabilities	141
7.3	Limitations and Potential Criticisms	142
7.4	Conclusion	143

V	Appendices	144
A	Living Computation	145
A.1	Introduction	146
A.2	What is Living Computation?	146
A.2.1	Programming Living Computers	148
A.3	Examples of Living Computation	148
A.3.1	Cell Cycle Control	149
A.3.2	Cellular Chemotaxis	151
A.3.3	Quorum Sensing	151
A.4	Conclusion	152
B	MCS-2	154
B.1	Introduction	154
B.2	Overview of modifications to MCS	156
B.2.1	Relevance to Wet-Lab Protocells	157
B.2.2	Variability of Fission Condition	157
B.3	MCS-2: The Modifications	158
B.3.1	Modifications to Molecular level chemistry	158
B.3.2	Modifications to Protocell level chemistry	160
B.3.3	Modifications to the MCS Reaction Algorithm	160
B.4	MCS-2: The Experiments	170
B.4.1	Evolution of Protocell-embedded molecular computation	171
B.5	Detailed Analysis of Figure B.4	176
B.5.1	Discussion of Results	181
B.6	Conclusion	182
	Bibliography	184

Abstract

In this thesis, a minimal artificial chemistry system is presented, which is inspired by the RNA World hypothesis and is loosely based on Holland's Learning Classifier Systems. The Molecular Classifier System (MCS) takes a bottom-up, individual-based approach to building artificial bio-chemical networks. The MCS has been developed to demonstrate the effects of hierarchical selection. Hierarchical selection appears to have been critical for the evolution of complexity in life as we know it yet, to date, no computational artificial life system has investigated the viability of using hierarchical selection as a mechanism for achieving qualitatively similar results. Hierarchy in MCS is enforced by constraining artificial molecules, which are modeled as individuals, to exist within externally provided containers—protocells. This research is focused on the period of time surrounding the conjectured first Major Transition—from individual replicating molecules to populations of molecules existing within cells. Protocells can be thought of as simplified versions of contemporary biological cells. Molecular replication within these protocells causes them to grow until they undergo a process of binary fission. Darwinian selection is continuously and independently applied at both the molecular level and the protocell level. Experimental results are presented which display the phenomenon of “selectional stalemate” where the selectional pressures are applied in opposite directions such that they meet in the middle. The work culminates with the presentation of a stable artificial protocell system which is capable of demonstrating ongoing evolution at the protocell level via hierarchical selection of molecular species. Supplementary results are presented in the Appendix material as a set of experiments where selectional pressure is applied at the protocell level in a manner that indirectly favors particular artificial bio-chemical networks at the molecular level. It is shown that a molecular trait which serves no useful purpose to the molecules when they are not contained within protocells is exploited for the benefit of the collective once the molecules are constrained to live together. It is further shown that through the mechanism of hierarchical selection, the second-order effects of this molecular trait can be used by evolution to distinguish between protocells which contain desirable networks, and those that do not. A treatment of the computational potential of such a mechanism is presented with special attention given to the idea that such computation may indeed form the basis for the later evolution of the complicated Cell Signaling Pathways that are exhibited by modern cells.

Chapter 1

Introduction

1.1 Context

Life, as we know it, is thought to have emerged on Earth about 3.5 Billion years ago, though we know for certain that living organisms existed at least 2.5 Billion years ago (Lepot et al., 2008). Evolution, by natural selection, is thought to have been the driving force behind the subsequent growth of complexity of these early living organisms which gave rise to the rich diversity of Life today (Darwin, 1859; Dawkins, 2009). Unfortunately, we do not know very much at all about the earliest history of life on Earth, as those initial single-celled ancestors left behind no fossils. Origin of life research has produced many interesting candidate theories of how life *might* have begun, but it seems unlikely that a definitive answer will ever be possible. Meanwhile, the field of artificial life has emerged with a focus on the investigation of Life, as it *could* be, where research is typically carried out using digital computers and software simulations (Langton, 1989a) as opposed to wet-lab experiments using chemical apparatus and processes. The core methodology of the field of artificial life then is the creation of artificial, virtual universes in which different forms of life might arise and/or evolve.

Evolution, acting upon life as we know it, has succeeded in producing progressively more “complicated” organisms, in at least some lineages (Maynard Smith, 1969). In terms of humanity’s best engineering efforts, we have yet to incorporate a similar automated growth of complexity in the objects that we build—the manufacturing plant is still considerably more complicated than the objects it produces. According to McMullin (2000), the problem of the “evolutionary growth of complexity” is precisely the problem that John von Neumann, a pioneer of digital computing and the theory of self-reproducing automata, formulated and addressed. Von Neumann’s contribution to the solution of this problem (Von Neumann and Burks, 1966) was quite elegant: a “general constructive automaton” which might be

capable of building *any* other machine (within the “physical” constraints of the “universe” in question), assuming it was given the appropriate instructional blueprints for the target and sufficient raw materials. As a special case, self-reproduction would be theoretically possible if the instructional blueprints described the general constructive automaton *itself*. Von Neumann further recognised that the reproduction process would also require the copying of the instructional blueprints to the “offspring” and that errors in the copying of the instructions could result in *heritable* modifications. The implication of this is that it is at least *possible* for such automata to produce offspring which are incrementally but indefinitely more complicated than themselves, and thus potentially to support an open-ended evolutionary growth of complexity, even in a completely artificial, “automaton” world. Maynard Smith and Szathmáry (2000, pp. 73) distinguish between systems of “unlimited heredity” such as von Neumann’s which use a “modular” hereditary information carrier of indefinite length and systems of “limited heredity, whereby a replicator can exist in a relatively small number of states, each of which reproduces its kind”. Life, as we know it, demonstrates the potential for unlimited heredity through the mechanism of a DNA-based genetic information store.

As far as artificial life models go, Ray’s Tierra (Ray, 1991) has arguably been the most successful demonstration of von Neumann-style automaton evolution to date. The automata in this case were computer programs, which were interpreted and embodied by the Tierran virtual machine. Tierra successfully demonstrated some substantive and interesting evolutionary phenomena, namely a rich ecology of replicators and parasites, but fell short of demonstrating the kinds of sustained or open-ended evolutionary growth of complexity that can be observed in life as we know it, instead consistently reaching a “plateau of complexity” which would appear to be intrinsic to the Tierran architecture.

Maynard Smith and Szathmáry (1997) argue that the evolutionary history of life as we know it has been punctuated by “major evolutionary transitions”, and further, that such transitions have been *critical* to the evolution of complexity that characterises life. They claim that such transitions occur when a new hierarchical level of selection arises from the combined actions of lower level entities. Tierra explicitly incorporates only a single level of selection, with no apparent mechanism for entities at that level to aggregate and organise themselves into a higher level of selection. The theory of major evolutionary transitions holds:

1. that evolution, by natural selection, can continue to act separately upon entities at all levels of the hierarchy , and,
2. that the evolutionary dynamics of entities lower down the hierarchy strongly influences the evolutionary dynamics of the entities at higher levels.

Bedau et al. (2001) present a list of open problems in the field of artificial life. From the point of view of this thesis, the most relevant problems from that list are to:

1. “Achieve the transition to life in an artificial chemistry in silico.”
2. “Create a formal framework for synthesizing dynamical hierarchies at all scales.”

While Maynard-Smith and Szathmáry describe eight theoretical major transitions, it is the first major transition “From individual replicating molecules to populations of replicating molecules in compartments” which most readily aligns with these two “grand challenges” facing artificial life. For the purposes of the work carried out in this thesis, the “replicating molecules” will be implemented in an artificial chemistry (Dittrich et al., 2001) which will describe the kinds of molecules which can occur, and also define the mechanisms by which those molecules interact. “Dynamical hierarchies” can be viewed as a generalisation of the hierarchical evolutionary actors envisaged by the theory of major transitions, though the work that is presented in this thesis does not aim to create the “*formal* framework for synthesizing dynamical hierarchies at *all* scales” that Bedau et al. describe, but rather a demonstration of a dynamical hierarchy of two levels of selection in the form of artificial protocells. Protocells are a hypothesised transitional phase in the origin of life and can be thought of as simplified versions of contemporary biological cells. Molecular replication within these protocells causes them to grow and eventually split into two by a process known as binary fission. Heritable information is not centrally mediated by a genetic mechanism like DNA, however: rather protocell “species” can be differentiated by the composition of the populations of replicating molecules which drive their growth.

Origin of life research has proposed many interesting candidate theories for the earliest chemical systems which might plausibly have been the basis for life as we know it. This research has primarily focused on three chemical sub-systems (information, metabolism, containment) which are thought to be required for the evolution of a minimal protocell (Luisi, 2002). The RNA-World hypothesis (Kruger et al., 1982; Gilbert, 1986; Joyce, 1991) describes a pre-biotic “replicator world” where RNA molecules act as both information carriers, and enzymatic metabolites, which might later serve as a platform for incorporation into compartments during something like Maynard-Smith and Szathmáry’s first major transition. Manfred Eigen and collaborators have shown however that while such RNA-based replicator systems are theoretically capable of supporting collective autocatalysis through mechanisms such as the hypercycle (Eigen, 1971; Eigen and Schuster, 1977), their susceptibility to parasitism is a significant weakness in the theory. Though spatial segregation through containment or some other mechanism has been shown to be capable of

reducing the impact of parasites in hypercycles (Hogeweg and Takeuchi, 2003), other models such as Szathmáry’s “Stochastic Corrector” (Szathmáry, 1986) have shown that the problem of “information integration” can be solved without resorting to hypercycles at all. Modelling toolkits such as Dissipative Particle Dynamics (DPD) allow for the direct simulation of chemical interaction systems, but due to time and space constraints, the simulation of anything approaching the complexity of an integrated (proto)cell make this approach infeasible for our purposes. Furthermore, such a simulation would require a complete understanding of the molecular dynamics of such a cell—something that we still don’t have. The *artificial* chemistry which is presented later in this thesis is inspired by this kind of replicator world, though the precise chemical interactions that will occur are not directly modelled on any “real” chemical system that we know about. As we have seen, there are a number of reasons why this approach is desirable, not least of which is that our goal is exploring life as it *could* be rather than the specific incarnation of Life that has occurred on Earth.

Due to these simulation limitations, it is clear that any computational model of artificial life will contain some abstractions and idealisations. The potential for evolving complexity in Ray’s Tierra appears to have been limited by the Tierran architecture itself, suggesting that there are at least *some* idealisations which are *not* desirable in artificial life systems. As yet, we simply don’t know which critical abstract features are necessary for life, in the sense of open-ended evolutionary growth of complexity, in *any* medium. At the same time, we have seen that evolution acting at all levels of “dynamical hierarchies” is thought to have been critical for the evolution of complexity for life as we know it. This thesis shall explore the phenomenon of hierarchical selection, which is potentially one of the above mentioned critical features.

As mentioned previously, it is not the goal of this thesis, nor artificial life in general, to build a system which models any incarnation of Life as we know it, but rather to explore the kinds of Life that might be possible, given a particular artificial “universe” as a starting point—“life as it could be” (Langton, 1989b). Wet-lab artificial life research projects have recently been focusing their efforts to fabricate synthetic protocells in-vitro, for example the Programmable Artificial Cell Evolution (PACE) (EU-FP6-IST-FET-002035) project, which partially funded the research presented in this thesis. Such wet lab protocell experiments cannot avoid the “dynamical hierarchies” problem that has been conveniently side-stepped by artificial life models like Tierra. The computational artificial life model presented in this thesis is a simulation of such PACE-style protocells, specifically embracing the kinds of hierarchical selection that occur in dynamical hierarchies, while idealising other, arguably less

critical features of the underlying chemistry. It is hoped that this kind of modelling approach might be qualitatively validated against the *kinds* of life that we can observe in the biosphere, as opposed to being quantitatively validated against specific wet-lab experiments.

1.2 Research Problem and Thesis Contribution

The central research problem addressed in this thesis is as follows:

To design and build a minimal artificial life system incorporating agents operating at at least two interacting hierarchical levels of selection, exhibiting unlimited heredity at both levels; and to test and characterise the resulting evolutionary dynamics.

In this thesis, I combine research in the interdisciplinary field of artificial life with theories about the origin of life on Earth, to construct an “artificial chemistry” (Dittrich et al., 2001) which supports hierarchical selection in the spirit of the “dynamical hierarchies” described by Bedau et al. (2001) and Maynard Smith and Szathmáry (1997). Specifically, this platform attempts to model an idealised abstraction of a “protocell world”. Protocells can be thought of as simplified versions of contemporary biological cells. The abstract protocells which are investigated here exhibit hierarchical selection—that is, selectional pressure is applied both at the level of protocells competing with other protocells, and internally, with molecules competing with other molecules. The following list highlights the key contributions of this thesis:

- The formulation of a framework which is devised to incorporate selectional dynamics at two hierarchical levels, where the lower level is an abstraction of the (hypothetical) self-replicase ribozyme molecules of the RNA-world, and the higher level is an abstraction of the (hypothetical) protocells, arising in the (hypothetical) first major transition, and composed of “cellular aggregates” of the molecular replicators. A key distinction between the molecular chemistry assumed here and in other artificial life models is that molecular replication here is declared to be *trans*-acting (two entities react to produce a third). This is in contrast to the *cis*-acting replication which is typical of von Neumann’s self-reproducing automata and computational artificial life models such as Tierra, where entities are capable of *self*-replication (one entity directly gives rise to a copy of itself, by a form of self-inspection).

- The presentation of a generic ODE analysis of the entire set of molecular reactions¹ possible in the lower-level, *trans*-acting, artificial chemistry presented in this thesis. This analysis will be used to validate the molecular level experimental results that follow in the thesis, and also as a basis for describing some of the evolutionary dynamics that arise when hierarchical selection is later applied. The artificial life model presented here then diverges from previous, comparable, artificial life models in two important ways:
 1. the model incorporates interacting evolutionary dynamics at multiple hierarchical levels, which, according to the theory of major transitions, is a key feature of the evolutionary growth of complexity in real biology, and
 2. the lowest level replicators rely on *trans*- rather than *cis*- replication, which may or may not be critical to the evolutionary growth of complexity, but *is* a known characteristic of the RNA-world hypothesis, and therefore of the hypothesised *first* major transition in the origin of life, and which has *not* been previously investigated as a feature of *artificial* evolutionary systems..
- The demonstration of a canonical example of the effects of hierarchical selection. This is deliberately contrived to apply selectional pressure in opposite directions for each of two hierarchical levels to make the effects of hierarchical selection as clear as possible. This results, in this exemplar case, in a robust “selectional stalemate”. This stalemate can only be explained by examining the evolutionary dynamics at both hierarchical levels and how they interact.
- The extension of the system to demonstrate *ongoing* evolutionary dynamics through varying interactions between the two levels of hierarchical selection. This is achieved by modifying one aspect of the molecular chemistry, namely the binding rule which determines whether an interaction between two given molecules will result in a molecular replication event or not.

1.3 Thesis Overview

The thesis is divided into five parts which are outlined below. Part I (Chapters 2 and 3) contains a number of chapters which review the prior scientific work upon

¹Note that, as is normal usage in the literature on artificial chemistries (Dittrich et al., 2001), molecular biochemical terminology shall be used as appropriate when referring to the abstract counterparts of these terms in the artificial system. In general, these terms will *not* carry all the same connotations in such an artificial system as they would have in real biology, but rather they will have exactly and only the meanings that will be axiomatically associated with them later in Part II of the thesis.

which the thesis is built.

In Chapter 2, I present a review of some of the key elements which compose our current understanding of life on Earth and a high-level introduction to the theory of evolution by natural selection. This chapter draws on research from a highly inter-disciplinary spectrum of fields such as Chemistry, Biology and Physics. Included is a description of current scientific understanding of the origin of life, with particular attention given to presenting some of the systemic problems that life would have had to overcome, independently of the physico/chemical problems presented by the environment. The purpose here is to highlight some of the kinds of problems that life has succeeded in solving. Chapter 2 also describes a framework of “major evolutionary transitions” via hierarchical selection, which is hypothesized to be critical to the evolution of complexity in life as we know it, and is the core motivation for the artificial evolutionary system that is presented later in the thesis.

Chapter 3 introduces the field of artificial life. This chapter surveys how scientists use the digital computer to carry out simulations of living systems in an attempt to explore the key organisational features and requirements of such systems. This chapter also describes the emerging field of building artificial cells mentioned earlier, namely the creation of protocells in-vitro. Protocells are hypothesized as a transitional phase in the early evolution of life and as such, recreating protocells in the laboratory may enable the later creation and evolution of new man-made life-forms. This wet-lab work has been approached from both the top-down perspective and the bottom-up perspective. A key difficulty faced by scientists in these areas is that both the physico-chemical and organisational problems must be tackled simultaneously. There is an opportunity in the kinds of computational universes described in the early part of this chapter to simplify these complexities but to retain the large scale organisational properties of the systems in question. It is argued in this chapter that none of the computational universes investigated exhibit open-ended growth of complexity and that these universes share the property that they implement only a single level of selection. Given that the theory of major evolutionary transitions proposes that interactions *between* hierarchical levels of selection has been critical to the evolutionary growth of complexity in real biology, it is at least worth investigating the implementation of some abstract version of this kind of dynamical hierarchy in artificial evolutionary systems.

The purpose of Part II (Chapter 4) is to detail the construction of the experimental modelling platform which was designed and built for the purpose of carrying out the experiments described later. Chapter 4 outlines the design specification for this new artificial life platform, called “Molecular Classifier System” (MCS). The MCS is an artificial chemistry inspired by the origin of life as we know it and the

current efforts of wet-lab chemists to fabricate artificial protocells in-vitro. More specifically, the MCS is inspired by the RNA-World hypothesis (Gilbert, 1986) for the origin of life and is loosely based on the “learning classifier systems” model presented by John Holland (1976). It is intended that this modelling platform will be able to simulate the kinds of protocells which might be expected to occur at or near Maynard-Smith and Szathmáry’s hypothesized first Major Evolutionary Transition—from populations of “naked” replicators to populations of such replicators constrained to co-exist within cells which compete with other similarly formed cells. It is argued that in order to fully understand the potential of this hierarchical system, one must first understand the evolutionary dynamics of the pre- and post-transition worlds independently, and recognise that at the very least during such a transition, the evolutionary dynamics observed at the cellular level are driven for the most part by the evolutionary dynamics at the sub-cellular level. The MCS system as presented in this thesis provides a platform for the investigation of the evolutionary dynamics at the protocell level which is only feasible due to drastic simplification of the chemistry that is involved.

The results of a sequence of experiments which were carried out using MCS are presented in Part III (Chapters 5 and 6). Chapter 5 describes a set of experiments which were carried out to validate the MCS model. In particular, this chapter is focused on the dynamics of individual replicating molecules and how they interact and compete with one another. This early implementation of the MCS platform is drastically simplified in comparison to any plausible real-world chemical system, and stands mainly as a proof of concept for the platform. Once the evolutionary dynamics at the pre- transition level have been explored, hierarchical selection is then applied in the form of artificial protocells, where the novel phenomenon of “selectional stalemate” occurs. In spite of the fact that this stalemate is entirely due to the simplifications employed by the MCS, a demonstration of such a phenomenon has not been demonstrated by other artificial system to date, and is an interesting result in its own right for the MCS toy-model.

Chapter 6 describes the effects of some modifications to the molecular chemistry which were applied to make the molecular level evolutionary dynamics slightly more realistic, though it should once again be emphasised that even in its most advanced form, the MCS is still drastically simplified in comparison with plausibly realistic chemical systems. These modifications result in the achievement of a stable protocell platform, MCS-1, which demonstrates ongoing evolution at the protocell level, driven by the dynamics of the molecules contained within them. A side-effect of these modifications is the disruption of the selectional stalemate that was discussed in Chapter 5. This chapter culminates with the presentation of a detailed descrip-

tion of three qualitatively distinct protocell level dynamics that have been observed using the platform, and that would not have been possible without understanding the molecular level dynamics in its own right and separately from the behaviour observable at the protocell level.

Part IV (Chapter 7) contains some further discussion of the results presented, describes some of the potential avenues for further study that were identified during the course of the research and provides some concluding remarks about the scientific contribution of the work.

Part V contains some appendix material which describes an extension of the MCS artificial protocell system. Appendix A presents a supplementary extension of the literature discussed in Part I. In this appendix, the computational potential of living systems is explored. A “living-system” based computational device would presumably have the capability of performing true parallel computations and would have all of the robustness of a self-repairing, self-maintaining system. Cell signalling networks are described here as a biological example of the complex types of computation that can emerge through the process of evolution. One of the primary objectives of the PACE project, which partially funded the research presented in this thesis, was to investigate the computational potential of artificial protocells, in anticipation of their realisation in the laboratory. The literature reviewed in this chapter is intended to guide expectations as to the kinds of computation that may be achieved by artificial protocells, and to hold up evolution by natural selection as one method for programming such devices. Following that, in Appendix B, the results of a series of experiments are presented in which the MCS-1 artificial protocell system is extended to support a very primitive computational function. It is shown that the protocells demonstrated the ability to “compute” a property of their environment, which was not directly observable from the molecular level of the hierarchy. Cells which are better at carrying out this computation are “fitter” than those that are not, and evolution against an externally provided fitness function sees to it that the population is dominated by cells which are best at carrying out this computation. Such simple forms of computation may be the origin of the complicated Cell Signalling Networks (CSN) that modern biological cells use to regulate the different sub-systems that they need to stay alive.

1.4 Associated Publications

The following is a list of publications that have been produced during this programme of research.

1.4.1 Conference Proceedings

- “Enrichment of Interaction Rules in a String-Based Artificial Chemistry”, *Eleventh International Conference on Artificial Life, (ALIFE XI)*. (Kelly et. al. 2008)
- “Multi-Level Selectional Stalemate in a Simple Artificial Chemistry”, *European Conference on Artificial Life, (ECAL9)*. (McMullin et. al. 2007a)
- “Preliminary Steps Toward Artificial Protocell Computation”, *International Conference on Morphological Computation*. (McMullin et. al. 2007b)

1.4.2 Posters

- “The Evolution of Complexity in a Multi-Level Artificial Chemistry.” *European Conference on Complex Systems, 2007 (ECCS-07)*. (Kelly et. al. 2007)
- “Cellular Computation Using Classifier Systems”, *International Workshop on Systems Biology, 2006*. (Kelly et. al. 2006a)
- “On Protocell Computation”, *European Conference on Complex Systems, 2006 (ECCS-06)*. (Kelly et. al. 2006b)

Part I

Literature Review

Chapter 2

Life

Life, as we know it, has evolved over billions of years into the rich diversity of the biosphere today. In this chapter, I present a review of some of the key elements which compose our current understanding of life on Earth. The study of life as we know it is a highly inter-disciplinary area drawing upon research from fields such as Chemistry, Biology, Physics and Informatics. The following sections present the key supporting literature in this field, particularly that examining what exactly “life” is and why we should be interested in understanding the key sub-systems that give rise to life. After a brief philosophical examination of the lay-view of life, I review some candidate definitions that have been proposed for life itself, one of which is adopted throughout the remainder of this thesis. It will be argued later in Chapter 3 that it has been an important goal of artificial life to reproduce the phenomenon of unambiguous evolutionary growth of complexity in computational models of living processes and in other artificial living systems. This chapter will explore and investigate some characteristics of life as we know it that have not previously featured in artificial systems, particularly those characteristics which are believed to have been important contributors to open-ended evolutionary growth of complexity.

The theory of major evolutionary transitions is presented with a focus on arguably the most interesting transition—from non-living to living systems. As we will see, it is notoriously difficult to attribute a precise description of what it means to be alive, though it is not so difficult to label things as alive or not on an individual basis. The first hypothesized major transition, from replicating molecules to populations of replicating molecules constrained to “live” together inside containers (cells) and to share a common fate, seems to span this conceptual gap between the non-living and the living. For life as we know it, this transition would have been continuous: at one point, there were individual naked replicating molecules, and at some later point, there were cellular structures which contained populations of

replicating molecules and which were *themselves* subject to evolution by natural selection. This thesis argues that in order to better understand the dynamics *during* such a transition, it is necessary to separately understand the dynamics of the pre-transition and post-transition worlds. The final sections of this chapter explore some of the key systemic problems that life necessarily overcame, along with some of the proposed solutions to how evolution may have solved these problems.

2.1 What is Life?

Life, as we know it, is reliant on DNA-based chemistry. The field of genetics, the study of heredity by the mechanism of DNA, emerged from man’s desire to understand the sub-systems of living processes. The elegant “machinery” of DNA-based chemistry arose from an evolutionary process, and is ubiquitous across all living things. Unfortunately, the fossil record cannot illuminate the earliest phase of evolution—the origin of life. In fact, there are many competing theories for how life might have begun. After a brief discussion about the conditions of early Earth, the key candidate theories for the origin of life will be presented, though it is argued by some that a synthesis of two or more of these theories can be used to tell a more plausible origin story (Martin, 2003). It is not the purpose here to add to the origin-of-life debate, but rather to take inspiration from that debate and apply it to understanding *artificial* life.

2.1.1 The “definition” problem

When first presented with the task of defining life, or to state the most important characteristic of life, it is easy to resort to what we know best, i.e., it needs to grow, or it needs to breath air, or it needs to reproduce. There are of course many examples of living things which do not meet these or other similar criteria. These are common fallacies in defining life (Farmer and Belin, 1990). In spite of these common fallacies, it is usually straightforward, even for a lay-person to choose one of the following labels for an object:

- living,
- not living.

Furthermore, in the case of objects labeled “not living”, it is fairly straightforward to state whether at any historic point that particular object could *ever* have been labeled as “living”. Fossils, which after all are just pieces of rock, can easily be seen as representing something that at one stage was alive, and by a process of fossilization

has been preserved in rock. On the other hand, a plain old rock would not inspire an explanation which involved living origins—it’s just an inanimate piece of rock that never lived, and does not represent anything that *ever* lived.

Bedau(1996; 1999) asks whether it is even important to have a precise definition for life, and further questions whether one is even possible. Life is a qualitative property—we cannot measure “how alive” something is—it is either alive or not—and further, it either was alive at some stage or it was not. Farmer and Belin (1990) maintain that for every acceptable definition, under *some* criteria, there is an example of a living thing which breaks the definition. Objects can easily be categorized on a piecemeal basis, but formulating a precise definition to cover all cases is very difficult. In other words, we can identify life when we see it, but before that, it seems impossible to have a yard-stick against which we can measure how alive something is or how close something is to being alive.

2.1.2 Definitions can be restrictive

It would be extremely naïve to suggest that life can only exist in the way that we currently know it. The history of life on this planet coupled with human knowledge of the science of biochemistry strongly suggests that life¹ came about as a series of “accidents”. For example, it would seem that the particular genetic code used by life to translate genomic information into proteins could have taken many different forms, and still produce the same qualitative result (Stegmann, 2004). Interestingly, almost all living things share the *same* genetic code, and all of the genetic codes are qualitatively similar (Osawa et al., 1992; Jukes and Osawa, 1993). This fact raises two questions:

1. whether there was at one stage in the history of life, competition between life forms which had differing genetic codes.
2. whether all surviving forms of life share a common ancestor (Last Universal Common Ancestor, *LUCA* (Forterre and Philippe, 1999)).

In theory then at least, it seems life could exist in fundamentally different material forms than those we currently know about. Deoxyribonucleic acid (DNA) is the common hereditary information carrier, and the fact that all life uses it suggests that it is particularly suited to this role. However, there is nothing to suggest that living things *in general* require a DNA-based information carrier. Peptide Nucleic Acid (PNA) is just one example of another family of polymers which seem to be capable of serving the purpose (Nielsen, 1993). Any satisfactory definition of life

¹From this point forward, unless otherwise stated, “life” is synonymous with “life as we know it”.

must therefore respect the possibility that all living things do not necessarily need to share the same basic components as the life-forms we currently know about.

2.1.3 A satisfactory definition

Maynard-Smith and Szathmáry(1997) present the following as prerequisites for an entity to be either capable of, or the result of, a process of evolution by natural selection. The entity in question must:

- be capable of producing other things like it,
- be capable of passing on differences in the instructions which are used to “build” these offspring,
- *or*, at the very least, be descended from such an entity.

These criteria are true in the case of *everything* that we know to be alive. Perhaps then the ability to evolve, or being the result of an evolutionary process is sufficient information to distinguish between “living” and “non-living”. In the next section, we shall see the importance of the process of evolution by natural selection to the growth of complexity and production of diversity of life on Earth.

2.2 Darwin’s Insight

Life exists in many different shapes and forms. Biologists categorize and classify life-forms into groups to make it easier to understand the lineages and common ancestry of all organisms on Earth. The end result of such classification is a “family tree” of life, with related life-forms appearing close to each other and more distant relatives appearing further apart. As might be expected, a “family tree” which related all life on Earth would have many “leaves” representing the diversity of life, but what of the “roots” of this tree? As we travel back in time from leaf to branch to trunk to root, we strip away the differences between the organisms and the similarities become apparent. Charles Darwin wondered how such a “tree” might have grown and branched over the millenia since life began. Darwin’s ideas developed during his 5-year round-the-world voyage on HMS Beagle between 27 December 1831 and 2 October 1836, Darwin (1845). During that time, he collected specimens of interesting creatures which he sent back to England for expert appraisal—Darwin himself being an amateur naturalist at best. One such appraisal was of a set of 14 birds that Darwin had collected on the Galapagos Islands. He had thought them to be a mixture of different variations of the same species of bird, but upon examination by a suitable expert, the birds were found to be closely related but nonetheless

distinct species of finch, classified by the size and shape of their beaks. Darwin, upon learning that they were to be considered separate species rather than varieties of the same species, wondered how such profound differences could manifest. He concluded that due to the physical separation of their habitats—the finches lived on the Galapagos Archipelago—adaptation to differing food-sources caused the finches to evolve differing beak sizes and shapes. Darwin reasoned that over time, the finches which were better suited to eating the particular kinds of food available to them were better at “being finches” in the given environment, and therefore would have had more offspring who also would have been good at eating that particular type of food. It was commonly accepted that children would resemble their parents and siblings—the mechanism for these heritable variations however was still not understood. Darwin’s initial insight was to suggest that physical separation between similar groups of individuals would eventually lead to the emergence of new distinct species. The crowning glory of Darwin’s contribution however was his further suggestion that the diversity of *all* things that have ever lived could be explained by a sequence of speciation events and that all living things therefore shared a common ancestor. Darwin (1859) proposed “natural selection” as the process by which these speciation events were driven, and compared it to the artificial selection that farmers had used since neolithic times to improve the quality of their breeding stock.

2.2.1 Evolution by Natural Selection

The identification of “natural selection”, or the “survival of the fittest” as a critical factor behind the diversity of life was Darwin’s greatest insight. Evolution by natural selection is the best known explanation for the emergence of new traits in species and how species arise. Evolution by natural selection relies on some important pre-suppositions, which are closely related to Maynard-Smith and Szathmárys framework presented earlier:

- It must be possible to produce offspring, and these offspring must be similar to their parent. This point means that an offspring should be more similar to its parent than it is to an unrelated member of the same species.
- It must be possible to have variation in the offspring, and that at least some of these variations are heritable. In other words, the offspring should mostly resemble the parent and for the ways in which it differs, there must be a way to pass on such differences to the next generation of offspring otherwise selective differences could not arise.

Individuals of a given species generally resemble each other much more than they would resemble members of another species. However, in spite of their similarities, it

is the differences between these individuals that provide the grounds for determining “fitness” for survival. It is not controversial to suggest that the differences between individuals of a given species are irrelevant for the most part. Opponents of the theory of evolution by natural selection do not believe, however, that these small differences between closely related individuals is enough to cause the great diversity of life we see on Earth. Examples such as the mammalian eye are held up as “proof” that some structures are so complicated and intricate that there could not possibly be any intermediate between “no eyes” and “having eyes” that could serve as a stepping stone for the evolutionary process—“What use is half an eye?”

Responding to such claims, evolutionary biologists, among others, have argued that the fitness difference between an organism that has no mechanism for detecting the presence or absence of light and one that has a single instance of light responsive tissue is enough to initiate an evolutionary process (Dawkins, 1991). Again, this is uncontroversial, since the ability to notice small changes in light levels might be enough to give a slight advantage to those who possess this ability, allowing them to produce more offspring similar to themselves. Supporters of the theory of evolution by natural selection understand that this small difference alone may not be enough to ensure the *absolute* survival of the “light-sensing” creatures, but also understand that *if*, on average, the creatures who can sense light do better than those who don’t, the population will tend to be dominated by the lineage of “light-sensing” creatures. It is clear that this rudimentary “light-sensing” ability is far inferior to fully developed eyes as we know them, but by this process of incremental improvement, evolution by natural selection can “design” complicated structures such as eyes if given enough time. It is important to remember that evolution produces novelty by a process of *blind* variation. There is no higher level process which decides what areas need further “design tweaks”, but rather there is an ongoing production of novelty through small differences in the heritable information that is passed from generation to generation. Under the conditions of limited resources, for example, limitations of space, food or mating partners, heritable characteristics which improve the survival chances of individuals will increase in frequency across the population over successive generations. Conversely, heritable characteristics which reduce the survival prospects of individuals will decrease in frequency.

2.2.2 Heritability

Children resemble their parents and siblings much more closely than they do unrelated individuals. This is another uncontroversial statement, and there are many examples of such inherited traits in humans:

- height

- eye-colour
- hair-colour
- skin-colour

Modern science attributes the first systematic study and subsequent documentation of this phenomenon to an Augustinian Friar, Gregor Mendel (1866). Mendel studied the hybridization—crossing of species—of Garden Peas and is known as the “father of modern genetics”. Mendel noticed that rather than the traditionally accepted “blending” theory of inheritance which stated that offspring would manifest at a point on a continuum between the extremes of the parents, a more discrete, “particulate” theory of inheritance seemed to be underpinning the hybridization results he was witnessing. Mendel also documented a dominant/recessive gene division to explain how subsequent generations could display a trait that neither parent displayed. Thanks to Mendel, the gene theory of information heritability was understood, and the mechanism by which fitness differences could be passed down to offspring was discovered. It would take another 50 years, however, before the physical nature of genes was discovered.

2.3 Information Chemistry

This is a thesis about *artificial* Life—life as it could be. Life, is invariably based upon the intricate chemistry of DNA-based genetic information storage. The genome of an organism—the list of instructions which define the organism—is represented in a quantity of DNA which essentially contains all of the information required to reconstruct the organism, given the appropriate supporting chemistry. Instructions are stored in DNA in the form of genes on chromosomes, and the DNA is made up from pairs of nucleic acid molecules which are bonded in the now familiar double helix shape, as discovered by Watson and Crick (1953). This discovery uncovered the mechanism by which heredity worked. Since Mendel and his peas, scientists had been trying to understand *how* traits were passed on from parent to offspring—DNA provided the answer. As we saw earlier in Section 2.1.3, the transmission of hereditary information between successive generations is the cornerstone in our satisfactory definition of life. Later in Chapter 3, it is argued that artificial life models have typically given a lot of attention to the abstraction of the procedures used for this transmission of information between generations. Of course, in computational models, these procedures can be made as simple or as complicated as researchers desire. In this thesis, I explore the possibility that there are other more important

abstractions that can be made, taking inspiration from the major evolutionary transitions framework in Origin of Life research which suggests that natural selection acting simultaneously at multiple hierarchical levels has been critical to the evolutionary growth of complexity in Life as we know it. In the following sections I will give a brief description of the complicated information chemistry which underpins the transmission of genetic information in all living things. The point of this section is both to highlight the intricate beauty of this system and to argue that although Life has evolved this complicated mechanism for heredity, there is nothing especially “vital” about it, and that artificial life models might safely abstract most, if not all of the details.

The nucleic acid monomer bases which make up DNA are Adenine, Cytosine, Guanine and Thymine (Watson and Crick, 1953). Adenine bonds to Thymine and vice-versa, and Cytosine bonds to Guanine, and vice-versa. DNA is stored in the form of chromosomes in cells and, in eukaryotes, is housed inside a cell nucleus. When information is needed from the DNA, helper enzymes cause the relevant part of the DNA to unwind, allowing another enzyme to make a copy of whatever portion of the DNA is required, after which the DNA recoils into its double helix form. The copied segment can then be transported to another part of the cell to be used, while the original structure remains unmodified.

In organisms with diploid cells (cells with two complete sets of chromosomes), offspring are produced when an egg from the female is fertilized by a sperm from the male. Each chromosome has two halves, one of which is inherited from the father, and one from the mother. When an egg is created, it will be given half of the DNA needed to build the new offspring, and in a similar fashion, each sperm will have the other half of the DNA needed. In this way, half of our DNA comes from our father, and half from our mother, ensuring that we inherit our parents’ genetic information in roughly equal proportions. Upon fertilisation, the internal state of the resulting *zygote* will prevent entry by other sperm, and the *zygote* will begin to grow and divide, replicating and expressing the inherited parental DNA at each cell division.

2.3.1 DNA Replication

DNA replication is a high fidelity process. It has the ability to recover from copying errors such that on average, only about one base pair per billion copied is in error. For example, the human genome consists of 2×10^9 base pairs, so on average, there will be two copy errors each time our DNA is replicated (Nachman and Crowell, 2000).

In order to carry out the complicated process of replication, DNA relies on a network of enzymes which assist at various stages along the way. The process first

sees the double-helix structure being unwound to expose the individual nucleotides in what is known as a replication-fork. Further enzymes are then called upon to promote the uptake of raw material nucleotide bases and cause them to bond to the exposed strands of the unwinding helix. The unwinding of the helix would ordinarily lead to tangles and knots in the strands which would potentially bring the process to a halt. However, there is also a set of enzymes, known as “topoisomerases”, which solve this problem by carefully cutting and rejoining the DNA strands at strategic locations to prevent the build up of these “tangles” (Wang, 1991).

2.3.2 Transcription

The genome contains the “recipe” for the components of the organism it belongs to in addition to specifying where and when to build each component. It is clear that the chemistry involved will need some way of extracting the appropriate information at the right time to carry out whatever task needs to be carried out. For example, skin cells and neurons, the cells which make up brains, contain identical copies of the genome, yet both types of cell are required to perform very different tasks for the organism. Furthermore, as mentioned earlier, in eukaryotes the DNA is sequestered within the cell nucleus. Therefore, a clear requirement exists for firstly extracting appropriate information from the genome and secondly, transporting that information outside the cell nucleus to direct further processes as and where they are required. This process is known as *Transcription*, and sees the genomic information being “transcribed” into an intermediate form using mRNA (messenger RNA) (Brenner et al., 1961). RNA (Ribonucleic Acid) is similar to DNA in that it is a polymer made up from chains of nucleotides. There are some distinctions however, most notably, RNA is a single stranded molecule which forms distinctive hair-pin loops, and the nucleotides which compose RNA see Uracil replace the Thymine that DNA uses.

Once again, a complicated set of enzymes are called upon to direct the process. Firstly, the appropriate segment of DNA must be located and “primed” to form an “action bubble” in a similar fashion to the formation of the “replication-fork” during DNA Replication. This stage is known as the initiation stage. Following this, an enzyme known as *RNA polymerase* traverses one strand of the DNA, unzipping the DNA in front of it, re-zipping the DNA behind it, and using complementary base-pairing to make an exact RNA copy of the DNA, with the exception that Thymine is replaced by Uracil, as mentioned above. This proceeds until the polymerase reaches a “stop” marker in the DNA, at which point the polymerase detaches from the DNA, ensuring that it is re-zipped and the newly formed RNA chain leaves the nucleus and travels to some other part of the cell where it can be translated into the amino

acid chains that form proteins.

2.3.3 Translation

The process of transcription gives rise to sequences of mRNA which contain representations of genes and are *outside* the cell nucleus. These sequences are now available for use by the “translation” process which initiates the production of the various proteins that are the building blocks of living things. This translation process happens by way of ribosomes, which are made from RNA and proteins, that surround the mRNA strand and process it sequentially in accordance with the genetic code, which defines a mapping between nucleotide triplets—codons—and amino acids (Nirenberg et al., 1962). In fact, the amino acids are bound to transfer RNA (tRNA) molecules which present a complementary structure to the mRNA which has just been read by the ribosome. The pattern described by the tRNA is known as an anti-codon. The anti-codon binds to the codon that is being processed by the ribosome, by complementary base-pairing, and the amino acid which was bound to the tRNA becomes attached to a growing chain of amino acids at the ribosome. This amino acid chain will eventually become folded into the appropriate protein which was described in the genome.

2.3.3.1 Genetic Code

The genetic code is the mapping between mRNA codons and amino acids. This mapping is facilitated by tRNA anti-codons which bind by complementary base-pairing to the mRNA codons. There is a “many-to-one” mapping between amino acids and tRNA anti-codons. There are 20 amino acids in total, and $4^3 = 64$ distinct anti-codons since they are formed by triples over the set {Adenine, Cytosine, Guanine, Uracil} (Brimacombe et al., 1965; Watson and Berry, 2003).

2.3.4 DNA—A Ubiquitous Language?

All life depends on storing heritable information in a DNA-based genome, which is replicated, transcribed and translated into proteins by the mechanisms just described. The chemistry of DNA is the same for all things that use it. Squid DNA, Human DNA, Cabbage DNA—in each case, the information contained in the genome is different but the underlying chemistry is identical.

The combination of the facts that all living things use DNA in the same way, and are built from the same protein building blocks strongly suggests that all of the living things we see share a common origin. Considering the complicated nature of DNA chemistry, and taking into account the fact that the mapping from DNA

base-pair to amino acid, the genetic code, could theoretically exist in numerous different ways: to suggest that such things could have happened by chance alone is quite implausible due to the vast number of living things that exist, and that have existed.

2.4 The Origin of Life

According to Darwin, it might be possible to trace back through the tree of life until we reach a “common ancestor” that is shared by all living things (Darwin, 1859). Unfortunately, rigorously testing this theory has proven to be a difficult task. In terms of evidence, the fossil record, at best, can take us back 2.7 billion years (Lepot et al., 2008) to a time when single-celled cyanobacteria which gathered energy by photosynthesis (Olson, 2006) were flourishing. This was a time when Earth was roughly half as old as it is now. It is thought that these cyanobacteria had the side-effect of oxygenating the Earth’s atmosphere, which at the time was poisonous to the existing single-celled organisms. However, by this stage, we already had living entities, so we need to go back even further to examine the “*origin* of life”. The fact that there is no hard evidence of what life was like before these cyanobacteria were fossilized means that there is a great deal of uncertainty about the precursors of these primitive life-forms. The approach has therefore been to hypothesize the conditions of early Earth, and then to explore plausible means for initiating an evolutionary process from simple beginnings which could potentially give rise to the evolutionary growth of complexity which produced the diversity of life on Earth.

2.4.1 Overview of timescale of Origin of Life

Though there is a great deal of uncertainty surrounding the origin of life, it is possible to build a time-frame of important events on Earth that would certainly have had an effect on the earliest forms of life, without needing a clear picture of what these early organisms were. The time-frame begins with the formation of the Earth about 4.5 bya (billion years ago). It is estimated that it would have taken about 400 million years for the Earth to cool enough to have a solid surface, oceans and a gaseous atmosphere. However, the Earth remained under heavy bombardment by meteorites until about 3.9 bya, causing a continuous evaporation and re-condensation of the oceans—conditions which certainly were not conducive to life. The oldest known fossils are approximately 2.7 billion years old so we know that life existed by then at the latest. In light of these uncertainties, we are forced to make a best-guess that life appeared on Earth somewhere between 3.9 bya and 2.7 bya.

2.4.2 Major Evolutionary Transitions

Primitive cyanobacteria are the earliest *known* life-forms that existed on Earth, approximately 2.7 bya according to Lepot et al. (2008). We can examine the key distinguishing features of such early life-forms to create a picture of what the key components are which they possess that bestow the characterisation of “living” upon them. Once these key components have been identified, the process of determining which components could exist as entities in their own right is begun. Miller (1955) demonstrated that the environmental conditions of prebiotic Earth could be simulated in-vitro and could lead to the production of various organic compounds such as amino acids, carboxy acids and hydroxy acids. This is a very important finding, as these are structurally complicated macro-molecules which are critical to life. However, leaving aside the fact that this was just a small subset of the organic molecules that life relies upon, these molecules in their own right could not be said to constitute living units.

JD Bernal presented a view of the origin of life which he called “biopoesis” (Bernal, 1949, 1959) for which there were three distinct stages:

1. The origin of organic monomers,
2. The origin of organic polymers,
3. The transition from individual molecules to populations of molecules within cells.

Perhaps Miller’s approach can help account for some of the advances of the first two stages of biopoesis, but the third stage has yet to be demonstrated by anything other than life as we know it.

Maynard Smith and Szathmáry (1997) described a framework of *major evolutionary transitions* in their seminal work. It was their opinion that there have been a series of major transitions in evolution which have been critical to the evolutionary growth of complexity which life has achieved. According to Maynard-Smith and Szathmáry, major evolutionary transitions are characterised by the emergence of a new higher level Darwinian *actor* composed of an aggregation or cooperation of lower level components. This implies that there is some notion of hierarchy involved. Essentially, a Darwinian process at a lower hierarchical level is subsumed into a higher level functional unit, resulting in a nett fitness improvement that would have been difficult, if not impossible to achieve without such a transition. They list the following major transitions that potentially played a crucial role in the evolution of life:

1. From replicating molecules to populations of molecules in compartments

2. From independent replicators to replicators bound together in chromosomes
3. From RNA acting as genetic material with enzymatic function to DNA as genetic material and separate proteins as enzymes
4. From Prokaryotes (cells with no nucleus) to Eukaryotes (cells with a nucleus)
5. From reproduction by asexual cloning to sexual reproduction
6. From single-celled protists to multi-cellular organisms
7. From solitary individuals to colonies with non-reproductive castes
8. From primate societies to human societies with language

The first major transition in the framework aligns with the third stage of Bernal’s biopoesis. Maynard-Smith and Szathmáry would therefore place the requirement of “replication” upon the first two stages of “biopoesis”, Bernal himself recognising that natural selection would probably have acted upon molecules at the first two stages of biopoesis. So, in terms of the hierarchy that was referred to above, and speaking specifically about the first major transition, it is clear that there will be a process of Darwinian selection being applied directly to the molecules as they compete with one another, but also there is a higher level of Darwinian selection being applied indirectly to the molecules. Assuming limited resources at the cellular level, then cells will also compete with one another. In Section 2.3, we had a flavour of the complicated chemical processes that are at work in modern biological cells. Presumably, at or near the first major transition, cellular chemistry would have initially been uncoupled with the molecular chemistry that it contained, and so the fate of the cell is highly dependent on the activity of the internal molecular population. Selectional pressure at the cellular level might be expected to favour molecular aggregations which were somehow better for the collective. This kind of *hierarchical selection* (Buss, 1987) is a key feature of the artificial life system which is described in the later chapters of this thesis.

2.4.3 The Transition to Cells

Life, as we know it, is an example of multiple interacting, and in many cases, hierarchically inter-dependent, systems. Given that, then the more interesting phase of the origin of life was the origin of “cellular life”, as described by the third stage of Bernal’s “biopoesis” and the first of Maynard-Smith and Szathmáry’s “Major Transitions”, as it is at this point that the effects of hierarchical selection begin to take hold.

Cellular life is characterised by the presence of a self-sustaining, self-maintaining, self-replicating chemistry within the boundary of a containing membrane. The interesting question for origin of life researchers then is to explain how the transition to “cellular life” might have occurred, since this arguably marks the transition between a collection of un-integrated, non-living components and a fully integrated collection of components which can undergo a process of Darwinian evolution in its own right—a major evolutionary transition. One systemic view of cellular life (Rasmussen et al., 2008) identifies three key sub-systems which interact:

1. metabolism
2. information chemistry
3. containment

The upcoming sections will describe research into the capabilities for each of the above identified sub-systems in isolation. Furthermore, there has been speculation as to the various pathways that might have been followed to get from one or more of these sub-systems existing in isolation to collections or aggregates co-existing as examples of what would be recognisable as cellular life. Martin (2003) presents a model of the evolution of cellular life which allows for the co-evolution of the various sub-systems in the vicinity of deep-sea hydrothermal vents, and then predicts the subsequent co-location and aggregation of these systems into primitive “protocells”. The following sections highlight key research in the areas of each individual sub-system, and form the basis of our current understanding of the kinds of chemistry that may have led up to the first major transition.

2.4.3.1 Metabolism First?

The ability to extract energy from the environment and convert it into a usable form is a key component of life, and moreover, it is difficult to imagine any form of *life-as-it-could-be* which did not require some ability to process or use energy. Metabolic systems can therefore be seen as one of the primitive components of living-systems. Wächtershäuser (1990, 2007) proposed that such primitive metabolic systems may have arisen spontaneously at deep-sea hydrothermal vents, under high temperatures and pressures. Specifically, Wächtershäuser mentioned specific iron and sulfur compounds that underwent chemical reactions on mineral surfaces near these deep-sea hydrothermal vents. The products of these chemical reactions would in turn undergo further reactions, and the number of molecules would grow via an auto-catalytic process. Depending on the individual chemical reactions, it might also be possible for the number of types of molecule to also grow, by an evolutionary

process. In addition, the products of these reactions would then be available to external (external to the “metabolic system”) molecules for use as an energy source. Martin (2003) proposes that these metabolic systems could have arisen on the surface of micro-porous rocks at the hydrothermal vents, in effect benefiting from an artificial container.

2.4.3.2 Replicators First?

Information representation and replication is perhaps the most important sub-system of life as we know it. Evolution by natural selection relies on the heritability of information—a property which is accomplished through genes and chromosomes built from DNA in the case of life as we know it. Section 2.3 showed that the DNA chemistry which maintains genomic information in all known living things is extremely complex, and relies on a large set of enzymes and proteins to assist and enable the various error-correcting and proof-reading tasks associated with high-fidelity information replication, and another large set of helper molecules to carry out the transcription and translation of the genomic information into the protein building blocks its living host requires. Theories which claim that life began via replicators such as these are known as “replicators first” theories, and the following are two examples of such.

RNA-World RNA plays a key supporting role in the DNA chemistry which supports life, and there is significant support for the hypothesis that RNA played an important role in the pre-protein, pre-error-correction history of the Origin of Life (Kruger et al., 1982; Joyce, 1991; Bartel and Unrau, 1999; Orgel, 2004). The importance of RNA to the DNA replication/error-correction apparatus suggests perhaps that evolution has favoured DNA as a more stable information carrier and relegated RNA to a support role. This theory has gained recognition as the “RNA-World Hypothesis” and it holds that there was a transitionary phase in the origin of life where a substance, RNA, acted as both information carrier and enzyme.

Clay-Crystals In an alternative “Replicators First” theory for the origin of life, Cairns-Smith (1982) hypothesised that the earliest replicators were clay-based crystalline structures. His theory was that such crystals would spontaneously form on riverbeds, and grow through the ongoing process of sedimentation. By their nature, crystals are self-similar—like-creates-like. Replication of such a crystal might then occur when it gets too large, and breaks into separate crystals, which in turn would undergo a similar growth and division process. These crystals could have an environmental effect on the riverbeds where they found themselves, for example, if they clumped

together, they might cause heavier sedimentation, leading to the creation of pools of water. Crystals downstream of these “blockages” would then experience a shortage of water, or even complete drought. This would lead to the drying up of the river bed and the potential for crystals to be swept away as dust in the wind, to potentially colonise a different suitable environment. A further elaboration of this theory suggests that these clay crystals may have acted as scaffolding for the construction of more complicated molecules. The shape of the crystal might hold certain molecules near to each other for long enough that they might form a bond between themselves, and subsequently move away from the scaffolding crystal, freeing up the original holding positions for further construction operations.

2.4.3.3 Containment First?

The final sub-system that is considered essential for life is that of containment. For contemporary biological cells, this containment comes in the form of lipid membranes (Overton, 1895). Lipids are long oily molecules, and have the property of self-organization under certain environmental conditions (i.e., pH, temperature, pressure). Some of these molecules are *amphiphilic*, meaning that they have both hydrophilic and hydrophobic tendencies. In other words, for each of these long oily molecules, one end favours contact with water, while the other end avoids it. It is this property which enables the self-organization of lipids into various structures, for example, micelles, bilayers and vesicles. Micelles are typically spherical and form when a group of lipids merge to present their hydrophilic tails to an aqueous environment, which causes the hydro-phobic tails to meet in the center—an “oil-in-water” micelle. Bilayers are formed from two layers of lipids, joined by their hydro-phobic tails such that both sides of the sheet present hydrophilic heads. Such a bilayer may form into a spherical shape, a vesicle, which is essentially the type of membrane that is utilised by contemporary biological cells.

The necessity of a containment mechanism for the success of fledgling life coupled with the self-organizational properties of lipid molecules inspired the theory that such chemistry might have had a contributory influence on the origin of life (Segré et al., 2001). Lipids have been found to have catalytic properties, a characteristic that seems to have worked well for the RNA world hypothesis. Furthermore, lipid membranes tend to give rise to more copies of themselves, thanks to these catalytic and self-organizational properties. In this way, heritable information can be passed between parent and daughter structures, in a fashion similar to the mechanism proposed by Cairns-Smith for his clay-crystals.

2.5 Challenges to Life

In the following sub-sections, some of the problems that shaped the early evolution of life on Earth are described. These problems are separated into environmental challenges and organisational challenges, and are followed by some proposed solutions from the literature.

2.5.1 Environmental Challenges

There were many key milestones in the evolution of *prebiotic* Earth that would later be very important for the appearance of life, for example:

1. the many volcanoes which brought up important chemical elements from deep beneath the Earth's surface,
2. the sufficient cooling of the Earth's surface to allow the condensation of the oceans,
3. the emergence of deep-sea hydro thermal vents which arguably have played a central role in abiogenesis and the evolution of early life (Section 2.4.3.1),
4. the availability of suitable building-block molecules to permit the complicated chemistry that life requires.

Miller (1955) showed that even taking fairly advanced starting points, we can still end up falling foul of chance. Replaying the tape, as best as we can, might still produce results that are unexpected. Miller's experiment produced some, but not all, of the organic compounds that are necessary for life. He attempted to simulate the conditions of pre-biotic Earth and explore the production of the organic compounds that life requires, and had at least some success. In spite of the progress made in understanding life, it is still not understood to the point where all of the required steps can be reproduced in a laboratory.

2.5.2 Organisational Challenges

If it were possible to abstract the nature of life away from the confines of Physics and Chemistry we might be able to reduce the problem of the origin of life to a problem of organisations, systems and interactions. In many ways, this abstraction makes the challenge of understanding life even *more* difficult, as it begins to approach the territory of the general definition of life, which has already been identified as an area which has many open questions (Section 2.1). Life, it would seem, had many more problems to overcome than the ones it faced due to the precise details of the

chemistry and physics involved and in particular, the general problems surrounding the reproduction of information. Molecular replication events can be generalised into two categories: “trans-”acting replication and “cis-”acting replication (Altmeyer et al., 2004), where “trans-” and “cis-” come from the Latin prefixes meaning “on the opposite side” and “on the same side” respectively. These two classes of replication produce dramatically different growth rates for the molecular species involved—cis-acting replication leads to exponential growth, whereas trans-acting replication leads to hyperbolic growth.

2.5.2.1 The Information Integration Problem

In Section 2.3, the intricate chemistry of DNA was explored at a high level. We saw then that the high-fidelity replication abilities of DNA were enabled by a number of other compounds. Upon closer examination, it becomes clear that not only does DNA rely upon the support of these other compounds, it actually encodes the instructions by which to make these helper molecules. The bacterium *Mycoplasma genitalium* is the smallest known free-living bacterium and was chosen by The Minimal Genome Project (Glass et al., 2006) for its study to find a minimal genome which had all the necessary information to support and sustain life. Even though this is the smallest known set of naturally occurring genetic information which can support life, it consists of 582,970 base-pairs making up 521 genes (Fraser et al., 1995). Nevertheless, this set of genes contains the information required for DNA replication, error-correction, DNA transcription and translation into proteins which keep the organism alive through functions not limited to energy metabolism and maintenance of containment. Eigen (1971) presented a “chicken and egg” paradox for the origin of this genetic system, which has become known as Eigen’s Paradox. The paradox arises from the interdependence of information preservation and long genomes. Eigen showed that for an information sub-system based on nucleic acid base-pairs *without* error-correction enzymes or proteins, the maximum sustainable genome length is of the order of 100 base-pairs. On the other hand, the ability to encode the necessary error-correcting mechanisms requires significantly more than 100 base pairs, as the *Mycoplasma genitalium* genome demonstrates.

2.5.3 Proposed Solutions

Perhaps the most tantalizing thing about Eigen’s Paradox is that we know that there *is* a solution—life as we know it has already faced and dealt with the problem. This has given rise to a number of proposed solutions to the paradox, some of which are explored below.

2.5.3.1 Hypercycles

Eigen and Schuster (1977) proposed the “hypercycle” as a potential solution to the information integration problem. They hypothesised collectively autocatalytic sets of replicators, individually capable of maintaining their own integrity without complex enzymatic support (i.e., each < 100 base pairs), each of which could catalyse the replication of its neighbour, and relied on catalysis by another neighbour in the cycle to replicate. In this way, the total amount of information maintained in the cycle could increase linearly with the number of hypercycle members. However, the hypercycle model is vulnerable to parasitism—if one member of the cycle does not pass on the catalytic “favour” to the next member, then the cycle breaks down. Some form of spatial organisation, up to and including cellular containment has been proposed as a solution to the problem of parasitism and the hypercycle model itself has been the subject of significant further study since its inception (Boerlijst and Hogeweg, 1991; Hogeweg, 1994; Hogeweg and Takeuchi, 2003).

2.5.3.2 Stochastic Corrector Model

Szathmáry (1986) proposed a model called the Stochastic Corrector Model. This model again relies on rudimentary compartmentalisation and works on the preposition that there is an optimal composition for compartments, and that compartments which are further from the optimum will be less fit than those close to it. Compartments contain mixtures of altruistic and parasitic molecules. In the absence of either type, compartment growth cannot occur. Compartments divide by binary fission into two daughter cells at the point when they reach a certain pre-defined size. Upon division, internal molecules are assigned to either of the two daughter cells by random assortment. Other variations of this model propose that shallow tidal rock-pools could have acted as fixed-size containers, and as the tide periodically washed over them, their contents could be mixed and reassigned to different rock-pools (Maynard Smith and Szathmáry, 1997).

Szathmáry showed that under constant selectional pressure, at least some compartments in the population would be at the optimal composition. In some ways, the relationships between the molecules can be seen as a fully connected hypercycle, where everything can catalyse the replication of everything else. Parasitism is therefore embraced by the model rather than reacted to. Information may then be integrated within a compartment as a function of the types and concentrations of the molecules present. A compartment which has a rich diversity of different molecular species can qualitatively be said to store more information than a compartment which has a single molecular species.

2.6 Conclusion

Life arose on Earth in spite of harsh environmental conditions, seemingly intractable paradoxes, and chemical challenges that are still not fully understood. Every living thing around us shares a common DNA-based information chemistry. We do not, however, state that living things must be DNA-based as if it were some defining characteristic. Living things were defined as entities which can undergo an evolutionary process, or are the result of such a process. By looking further than the constraints imposed by chemistry and physics, we can identify key organisational challenges that would need to be overcome by any form of life, DNA-based or not. By focusing on these organisational, or systemic challenges, we may be able to improve our understanding of how and where evolution by natural selection works best, and furthermore, how we might harness that power to build complicated things from less complicated things. In the next chapter, we review research that has focused on these very challenges through exploring not biological but *artificial* life.

Chapter 3

Artificial Life

Exact details of the origin of life on Earth may never be discovered. The many different views of how life *may* have evolved on Earth presented in Section 2.1 highlight some of the current popular theories in the scientific community about the most probable candidates, but the debate over which, if any, of these pathways was followed is ongoing. This thesis does not aim to contribute to the origin of life debate: rather to assume that life, by whatever definition (Section 2.1.2), has originated at some point. From the point of origin onwards, evolution by natural selection has continued to produce novel changes to the living organisms by a process of *blind variation*, leading to a readily observable growth in the complexity of these organisms. Research in the field of artificial life is focused on exploring the potential pathways through which this growth of complexity may have occurred—“life as it could be” (Langton, 1989a). This chapter highlights some of the key literature in the field of artificial life, from the early work of von Neumann (1966) on constructive and evolutionary automata to current approaches to building artificial cells (McCaskill et al., 2008; Szostak et al., 2001; Glass et al., 2006; Gibson et al., 2008). Computer models of artificial life such as Tierra (Ray, 1991) are examined with respect to their ability to demonstrate an unambiguous growth of complexity in the digital “organisms” that they instantiate.

In Chapter 2, the key life sub-systems of containment, metabolism and an information carrier were discussed from the perspective of the origin of life. In this chapter, it will be argued that artificial life models typically focus on one or more of these sub-systems individually, and that this approach has so far yielded no satisfactory demonstration of the open-ended evolutionary growth of complexity. Any attempt to mimic living processes inside computers will necessarily require at least *some* abstractions from the kinds of system we can observe in the bio-sphere. It is argued here that since current artificial life models have failed to demonstrate a growth of complexity comparable to that observed in the evolution of the bio-

sphere, we need to revisit the abstraction process and investigate alternative ways to build our models. Later in this thesis, an artificial life model will be presented which explicitly deals with two interacting cellular sub-systems, though the details of each individual sub-system are significantly idealised. This approach is inspired by wet-lab efforts to build artificial “protocells”—a hypothesised transitional phase in the origin of life. The molecular entities which make up each of the sub-system components can be seen as individuals, and the entire collection can also be seen as an individual. Evolution by natural selection is acting upon all individuals, with the fitness of the collective being a function of some of the properties of the inner components. It is this “hierarchical selection” which forms the basis for the framework of Major Evolutionary Transitions proposed by Maynard Smith and Szathmáry (1997). Each of their hypothesised Major Evolutionary Transitions typically involve a reassessment of the definition of individuality in this sense, and a transition is usually deemed to have occurred when things that were once individuals in their own right can no longer exist independently of the collective. In the framework of the Major Transitions, the first transition, from individual replicating molecules to populations of replicating molecules inside cells is a good example of changing the scale of definition of the term “individual”, and is also clearly related to wet-lab efforts to create artificial protocells. In the post-transition cell, the internal replicating molecules are now dependent on the survival of the collective for their own survival. The most efficient way for such a collective to be managed would be under some form of centralised control, with a division of labour amongst the population of replicators. Evolution would therefore tend to favour populations which were more integrated than their peers. After repeated applications of this selective pressure, lineages of cells would emerge which contained highly coupled and integrated components, which truly depend on the other members of the collective to the point that they could no longer exist separately from the collective. If we were to examine each step along this evolutionary trajectory, we would find that the actions of every component were always “selfish”, regardless of their effect on the health of the collective. “Survival of the fittest” would ensure that fitter collectives will do better than weaker ones. By a process of division of labour then, selection will favour those collectives which have higher instances of cooperation, even though the actions of individuals in the inner population can still be described as purely selfish. For multi-cellular life, this evolution towards division of labour has happened on at least two levels: internally to each cell, and amongst cells making up the multi-cellular organism.

If Maynard Smith and Szathmáry are correct, it is this hierarchical selection, rather than the specific implementation details of the individual components, which holds the most promise for demonstrating the evolutionary growth of complexity in

artificial life models. The artificial life model presented later in this thesis provides a platform for the investigation of the potential of hierarchical selection to open up new evolutionary pathways as an incremental step towards, but stopping short of any attempt at demonstrating, full-blown evolutionary growth of complexity. In the following sections, I will present some of the most well-known artificial models in the context of their contribution to our understanding of the evolutionary growth of complexity, and highlight the important features of the biosphere that they have incorporated or omitted. These approaches will be contrasted with wet-lab efforts to build artificial protocells in the penultimate section.

3.1 Introduction

In this chapter, I will highlight the key milestones in the history of simulating life-like processes using computers. Artificial life is the study of the abstract organisational properties that constitute life. The basic tenet of artificial life is that life, as a process, is not necessarily constrained to operate within the bounds of terrestrial carbon-based chemistry. The question becomes one of exploring “life-as-it-*could*-be”. The umbrella of artificial life research has served to bring together researchers from many disciplines and focus them on the common goal of building living systems—*creatio ex materia*. Regardless of whether this field of research ever succeeds in producing entities which will be considered “alive” in their own right, taking the artificial life approach allows to investigate the *features* of real life in isolation. The particular feature of real life that is most relevant to this thesis is the unbounded evolutionary growth of complexity, and the potential for using artificial systems to explore this phenomenon is undeniable. In the following examination of the key artificial life models which have shaped the field, it will be argued that these models have failed to validate against the kinds of real life that are observable in the biosphere precisely because they have not demonstrated such an open ended evolutionary growth of complexity. This kind of qualitative validation is in contrast to a quantitative validation where the precise mechanisms used by life as we know it would be abstracted in the model, with later comparison of results obtained using the model with real world observations. In other words, rather than an attempt to model life, we are attempting to build an artificial system which might demonstrate qualitatively similar phenomena to life in the biosphere. By examining the successes and failures of current artificial life models with respect to their contribution to our understanding of the evolutionary growth of complexity and synthesising and comparing this to wet-lab approaches to building artificial cells, we hope to identify some potentially interesting areas for further exploration, and these will be applied

in the artificial life model that is presented later in the thesis.

3.2 Soft Artificial Life

A fine example of inter-disciplinary study in artificial life is the area of artificial chemistry. As we saw in Chapter 2, everything in our Universe is made up from combinations of elements from the periodic table, which arose from the “Big Bang”. In turn, these elements are all made up from more basic components and the related forces that hold them together. The Standard Model of Particle Physics is a particularly elegant example of a set of “simple” rules which is capable of describing the interactions of every particle in the visible Universe. Of course, being in possession of such a static list of rules and interaction descriptions is of little benefit to predicting what will happen when dynamic factors are taken into account. Chaos theory (Packard et al., 1980) is a field of study which is usually applied to such situations, and is particularly suited to the analysis of complex systems with many variables. Our purposes however are not to predict what *will* happen, but rather build an artificial system which has the potential to produce results which are qualitatively similar to our observations of life. The first step in building such an artificial system is to define the “chemistry” upon which it is based.

According to Dittrich et al. (2001), artificial chemistries consist of:

1. a list of simple rules which govern the interactions of certain elementary particles,
2. a (potentially infinite) set of valid elementary particles that may appear
3. a reaction algorithm, or chemical dynamics, which determines how the list of rules is applied to the set of elementary particles.

Furthermore, they argue that the use of artificial chemistries is the best way to study the early stages of the evolution of life as we know it. Holland (1998) spoke of “constrained generating procedures” which are essentially a super-set of artificial chemistries. According to Holland, “potentially *any* constrained generating procedure can exhibit emergent properties”, and the study of artificial chemistries has been focused on demonstrating and exploring such emergent phenomena.

In Chapter 2, the question of the physical nature of life was addressed, particularly the philosophical work of Bedau (1996; 1998). Bedau’s work questions whether life is substance-based or process-based. The latter theory is clearly favourable to the study of artificial life, especially in the instantiation of computer-based soft(ware) a-life—life in-*silico*. Since the birth of modern computers, there has been a constant

drive to use computing power to simulate complicated natural phenomena in an attempt to better understand them (Barricelli, 1957, 1963; Reed et al., 1967). This section offers a brief overview of the milestones in the history of the field of soft artificial life, or in the words of Bentley (2001): “Digital Biology”.

3.2.1 Von Neumann’s *Theory of Automata*

John von Neumann, a pioneer of computing, devoted a significant part of his research to the study of living-processes through what he christened “tessellation automata” (1966)—a concept that he developed without the aid of computers as a way of using pencil-and-paper to explore self-reproduction. Von Neumann hypothesised a self-reproducing system which employed a “tape” which stored the instructions required to build the self-reproducing automata, and a “universal constructor” which could read from the tape and produce whatever the tape described (Von Neumann and Burks, 1966)¹. Von Neumann recognised that if such a universal constructor could be designed and built, assuming suitable raw materials were available in abundance, a “tape” could be written to describe the universal constructor, thereby initiating a process of self-reproduction. Watson and Crick (1953), by discovering the structure of DNA, showed that life also uses the concept of a data-store, or “tape”, to encode the instructions needed to build self-reproducing entities (Section 2.3). Von Neumann also explicitly stated that the “tape” description should take the form of a list of building processes rather than a blueprint that described the direct connections between components. We know now that this is precisely the mechanism that the DNA based information chemistry of life adopts (Section 2.3). This view of self-reproduction is in contrast to the idea of reproduction by “self-inspection” which requires an entity to make an exact copy of itself by, presumably, disassembling itself to see how everything fits together and then rebuilding itself along with an identical copy.

Von Neumann also realised that if there should be any errors during the reproduction process, there might be serious consequences for the automata. For example, if there was an error during the production of the universal constructor, then that particular offspring may or may not be capable of acting as a universal constructor, which would potentially end that “lineage” of automata. If there was an error during the copy of the “tape” however, this error would be propagated throughout the lineage, assuming that the constructor described by the damaged tape was still capable of doing its job. In theory, such errors might then be subject to natural selection (Section 2.2.1), opening the doorway to a potentially unbounded evolu-

¹This work was posthumously completed by Arthur Burks, an early and frequent collaborator of von Neumann.

tionary growth of complexity for a lineage of such automata. Unfortunately in the case of von Neumann’s work, the parent and offspring tended to interfere with each other to the point that the entire system collapsed rather than produce anything more complicated.

3.2.1.1 General Cellular Automata

Von Neumann’s model was a specific example of a family of models known as Cellular Automata, which are based on a grid of cells. Cellular automata are useful mechanisms for abstracting “space” in artificial models of living systems. There is no limit to the dimensionality of the grid of cells, and each grid location is always in one of a finite set of states. Timesteps are discrete, and upon each timestep, the state of each grid location is updated based on rules which take into account both the current state of the cell, and the states of those cells in its neighbourhood. One such cellular automata model was that developed by John Conway (1970) which has become known simply as the “Game of Life”. It uses four simple rules:

1. Any live cell with fewer than two live neighbours dies—the *loneliness* rule.
2. Any live cell with more than three live neighbours dies—the *overcrowding* rule.
3. Any live cell with two or three live neighbours lives to the next generation—the *survival* rule.
4. Dead cells with exactly three live neighbours become live cells—the *birth* rule.

Conway initially hypothesised that it was impossible to create a pattern which could grow indefinitely—indeed, that was one of the “three desiderata” that Conway designed the system around in order to ensure unpredictability. Conway offered a prize of \$50 to the team who could prove this wrong within the year of publication, and this prize was claimed by a team led by Bill Gosper, who designed a “glider gun” which had the ability to periodically emit a moving pattern known as a glider. This discovery opened up the possibility of building a Turing Machine inside the “Life” universe, and such a machine was duly built (Rendell, 2002). The successful demonstration of a Turing Machine in “Life” proves that the “game” is theoretically capable of simulating the logic of *any* computer algorithm, but the difficulty in programming such a machine, or extracting meaningful results from it, not to mention its drastically reduced speed compared to the bare metal computer that it is running on has meant that as a computational device, this remains a toy model. More generally however, research using cellular automata has highlighted that spatial relationships can play an important role in producing interesting dynamical behaviours.

3.2.2 Genetic Algorithms, Classifier Systems & α -universes

Meanwhile, the field of computer science progressed as scientists explored how to write optimised computer programs which could exploit the power of computers for various tasks. John Holland developed the idea of genetic algorithms (Holland, 1975) which put evolution by natural selection forward as a mechanism for designing computer programs. The genetic algorithm approach has further been applied, for example, to design solutions to real world problems in the study of life (Theis et al., 2006).

3.2.2.1 α -universes

The α -universes presented by Holland (1976) were a theoretical presentation of a system capable of self-reproduction and self-maintenance. Later work by McMullin (1992) provided the first implementation of the system. This demonstration showed that in fact, the α -universes were more brittle than Holland originally assumed and that self-reproduction and self-maintenance were not so easily achieved. This was, for the most part, due to the presence of unanticipated side-reactions, though the model dynamics and lack of individual “containment” were also found to contribute to this brittleness.

3.2.2.2 Classifier Systems

Classifier systems were proposed by Holland as a further adaptation of the genetic algorithm. The goal of a classifier system is to classify elements of the environment in which it resides. Classifier systems consist of an environment which provides messages, or signals, and a set of rules which can react to some or all of these messages. These rules, or classifiers, are essentially “condition→action” rules—“if→then”. In the presence of environmental signals, each classifier in the system is checked to see if the signal can match the condition part of the classifier. If so, then this classifier is added to a working set of active classifiers. Once all classifiers have been compared to the environmental signal, an appropriate action is chosen by enumerating all the possible actions suggested by the classifiers in the working set, and choosing one based on various system parameters. Such a system can be seen as a type of artificial chemistry, though Holland himself did not explicitly refer to it in this way. The messages are the basic chemical molecules and the rules represent the different chemical reactions that can take place between given molecules. The experimental platform presented in later chapters of this thesis is inspired by this type of artificial chemistry with the major difference that rules and messages are no longer strictly demarcated from each other.

3.2.3 CoreWar, Tierra, Avida

As computer programming languages and techniques developed, and the availability of computing devices increased, the supply of available computing power began to outweigh the demand for it. The operators of these machines started to utilise these spare computing cycles to experiment with simulating living processes. As more computer operators became interested in this area, the increasing level of sophistication of their programs gave rise to direct competition between the computer operators to see who could dominate the “core” of the computer for the longest time. One of the earliest accounts of such recreational competition between digital organisms is that of McIlroy et al. (1972). There, they describe a game known as “Darwin”—“A Game of Survival and (Hopefully) Evolution” which they devised ten years earlier. Later, another flavour of this underground “sport” was presented by Dewdney (1984) as a game known as “CoreWar”. Modern day Internet worms and computer viruses are also inspired by these early games with the important exception that the earlier games were carried out inside specially controlled sandbox environments whereas Internet worms and computer viruses respect no such boundaries.

3.2.3.1 CoreWar

The CoreWar specification described the battlefield: the core-memory of a mainframe computer, the “warriors”: short computer programs written in a language known as RedCode, and the “rules of battle”: the MARS (Memory Array RedCode Simulator) which defined how the warrior programs were interpreted and generally enforced the rules of the “game”. CoreWar brought this concept of viral programs battling each other to an audience outside of the hardcore of System Operators who until now had been the only ones to play such games. CoreWar also standardised this game by specifying MARS and the RedCode language, which made experimenting with this kind of programming more accessible, since the MARS virtual machine could be ported to run on many different underlying architectures and ensured that these viral programs were not “running loose” in the core of expensive mainframe computers. Many different strategies evolved for building successful CoreWarriors, but evolution was not itself a component of the system. It would be six years before the concept of evolvability was applied to CoreWars by Rasmussen et al. (1990) when they created the CoreWorld system. The system was designed to explore the behaviour and possible spontaneous emergence of self-reproducing entities. This marked a significant departure from the CoreWars system which inspired this work, since CoreWars depended exclusively on hand-designed RedCode programs.

3.2.3.2 Tierra

In 1991, Tom Ray built a qualitatively similar system to CoreWar and called it “Tierra” (Ray, 1991). Tierra has huge significance for the field of artificial life because it had the specific goal of exploring the evolution of multi-cellular systems. “Organisms” are made up from sequences of simple computer instructions—similar to low-level assembly language—which are then processed by the model. This is a “Turing-complete” set of instructions, so Tierran organisms are capable of universal computation. The system essentially attempts to create a computer-based environment in which these digital organisms may evolve, presumably in the direction of increasing complexity. In contrast to the CoreWorld system, one of the key innovations of Tierra was the introduction of rudimentary memory protection as an approximation of cellular boundaries. This memory protection was variable, however, allowing creatures to “read” and use the code of other creatures, but preventing a creature from damaging or tampering with the code of another creature. Indeed, Ray was able to use Tierra to demonstrate the evolution of complicated eco-systems of replicators, parasites and even hyper-parasites—creatures which parasitised the parasites. In terms of the evolutionary growth of complexity however, Tierra falls short of the mark, typically reaching a complexity plateau which appears to be systemic to the architecture. CoreWar and Tierra inspired a family of models based on the same principal of using a computer programming language similar to assembly code to build the creatures. Avida by Adami and Brown (1994), Amoeba by Pargellis (2001) and Cosmos by Taylor (1999) are just some of the members of this “family”.

3.2.3.3 Avida

Adami and Brown (1994) presented the Avida system as an artificial life model which specifically targeted the evolution of complex *computational* behaviours, as opposed to producing a qualitative replica of life as we know it. Avida used a modified “Tierran” instruction set and required that the “creatures” actually carry out some pre-determined computation for survival. Creatures who could perform the computation better than the others would be more successful in the population. Avida short-circuited evolution in the sense that the optimum path for evolution to follow was pre-determined so that creatures which appeared to be somewhere along this path could be duly rewarded. Avida also included some modifications to the Tierran programming language, and introduced a 2D spatial structure to the environment. Tierra’s environment was linear, and creatures had the ability to interact with other creatures which existed within some parametrised distance from themselves. Avidan creatures existed on a grid, and had the ability to overwrite neighbouring cells by

reproduction, though unlike Tierra, they could not read or write to the memory of other creatures. Another key difference between Avida and Tierra is that in Tierra, there was essentially a single CPU, and a “slicer queue” determined which organism would “run” next, with the number of CPU cycles granted set as a parameter to the simulation. Each organism in Avida had its own CPU. Parallelism was simulated across all organisms, and the relative speeds of each organism’s CPU could be adjusted by a reward mechanism which externally judged how close each particular organism was to solving a pre-defined problem.

3.2.4 Algorithmic Chemistry

AlChemY, or Algorithmic Chemistry, was a system devised by Fontana and Buss (1992; 1994a; 1994b) as a direct attempt to build an “artificial chemistry”. Molecules in AlChemY are represented by expressions in the λ -calculus. The λ -calculus has a well defined structure and has the helpful (from the point of view of artificial chemistries) property that any λ expression can react with any other λ expression, and their product is the normalised form of their concatenation. Elastic collisions, collisions which do not result in a successful reaction, were introduced by placing a restriction on the normalisation process. The reason for this restriction was that the normalisation process was iterative and was not guaranteed to terminate. Furthermore, there is no general algorithm for discriminating expressions where normalisation would terminate from those where it would not. The restriction therefore was essentially a limit on the number of iterations of the normalisation loop. Fontana and Buss carried out sets of experiments to explore emergent organisations of molecules in AlChemY. It was found that in the presence of self-replication—the ability of one molecule to catalyse the production of a copy of itself—the reactor quickly became dominated by such self-replicators. By preventing self-replication, collectively autocatalytic sets of molecules were found to emerge, much like the Hypercycles hypothesised by Eigen and Schuster (1977) (Section 2.5.3.1). Dittrich and Speroni (2007), inspired by the work of Fontana and Buss, presented “Chemical Organisation Theory” as a generalised analysis for organisations and their components.

3.3 Wet-Lab Artificial Life

In the previous section, a brief history of artificial life was presented, specifically the computer-based life—life in-silico. In parallel with this research though, there has been significant work on the creation of artificial life *in vitro*. In spite of the problems highlighted in Section 2.5, wet-lab experimental work has continued to push ahead towards the goal of engineering a simplified artificial life form.

3.3.1 Protocells

Protocells are hypothesised as a transitional phase in the origin and evolution of life (Section 2.4.3). The earliest protocells would most likely not have had anything like the complicated scheme of interacting systems that exist in contemporary biological cells. In their most basic form, protocells are essentially vessels which can localise certain chemical reactions. This localisation is both in terms of reaction materials and reaction products (i.e., inputs and outputs). The vessels themselves may even have been externally provided, for example by the micro-porous surfaces of rocks at deep-sea hydrothermal vents (Section 2.4.3 and (Martin, 2003)). In these very early stages of protocell evolution, the interaction between the containment chemistry—the cell membrane—and the ongoing chemical reactions—RNA replication, redox reactions etc.—would have been minimal, and almost certainly nothing resembling a tightly coupled autocatalytic set like the Chemoton proposed by Gánti (2003). However, these rudimentary protocells can be seen as individual units of selection even though they are still not fully integrated.

Minimal protocells such as these can also be seen as basic living systems. Taking a reductionist approach and dismantling the protocell into its component systems will more than likely result in a collection of systems which by themselves could not be described as being alive. There is something special therefore about the organisation of collections of these systems, above and beyond the properties of the individual systems themselves. In their work on Major Evolutionary Transitions, Maynard Smith and Szathmáry (Section 2.4.2) explored some of the possible ways in which unconnected individual systems could combine into a higher level unit of selection and they point out key stages of the history of evolution which may in fact be due to such transitions. As they saw it, the early stages of each of these transitions was essentially a period of co-existence between a number of unconnected systems. In response to selectional pressure, they hypothesised that specific instances of these co-existing collectives may undergo more complicated interactions, and eventually evolve into a new combined system, which may result in the loss of “individuality” for some or all of the systems. Their key insight was to suggest that the complicated living systems that we see around us today may in fact be the result of a series of such “transitions” (Maynard Smith and Szathmáry, 1997).

Current research and experimental work aimed at creating artificial protocells *de novo* in the lab is heavily influenced by this “systems” view of life. The main idea is that any attempt to build a living entity must find a way to make the appropriate sub-systems integrate together in a reliable way. There are two main approaches to this: a top-down approach, which uses modern living cells as a starting point and strips away unnecessary components to approach a minimal set, and a bottom-up

approach, which tries to identify which systems are required for a minimal “living” set, and then combine them into something that could be described as a “living” unit by some definition.(Luisi, 2002; Szostak et al., 2001; McCaskill et al., 2008)

3.3.2 Bottom-Up Approach

The bottom-up approach to building life-forms in the lab is an attempt to test the various hypotheses for the origin of life (Sections 2.4.3.1, 2.4.3.2 & 2.4.3.3) in the hope of producing something which could be considered alive in the laboratory. The goal is to combine the various subsystems of metabolism, containment and an informational chemistry in just the right way that they form a self-sustaining, self-maintaining system.

Szostak et al. (2001) say that life is “a property that emerges from the union of two fundamentally different kinds of replicating systems: the informational genome and the three dimensional structure in which it resides”. The Chemoton proposed by Gánti (2003) is a theoretical example of such a tightly coupled union, though it is has not been implemented in “real” chemistry to date. The PACE (Programmable Artificial Cell Evolution) Project McCaskill et al. (2008) which funded the majority of the research supporting this thesis is a notable example of the bottom-up approach to building artificial protocells. These protocells would be created and managed inside a micro-fluidics system. Rasmussen et al. (2004b) proposed a simplified version of such a protocell, which heavily relied upon the natural self-organising properties of the different components of the system—the components were designed to self-assemble in water, and the system was eventually simulated using a model of dissipative particle dynamics.

3.3.3 Top-Down Approach

In recent years, considerable success has been achieved by those research teams who have taken the “top-down” approach to building artificial cells. This approach has also been described as “synthetic biology” by some researchers in the field (Szybalski and Skalka, 1978). Essentially, the approach involves identifying a “minimal genome” and then replacing the genome of a bacterial cell with the newly identified “minimal genome”. In 2008, a team led by Craig Venter, who until now has been widely recognised for his work with the Human Genome Project, succeeded in synthesising the genome of a minimal bacteria (Glass et al., 2006; Gibson et al., 2008), and have filed a patent claiming ownership of this “invention”. The minimal genome which the project began with was that of *Mycoplasma genitalium* and the resultant minimal set of genes which were identified to theoretically support life was dubbed

M. laboratorium (Reich, 2000).

3.4 Conclusion

Life, as we know it, is completely reliant on the DNA based physical chemistry that has evolved on Earth. However, there is nothing intrinsically special about DNA to suggest that it is the *only* mechanism that might support life. The field of artificial life—the study of life as it *could* be—emerged almost simultaneously with the birth of computers, and there is evidence to suggest that the earliest computing pioneers were interested in the idea of simulating living processes using their new machines. There have been many approaches to building artificial life forms in a digital medium, and in recent years, there have been numerous approaches to synthesising life in the laboratory. Digital artificial life platforms have the luxury of precise control over all aspects of the environment that the lifeforms will live in. In wet-lab chemistry however, some environmental parameters are less easily controlled. In the later experimental chapters of this thesis, an artificial life system, MCS, is presented as an idealised replicator world, where the majority of the environment is simulated as a black-box so that the key organisational properties of the system can be explored. This system is presented as a simple artificial chemistry and takes inspiration from Holland’s Classifier Systems.

Part II

Experimental System Design

Chapter 4

Building MCS

In this chapter, I describe the design and implementation details of the platform upon which the experiments supporting the discoveries reported in this thesis were carried out. The platform is called the Molecular Classifier System (MCS) and is essentially an agent-based artificial chemistry. MCS is an hierarchical artificial chemistry, which supports the simulation of primitive cells—protocells (Section 3.3.1). In the early sections of this chapter, I will present the basic assumptions which guided the design of MCS, including a synthesis of key research in the areas of origin of life and artificial life upon which this system is built. After that, the details of the two levels of the MCS hierarchy will be presented. Specifically:

1. the internal molecular chemistry of the protocells,
2. the coupling of this molecular chemistry to the growth and division of artificial protocells.

The chapter concludes with an attempt to produce a theoretical analysis of the system, though as we will see, a complete analysis of MCS becomes difficult due to the combinatorial effects of the increasing number of potential interactions between the molecular species in a reactor (protocell).

4.1 Methodology

In this section, I will clarify and justify the modelling methodology adopted throughout the remainder of this thesis. According to Bedau (1999), complex biological systems can be modelled computationally in at least two different ways: realistically and “unrealistically”. The first class of model focuses on “micro-mechanical realism”, while the latter class, common in the field of artificial life, are “intentionally [...] as *unrealistic* as possible”. At first glance, it might seem that Bedau is arguing that artificial life research is permitted to adopt an “anything goes” approach and that

such an approach might be excused from the burden of validation. On the contrary, Bedau points out that—

“... the micro-structure of such unrealistic models will not be *completely* unrealistic. Although vastly simplified, a model can explain the behavior of a real-world system only if the model’s micro-structure captures the *abstract form* of the target system’s micro-structure. The micro-structure of unrealistic models *is* realistic when viewed at an abstract enough level.”

In this chapter, the details of the MCS model will be described in the context of the biological literature reviewed in Chapter 2, though the MCS is not intended to be validated back against any particular target biological system in detail. Rather, the MCS is presented as an abstract, bio-inspired, virtual world which will be used to explore the effects of hierarchical selection in virtual worlds, as real objects in themselves, as opposed to “virtual worlds as ‘models’ of something else”. The “target system” in this case is the set of artificial life models reviewed in Chapter 3, and the MCS is intended to provide a platform for the preliminary investigation of the novel evolutionary dynamics which might arise in an artificial system which incorporates the phenomenon of hierarchical selection.

This “unrealistic” modelling methodology is quite different to the more traditional approach to building scientific models. Typically, the standard approach would be as follows:

1. choose some aspect of reality,
2. simulate this aspect in an abstract way,
3. validate the results obtained from the model simulation against those obtained from the real system,
4. adjust the model based on the previous comparison,
5. repeat steps 3 and 4 as required to improve the model.

In the traditional approach, the resulting model would then be used to make predictions about the real system. The accuracy of the model can then be tested by comparing such predictions to results obtained from the real system. The goal of realistic modelling then is to produce an abstraction of reality which can be used to make valid predictions about the real world. In contrast, the “unrealistic” approach adopted throughout this thesis takes *inspiration* from a real world system. The “model”, then, becomes a subject of enquiry in its own right, rather than informing our understanding of the system which originally inspired it.

4.2 Underlying Theory

Before describing the technical details of the implementation of MCS, I will briefly highlight the key concepts underlying the system. In Part I of this thesis, it was argued that life is the result of a series of complicated interactions between systems at various hierarchical levels. These living systems are constrained to be built from elements from the periodic table, and bound by the rules of physics.

The purpose of designing a system such as MCS then is to enable the simulation of *abstract* chemical systems. In particular, MCS is designed to support an abstract model of a system of replicating enzymatic templates. Biologically, these templates might equate to, for example, RNA or any other chemical compound capable of replication by a mechanism of template matching. In other words, the molecules in MCS will function both as templates to be replicated and as catalysts that can implement such replication. The benefit of a system such as MCS is that it allows a researcher to explore the organisational properties of these kinds of replicator systems with relative ease. Conversely, the wet-lab approach relies on the initial solution of many chemical problems before the organisational properties of the system can be analysed (Section 2.5). Computer simulations also give the researcher complete control over the reaction environment, whereas the variable factors that might affect a wet-lab chemistry system are physical properties of the environment, for example:

1. pH,
2. temperature,
3. pressure.

By abstracting away from these physico-chemical restrictions, systems such as MCS can focus on the higher level organisational, and potentially life-like, interactions that might occur.

4.2.1 Replicator Theory

Since the purpose of MCS is to simulate an abstract replicator world, it will be important to identify what exactly we mean when we speak of such systems. A replicator world is a very minimal model of a pre-biotic system. The operational entities are replicators—they are potentially “copiable” by for example, an autocatalytic process. Szathmáry and Maynard Smith (1997) proposed a classification mechanism for replicators in which they suggested that there are two core properties of replicators upon which life depends:

- unlimited heredity—replicators can exist in an indefinitely large number of forms, and
- modular rather than processive replication—information can be stored digitally.

Replicators in the MCS abstraction will be designed to meet these criteria.

4.2.1.1 Replicator Worlds

As we have seen, a “replicator world” can be seen as a plausible candidate for something resembling a pre-biotic world that might have led to the origin of life as we know it. Essentially, a “replicator world” could be described as a precursor to a “minimal living system”—one where the basic interaction dynamics and “living conditions” can be easily studied and modified. At the same time, it is also apparent that building such simple systems allows us to be extremely focused on a small number of variables, at least when compared to wet-lab attempts at building minimal living systems such as those presented by Rasmussen et al. (2004a); Luisi (2002); Hanczyc and Szostak (2004).

4.2.2 Catalysed Replication

Natural template replicator world systems assume the presence of at least some molecular species which are capable of catalysing or enabling the replication of other molecular species in the system—*replicases*. If such a species can catalyse the replication of another member of its own species, then we could call that species a self-replicase. In MCS, we make the further assumption that *every* molecule is capable of acting as a replicase—under certain conditions, they can cause the replication of other molecules. We also take the step of removing all intrinsic fitness differences between replicases in the system—all replicases are equally good at being replicases. Using this *ideal* system, we can therefore study the purely organisational properties of the interacting systems, since we have removed all other aspects of the system where Natural Selection may act.

4.3 Idealisation

The abstractions applied in the design of the MCS can be described as part of a process of “idealisation”. In much the same way as the ideal gas model is an important part of statistical chemistry, *idealised* systems such as MCS might prove to have an important role to play in the understanding of complex dynamical systems, of which

“life as we know it” is a good example. The ideal gas model allows the approximation of various properties of gaseous mixtures which prove to be qualitatively close to the true measurements. The process of idealisation identifies the core concepts that govern the system, and makes assumptions or approximations for some of the more complicated parameters. In the case of the ideal gas model, particles are assumed to interact with each other only through elastic collisions, and to only occasionally interact with their surroundings. In many cases, this is a good approximation of how real gaseous particles behave, and it is for this reason that predictions made based on the ideal gas model correspond well with the real life situation.

The MCS as presented in this thesis is the result of a similar process of idealisation—an idealised enzymatic replicator world. In particular, replication in MCS requires mediation by a replicase which also needs to be replicated. This trans-acting nature differs from other idealised replicator worlds that have been presented, for example, the RNA replicators of Eigen and Schuster (1977), which involved catalysis by an enzyme which itself was not replicated. As such, the MCS replication dynamics is that of hyperbolic growth (Szathmáry and Maynard Smith, 1997). As discussed earlier in the thesis, “replicators first” theories for the origin of life, such as the RNA World Hypothesis, have received considerable attention from the scientific community, though there is increasing support for a combined “replicators, metabolism and containment” origin scenario (Section 2.4.2). As a concrete example of a “replicators first” theory, the RNA World hypothesis places specific requirements on the physico-chemical properties of the world which would facilitate the origin of life:

1. the presence of enough nucleotides to support a rich RNA chemistry,
2. the existence of RNA-based replicases which enzymatically support the replication of some classes of RNA molecules, allowing closed, self-sustaining, replication networks.¹,
3. and of course, suitable reaction conditions (temperature, pressure, pH etc).

The MCS is presented in this thesis as a basic replicator world and is a good first step in the “idealisation” process. This first step involves abstracting away from the precise physico-chemical nature of the replicator world, as suggested in Section 4.2. This involves changing the particular “family” of polymers that the replicators are composed of. “RNA” can be said to denote the family of all polymers based on the set of nucleotide monomers [“Adenine”, “Cytosine”, “Guanine”, “Uracil”]; “DNA” denotes

¹While an RNA-*self*-replicase remains to be found in laboratory explorations, limited examples of RNA catalysed replication have been presented, including a demonstration of a cross-catalytic network of two RNA enzymes. (Johnston et al. (2001); Voytek and Joyce (2007); Lincoln and Joyce (2009))

the family based on a slightly different set of nucleotides [“Adenine”, “Cytosine”, “Guanine”, “Thymine”] and proteins are yet another family of polymers based on the set of 22 amino acids that are coded for in the standard genetic code. Replicators in the MCS system will come from a family of artificial polymers simulated in an agent-based computer program. A possible next step in the idealisation process would be to assume constant and appropriate conditions of temperature, pressure, pH and other environmental variables. We are left with a system which assumes that there is a population of polymer molecules drawn from an abstract “family”, some of which have the property of being able to catalyse the replication of some more-or-less specific polymers in the same “family” (possibly even catalysing other instances of themselves—“self”-replicases).

Perhaps the most striking idealisation in the system that is presented in this chapter is that it assumes not only the existence of simple replicase molecules but an abundance of such replicases to the extent that *every* molecule in the system has that capability. If the system cannot be made to work in this “best case scenario” where replication is relatively pervasive, then the probability of success using a more complicated system where replication is significantly more difficult must presumably be negligible. As Dawkins says in his 2009 book, “The Greatest Show on Earth” (p. 230): “If they could get realistic results even after throwing away some of the known properties of a cell, it would presumably be possible to get at least as good results with a more complicated model that left those properties in”.

4.4 System Design

MCS has been designed as a highly idealised artificial chemistry (Dittrich et al., 2001) loosely based on John Holland’s Learning Classifier Systems [(Holland, 2006; Holland and Reitman, 1977); Section 3.2.2.2]. The operation of the system depends on a population of “molecules”, which take the form of binary strings. Each molecule has an “informational” representation (primary structure, or monomer sequence) and an enzymatic representation (“folded” or secondary structure, or “shape”), as inspired by the ribozymes of the RNA world hypothesis (Joyce, 1991). The model also contains a rule-set which determines the potential for interaction upon collision of a particular pair of molecules. It was decided that there would be no intrinsic fitness differences between individual molecules. The core motivation for this was the fact that the MCS is intended as an investigation into the effects of hierarchical selection. By forcing molecular replication to happen at the same rate, regardless of molecular sequence, we hope to isolate and explore the phenomenon of hierarchical selection. Artificial protocells are then crudely modelled as containers of a dynamic

mix of these molecules, which continuously interact and exert enzymatic actions on each other. This “informational chemistry” ultimately may be evolved to realise some particular computation—provided that it is simultaneously capable of sustaining its own dynamic organisation. In particular, this informational or computational subsystem must grow (in absolute number of molecules) and divide in coordination with overall cell reproduction.

4.4.1 Modelling Platform

MCS is an agent-based, or individual-based model. This means that every entity inside an MCS simulation is modelled as an entity in its own right rather than being simulated as, for example, a numerical count of molecules which share similar properties. For example, rather than analytically simulate the expected results of 1000 interacting molecules, we instantiate 1000 molecules, allow them to interact for the required number of timesteps, and note the results. We can then re-seed the random number generator and run the simulation again to produce a different dataset.

4.4.2 MCS-“world”

The MCS world consists of one or more discrete “reactors”, each of which contains a “well-stirred” mixture of the artificial molecules already mentioned.

4.4.2.1 Reactor Instantiation

The individual reactors in the MCS world can be instantiated in two distinct ways:

1. a continuously stirred tank reactor (Stadler et al., 1993)—implemented with a variable dilution flow exactly matched to the production rate (once the maximum capacity is reached) The “variable dilution flow” is implemented as follows:
 - (a) the maximum number of molecules it may contain is decided.
 - (b) If this threshold is exceeded, a random molecule is removed from the reactor.
 - (c) To prevent any initial transient dynamics, all single-reactor runs will be initialised with the maximum number of molecules (normalised to an appropriate mutational load according to a given mutation rate).
2. a protocell reactor—protocells do not have a fixed molecular capacity, but rather a threshold which triggers protocell division by binary fission. With

the exception of the experiments in Appendix B, this threshold will be based simply on the total number of molecules inside the protocell². It is also worth noting that one *could* consider the case of a single such protocell. Whenever it “divides” we discard half the contents. Considering only molecular concentrations (which factors out the varying absolute number of molecules), the single protocell reactor and the single tank reactor would be equivalent. Of course, when we consider a population of reactors, the two cases would be quite different at the population level, but individual protocells can still be seen as tank reactors in their own right.

4.4.2.2 World Instantiation

Given the above reactor types, the MCS world can be instantiated in two different ways, one which consists of a single tank reactor, and one which consists of a population of protocell reactors. The latter can be further sub-divided into two cases depending on whether a maximum limit is imposed on the size of the protocell population or on the underlying molecular population.

A fixed size *protocell* population can be maintained as follows:

- Upon every protocell division, a random protocell is chosen from the population and removed to accommodate one of the new cells, with the other daughter cell replacing the parent, thereby keeping the total number of protocells fixed.

A fixed global *molecular* capacity can be maintained as follows:

- Upon every molecular replication event, if the maximum global molecular capacity is exceeded, a random *protocell* is chosen from the population of protocells and removed, along with its contents. This has the effect of causing the total number of molecules in the system to fluctuate by an amount roughly equal to the average size of a protocell in the system.

4.4.2.3 Tank vs. Protocell Reactors: Role of Variable Reactor Size

In later chapters, both fixed-size tank reactors and protocell reactors which vary in size as measured by the number of contained molecules will be considered. Analyses for each of these cases will be presented at the appropriate time. At this point, however, it is worth noting that it is possible, in some instances, to equate the analysis of both fixed-size reactors and variable-sized reactors. Under the following three assumptions, results from tank reactor experiments are comparable to results from protocell reactor experiments:

²An additional abstract molecular component, termed *protocell fission factor*, controls protocell fission in the experiments of Appendix B, the precise details of which will be presented there.

Equivalence of Reaction Rates In order to keep reaction rates comparable between tank reactor experiments and protocell reactor experiments, time will need to be normalised to the total number of individual molecules in the system. For a tank reactor, this will be the sum of all molecules in the tank, while for protocell reactors, it will be the instantaneous count of the number of molecules across the protocell population. This normalisation process results in a measure of time such that each molecule in the system is involved, as a substrate, in a reaction approximately once per unit time. We can use this measurement of time to achieve comparability between reaction rates measured for tank reactors and for protocell reactors.

Equivalence of Reaction Kinetics Reactors in MCS, tanks or protocells, are implemented to work as a well-stirred mixture. This means that every molecule in reactor has equal probability of acting as a catalyst for any interaction that occurs in that reactor. Applying mass action chemical kinetics, reaction rates will be proportional to the product of the *concentrations* of any two reactive species. Thus, for any given species, the rate of change in its concentration *due to reactions* will be a function only of the *concentrations* of all the species present regardless of the reactor size (as measured by the absolute number of molecules present).

Equivalence of Molecular Concentrations When dealing in terms of molecular concentrations rather than absolute numbers of molecules, using approximate continuous (ODE) methods, then the molecular dynamics in a fixed-size tank reactor will generally be precisely the same as the variable size protocell analysis. To illustrate, consider the case of a fixed-size tank reactor which contains two species, X and Y with concentrations x and y respectively. Taking unity as the total molecular concentration of the reactor, we have:

$$x + y = 1$$

In the case of a growing protocell reactor on the other hand, whose size in terms of absolute number of molecules is variable, we can measure the individual concentrations x and y by scaling the instantaneous absolute number of each species by the instantaneous total number of molecules as follows:

$$x = \frac{X}{X + Y}$$

$$y = \frac{Y}{X + Y}$$

where X and Y are the instantaneous numbers of molecules in species X and species Y respectively. Adding these concentrations together sums to unity as would be expected:

$$\begin{aligned} \left(\frac{X}{X+Y}\right) + \left(\frac{Y}{X+Y}\right) &= \left(\frac{X+Y}{X+Y}\right) \\ x + y &= 1 \end{aligned}$$

This equivalence of concentration will be maintained in the ODE analysis through the use of a dilution term. In the case of the fixed-size tank reactor, dilution is due to outflow which maintains a fixed maximum number of molecules. In variable-sized protocell reactors, dilution equates to a re-scaling of concentrations as the absolute number of molecules in the reactor increases.

So, it is possible to maintain equivalency between fixed-size analysis and variable-size analysis if we constrain ourselves by:

- scaling time to the total number of molecules in the system,
- applying mass action kinetics so that reaction rates are proportional to individual concentrations,
- applying a dilution term to the absolute numbers of molecules so that concentrations are equivalent.

4.4.3 Molecular Sub-System

Molecules in MCS are strings of symbols taken from the alphabet $[0,1]$. The molecules are of arbitrary length in general, but for the purposes of the work presented here, molecules are typically less than 30 monomers long. This basic structure, also known as the molecules “primary” structure, forms what was referred to earlier as the “informational” aspect of the molecule.

In Chapters 5 and 6, the only supported enzymatic function of molecules is to make an error-prone bit-wise copy of the primary, informational, structure of the bound, substrate, molecule; that is, a replicase function. More specifically, if a particular molecule has the ability to bind to molecules with the same molecular structure as itself, it will effectively be able to function as a self-replicase. The precise details of the folding and binding mechanism will be reserved for later discussion, as this varies across the experimental sets, though at this point, it should be noted that the probability of a molecule binding to another molecule in any given interaction or collision is always either 0 or 1—it is an “all or nothing” match. In Appendix B,

the enzymatic functionality of MCS molecules will be extended to support a basic computational function.

4.4.4 Protocell Sub-System

As described in Section 2.4.2, a set of major evolutionary transitions are hypothesised to have been critical for the evolution of complexity for life as we know it. Such transitions rely heavily on the concept of hierarchical selection which is intrinsic to the chemical organisation of life as we know it, and was also described in Section 2.4.2. Hierarchical selection is incorporated in MCS by constraining the molecular artificial chemistry outlined above to discrete, variable-sized reactors, referred to here as protocells, and inspired by wet-lab efforts to fabricate artificial cells *in-vitro*.

4.4.4.1 Description of Membrane Chemistry

In nature, cellular membranes are made up from oily lipid molecules called amphiphiles. These are long molecules which have two uniquely identifiable ends. One end is hydrophilic, the other is hydrophobic—one end is attracted to water, the other end is repelled by it. Given the correct temperature, pH level and other chemical conditions, these molecules can spontaneously self-organise into a variety of structures—including micelles, bilayers and vesicles. Such structures could form the basis for the containment of a replicator chemistry such as that proposed in the RNA World hypothesis. In MCS, we have taken inspiration from this lipid membrane chemistry for building our artificial protocells, though we do not explicitly model individual lipid molecules—we assume that there are always “enough” lipids in the system, and that their dynamics is such that protocells can grow and divide as necessary.

4.4.4.2 Protocell Growth

As reactions take place between the informational molecules in the MCS, new informational molecules are produced. In theory at least, this should cause the containing protocell to increase in size slightly. While there are various physical and chemically realistic methods to achieve this behaviour—for example, one could model a protocell as a sphere, and use the well-known formulae for spherical volume and surface area to calculate the various pressures that would act on the protocell—we will simply assume that a side-effect of a successful replication event is the catalysis of some abundantly available raw material to produce just enough new membrane material at just the right rate.

4.4.4.3 Protocell Division

As described in Sections 4.4.4 and 4.4.4.2, the protocells in the MCS continuously grow as replication events proceed. Protocells reach maturity and undergo binary fission when the amount of “fission factor”, F , they contain reaches a parametrised threshold, F_{max} . The informational molecules that were contained within the parent cell are each randomly assigned to one or other of the daughter cells and the “fission factor” which was contained in the parent is divided evenly amongst the two daughter cells. On average then, each of the daughter cells will be half the size of the parent, and will contain half of the “fission factor” of the parent, $\frac{F_{max}}{2}$. In the experiments presented in Chapters 5 and 6, fission factor will be produced as a side product of every successful replication, so that, in effect, a protocell will divide when the number of informational molecules reaches a fixed threshold, which is common to all protocell strains. The system is modified in Appendix B such that only certain replication reactions will be accompanied by the generation of fission factor, so that there may be variation among protocells in the sizes at which they fission, and this variation may be, to some degree, heritable.

4.4.5 Reaction Mechanism

Having presented the Molecular sub-system and Containment sub-system in detail, we examine below the mechanism for selecting reaction candidates on any given MCS timestep. An MCS timestep consists of:

1. selecting two molecules,
2. assessing whether these molecules can “bind” to one another (see Figure 4.1),
3. if necessary, splitting the protocell into two daughter protocells (see Figure 4.2),
4. handling any resource limiting that the “world” requires (Section 4.4.2).

4.4.5.1 Picking Reaction Materials

Since we are essentially trying to model parallel dynamics on a serial computer, an important aspect of the MCS simulation will be the mechanism that is used to choose which pair of molecules will be used for collision events at each timestep. There are various factors to consider when designing a fair algorithm to do this. Firstly, and most obviously, the two chosen molecules must be from the same reactor/protocell, otherwise they are physically not capable of reacting, regardless of any

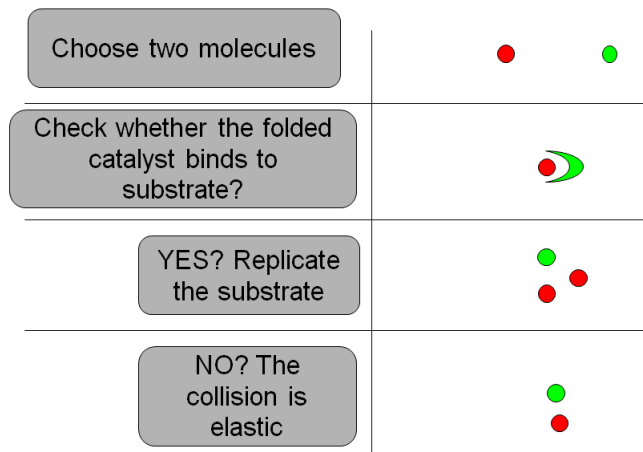


Figure 4.1: MCS Reaction Algorithm Flowchart

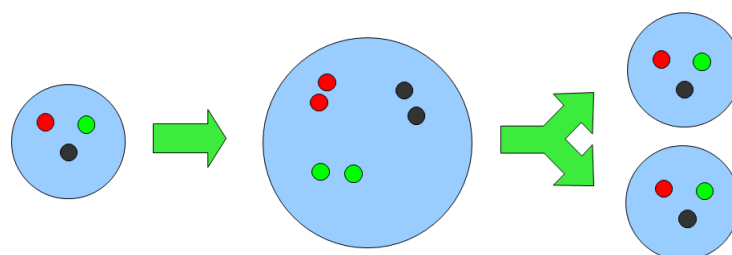


Figure 4.2: MCS Protocell Growth

other factor. Secondly, we must ensure that all molecules have an equal probability of being included in the reaction (notwithstanding the fact that they must share a common reactor/protocell), so we cannot, for instance, choose a random protocell, and two random molecules from within it, since that would imply that individual molecules which are contained within smaller protocells will have a higher probability of reaction than individual molecules contained within larger protocells. The implemented algorithm operates as follows:

1. Pick a molecule uniformly at random from all possible molecules (regardless of protocell boundaries), then
2. Pick another molecule from the same protocell as the first,
3. Test for a potential reaction, otherwise repeat step 1.

4.4.5.2 Error-prone Molecular Replication

Assuming a pair of compatible molecules has been selected, the simulation proceeds by making a bit-wise copy of the substrate molecule. It is assumed that *all* catalyst molecules which can bind successfully to the substrate are capable of instantaneously replicating that substrate—there are no differences in reaction rates among those reactions that take place. Replication rate is neither affected by the structure nor length of the substrate molecule and the catalyst molecule.

The presence of variation during the replication process is crucial to a process of evolutionary selection in MCS, so the process is made error-prone by the application of a per-bit (per-monomer) error rate. In other words, each bit in the substrate bit string, when copied, is susceptible to mutation at a particular rate. Mutation can be one of three types:

1. Bit-Flip: the bit is copied, in error, by the transformation $\{0, 1\} \rightarrow \{1, 0\}$.
2. Bit-Deletion: the bit is not copied $\{0, 1\} \rightarrow \{\emptyset\}$.
3. Bit-Insertion: a random bit, $\{0, 1\}$ is inserted before this bit is copied.

Once the entire substrate bit string has been processed, the transient binding between the catalyst, substrate and product is broken, and the newly produced molecule is added to the reactor (tank reactor or protocell). The various constraints on molecular population size described in Section 4.4.2 are then applied at this stage as appropriate.

4.5 Theoretical Analysis of Molecular Chemistry

Given the reaction rules laid out above, we can attempt to formulate a general theoretical analysis which can serve as a tool for exploring the data that are produced from the various experimental sets which follow in the later chapters. This analysis was first published in McMullin et al. (2008) in the final report on the EU FP6 funded PACE project (FP6-IST-FET-002035; McCaskill et al. (2008)).

4.5.1 General Dynamic Equation

In Section 4.4.3, we saw that molecular species i binds to molecular species j with a probability of either 0 or 1. Let c_{ij} then denote the rate with which species i , acting as replicase, will replicate species j . Given our assumption of all-or-nothing binding, and the further assumption of equal replication rates for all cases where binding takes place, we can say, without loss of generality, that every c_{ij} (“replication co-efficient”) will have a normalised value of either 0 (no binding) or 1 (binding). Under a continuous approximation, we will denote the concentration of any species i in the reactor by $x_i \in [0, 1]$ (with $\sum_i x_i = 1$). Applying mass-action kinetics, the replication rate of each species will be:

$$r(x_i) = \sum_j x_i x_j c_{ji}$$

and the total replication rate for the reactor will be:

$$R = \sum_k r(x_k) = \sum_k \sum_j x_k x_j c_{jk}$$

Under the condition that the total of all species concentrations is fixed at one ($\sum_i x_i = 1$), this total production rate must equal the total dilution rate, and the dilution rate for each individual species can be written:

$$\begin{aligned} d(x_i) &= x_i R \\ &= x_i \sum_k \sum_j x_k x_j c_{jk} \end{aligned}$$

Neglecting mutation (in the first instance) the dynamics of the reactor is then governed by the following system of ordinary differential equations (ODE) in each x_i :

$$\begin{aligned}
\dot{x}_i &= r(x_i) - d(x_i) \\
&= x_i \left(\sum_j x_j c_{ji} \right) - x_i \left(\sum_k \sum_j x_k x_j c_{jk} \right)
\end{aligned}$$

Under the conditions specified, we are therefore dealing with various cases of a catalytic reaction network, in the sense of Stadler et al. (1993). To actually examine the concrete dynamics of any given reactor, we must specify the total number of species, n , and the corresponding $n \times n$ matrix of replication co-efficients c_{ij} . Thus, even with the strong simplification of $c_{ij} \in \{0, 1\}$, for any given n there are 2^{n^2} distinct possible systems. Exhaustive analysis rapidly becomes infeasible. We shall limit our detailed analysis to considering all possible cases for $n \leq 2$. This provides a repertoire of “core” behaviours. We shall augment this analysis with more qualitative discussions of how these core behaviours may be generalised or combined for systems with $n > 2$ and/or perturbed by replication error (molecular level mutation).

4.5.2 Terminology

Although we are discussing molecular level evolution, it will be convenient to use the following, ecologically based, terminology. Consider two distinct molecular species i, j where $i \neq j$:

- If $c_{ii} = 1$ we say this is a self-replicase; otherwise it is *self-inert*.
- If both $c_{ij} = 1$ and $c_{ji} = 1$ we will say that i and j are *mutualists* relative to each other.
- Where $c_{ij} = 1$ but $c_{ji} = 0$ we will say that, relative to each other, i is a *host* and j a *parasite*.
- A mutualist, host or parasite may be said to be *facultative* if it is also a self-replicase; otherwise it is *obligate*.
- If $c_{ji} = 0$ we may also say that i is *inert* to replication by j (and vice versa); if both $c_{ji} = 0$ and $c_{ij} = 0$ we say they are *mutually inert*.

4.5.3 “Self”-Systems

For “self”-systems we have $n = 1$; that is, there is just one species with concentration x_1 , one co-efficient, c_{11} , and two possible systems:

- $c_{11} = 0$: The species is self-inert. Given that this is the only species present, then there are no reactions, and the system as a whole is inert ($\dot{x}_1 = 0$, $R = 0$).
- $c_{11} = 1$: The species is a self-replicase. Given that this is the only species present, then although there is constant turnover at the maximum rate ($R = 1$), in the absence of mutation this produces identical replacement molecules meaning $\dot{x}_1 = 0$. (The concept of mutation necessarily requires $n > 1$, so it is not meaningful in systems restricted to $n = 1$.)

4.5.4 Binary Replicase Interaction Systems

For pairwise systems we have $n = 2$; that is, there are two species with concentrations x_1 and $x_2 = 1 - x_1$; and their interactions are represented by the matrix:

$$\begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix}$$

As $c_{ij} \in \{0, 1\}$, there are $2^4 = 16$ possible such pairwise interaction matrices. However, some of these are equivalent under a relabelling of x_1 and x_2 ; as a result there are just 10 dynamically distinct classes of pairwise system. All of these will be considered in turn. In each case, we shall initially neglect mutation, so that the reactor operates under the constraint $x_1 + x_2 = 1$, and the (approximate) dynamics is fully characterised by the single differential equation:

$$\begin{aligned} \dot{x}_1 &= x_1^2 c_{11} + x_1 x_2 c_{21} - x_1 (x_1^2 c_{11} + x_1 x_2 (c_{12} + c_{21}) + x_2^2 c_{22}) \\ &= x_1^2 c_{11} + x_1 (1 - x_1) c_{21} - x_1 (x_1^2 c_{11} + x_1 (1 - x_1) (c_{12} + c_{21}) + (1 - x_1)^2 c_{22}) \end{aligned}$$

We shall also be interested in the overall replication rate, given by:

$$R = x_1^2 c_{11} + x_1 (1 - x_1) (c_{12} + c_{21}) + (1 - x_1)^2 c_{22}$$

For each distinct class we shall present the relevant pairwise interaction matrix or pair of equivalent matrices, the resulting simplified expression for \dot{x}_1 , optionally a graph of \dot{x}_1 , and a brief discussion of the resulting dynamic behaviour including example trajectories as appropriate. Finally we shall return to re-consider the general case ($n > 2$) as potentially approximated by a simple superposition of pairwise systems ($n = 2$), all instantiated in the same reactor simultaneously.

4.5.4.1 Analysis of Systems

Class-0:

$$\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

$$\dot{x}_1 = \dot{x}_2 = 0$$

Here the two molecular species are inert. The growth rate for both species is therefore, trivially, 0. The molecular concentrations of each species will therefore remain constant at the values at which they were initialised, and global reaction rate, R , will be 0. In fact this particular result holds for arbitrary numbers of molecular species (assuming they are all individually and mutually inert)

Class-1:

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$\dot{x}_1 = 0$$

The two molecular species are (facultative) mutualists. In other words, there is full cross-catalysis between the molecules. Since neither molecule has a distinct advantage over the other however, the growth rate for both species is again 0. In contrast to class-0, there is continuing turnover of the reactor contents at the maximum possible rate, 1.0. Accordingly, this would predict, under stochastic conditions, that there would be slow random drift in relative concentrations, with eventual takeover by one or the other. However, because the intrinsic dynamic here is essentially neutral, then mutation, even at a low level, would significantly modify this behaviour, effectively introducing negative frequency-dependent selection. In the simplest case, with just two possible species, each having an equal rate of mutation, this would actively stabilise the state $x_1 = x_2 = 0.5$ against drift. There would still be stochastic fluctuation around this state, with the level of fluctuation inversely related to the mutation rate (and, of course, the population size). With a somewhat larger number of mutationally related species, still all being pairwise facultative mutualists, this result may generalise to a stable equilibrium distribution across all species in the set (though with continuing replication/turnover at the maximum rate). Such a family of species would be somewhat analogous to the concept of a quasi-species in (externally catalysed) replicator systems (Eigen et al., 1989).

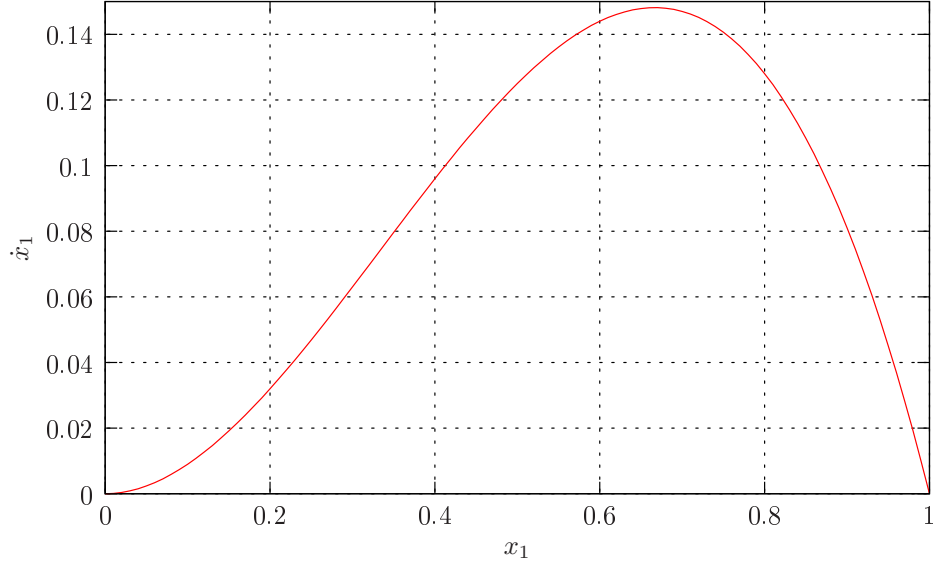


Figure 4.3: Class-2 Differential Equation

Class-2:

$$\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$

$$\dot{x}_1 = x_1^2(1 - x_1)$$

One species is a self-replicase and the other is completely inert. As seen in Figure 4.4, the self-replicase species can invade and displace a population of the inert species, even from an arbitrarily low initial concentration. Also evident in Figure 4.4 is a sharp “switching” dynamic in global reaction rate, R , as the self-replicase increases in concentration, with R eventually reaching the maximum possible value of 1.0.

Class-3:

$$\begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$$

$$\dot{x}_1 = -x_1^2(1 - x_1)$$

Neither species is capable of *self*-replication, though one species can act as a replicase for the other—the latter is an obligate parasite of the former. The ODE here shows that x_1 , the obligate *host*, is essentially a finite resource which will be irreversibly consumed by the obligate parasite, even from an arbitrarily low initial concentration of the parasite. As the concentration of the host species decreases, so too does overall reaction rate as seen in Figure 4.6. In the limit, when the final host molecule is eventually displaced, there will be no further reactions, at which

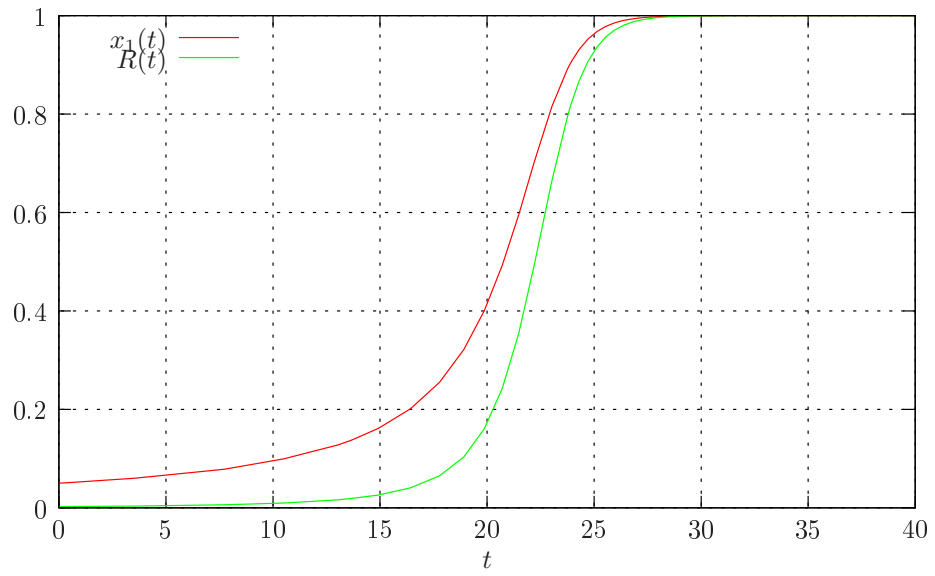


Figure 4.4: Class-2 Trajectory: Growth of self-replicase [$x_1(0) = 0.05$]

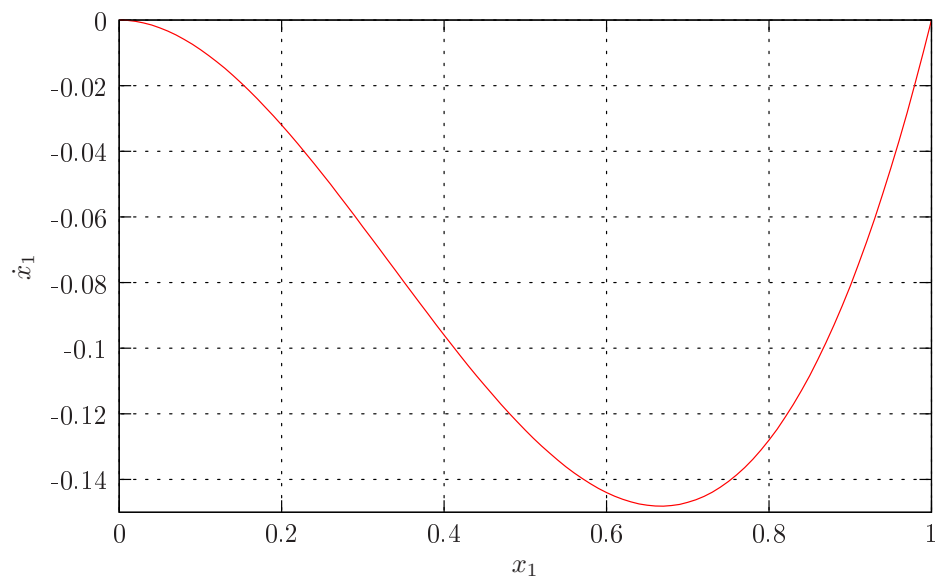


Figure 4.5: Class-3 Differential Equation

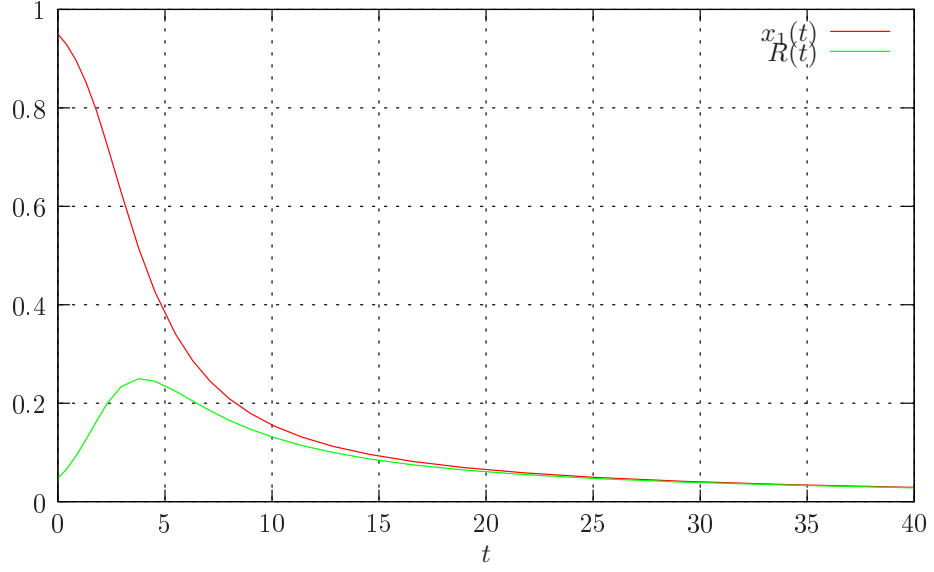


Figure 4.6: Class-3 Trajectory: Decay of Obligate Host [$x_1(0) = 0.95$]

point the reactor formally contains a self-system ($n = 1$) which is self-inert (Section 4.5.3).

Class-4:

$$\begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix}$$

$$\dot{x}_1 = 0$$

One species is a self-replicase and the other is an obligate parasite. This class-4 dynamic lies somewhere between class-2 and class-3 in terms of the predicted outcome for the host species. Since at least some of the interactions will involve parasitism by the obligate parasite, things are not as favourable for the host as in the case of class-2 dynamics. On the other hand, things are not as bad for the host as in the case of class-3 dynamics, since at least some of the interactions are spent by the host on replicating itself. In the ODE approximation, these effects exactly balance each other out, $\dot{x} = 0$, for all concentrations of x_1 , though the overall reaction rate, R , will indeed vary directly with x_1 , as it is x_1 which drives replication in the system.

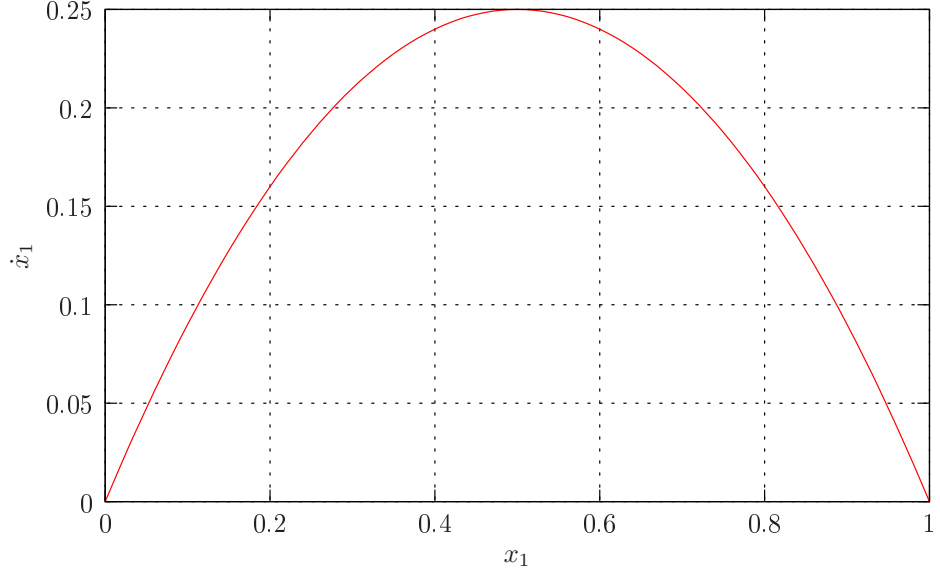


Figure 4.7: Class-5 Differential Equation

Class-5:

$$\begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}$$

$$\begin{aligned} \dot{x}_1 &= x_1^2(1 - x_1) + (1 - x_1)^2 x_1 \\ &= x_1(1 - x_1)(x_1 + (1 - x_1)) \\ &= x_1(1 - x_1) \end{aligned}$$

One species is a self-replicase and also a parasite of the other species; the latter is self-inert, and thus acts as an obligate host. As with class-2, the self-replicase has a positive rate of growth at all concentrations $x_1 > 0$ (see Figure 4.7) and will essentially deterministically invade and displace the second species, even if starting from arbitrarily low initial concentration. The difference when compared to class-2 dynamics is that, even if the self-replicase itself is initially in low concentration, it will still achieve an exponential replication rate (because it receives replication support from its obligate host, which is, by assumption here, in high concentration). Thus the displacement can be initiated more rapidly than for class-2 (see Figure 4.8); in this specific example, from the same starting concentration of 0.05, the takeover is complete by time $t \simeq 10$ for class-5, compared to $t \simeq 30$ for class-2. Formally, the growth law for the class-5 case is strictly logistic. Global replication rate, R , in the class-5 case is always equal to x_1 .

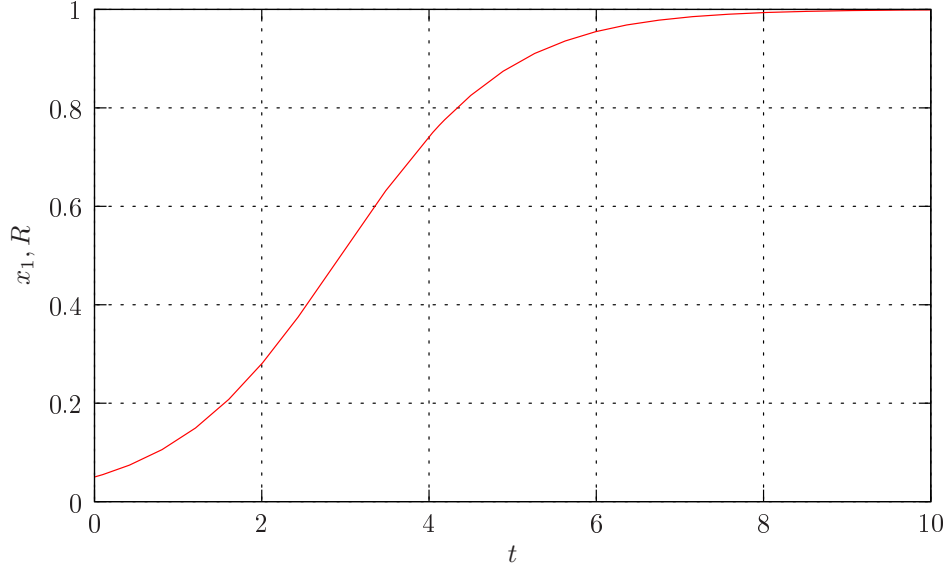


Figure 4.8: Class-5 Trajectory: Logistic Takeover [$x_1(0) = 0.05$]

Class-6:

$$\begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

$$\dot{x}_1 = -x_1(1 - x_1)^2$$

In a class-6 pairwise system, both species are self-replicases, though one is also a facultative parasite of the other, which could be called a facultative host. Figure 4.9 shows that $\dot{x}_1 < 0$ for all $x_1 > 0$, so x_1 will inevitably be displaced in the reactor. Since the displacing species is also a self-replicase, however, the global reaction rate, R , only suffers slightly during the transition between species, and eventually reaches its maximum value, 1.0, as the invading species completely takes over the reactor. This is a properly “selective”, quasi-deterministic, displacement of one self-replicase by another, because, at all relative concentrations, the latter achieves a higher replication rate. The displacement will take place even if x_2 initially has essentially arbitrarily low concentration. Figure 4.10 shows an example trajectory for $x_2(0) = 0.05$. While there is an initial period of slow growth, the parasite does inevitably achieve a “critical” concentration after which the displacement is completed very rapidly. Since there are no intrinsic fitness differences between the molecular species, however, each would appear identical to the other if incubated in separate reactors. It is only when incubated together that the class-6 dynamic is apparent. When we move from a single instance of a fixed size reactor to a population of variable sized reactors, protocells, (Section 4.4.2), then class-6 dynamics will become a useful mechanism for upward propagation of molecular level dynamics into mutation events visible at the protocell level. This mutation event may in

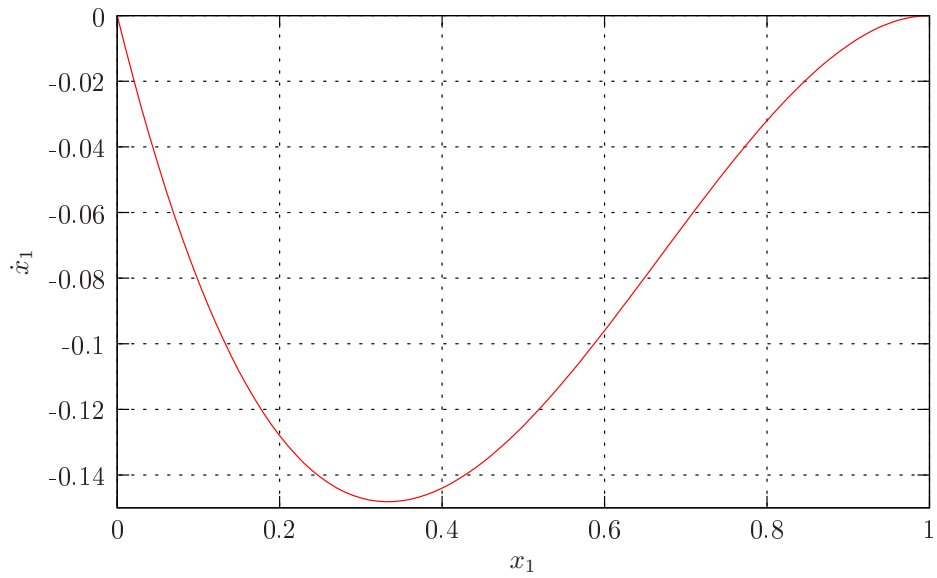


Figure 4.9: Class-6 Differential Equation

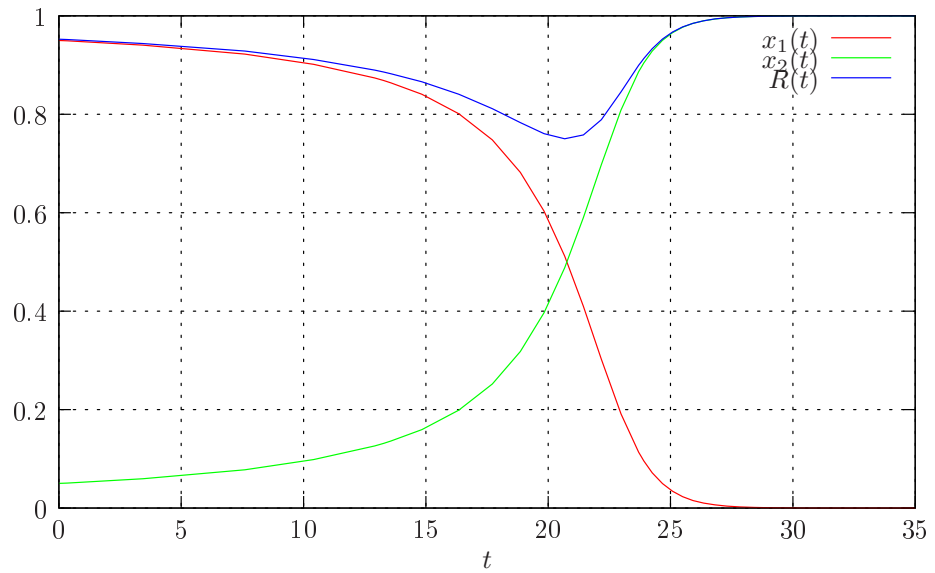


Figure 4.10: Class-6 Trajectory: Selective displacement by facultative parasite [$x_1(0) = 0.95$]

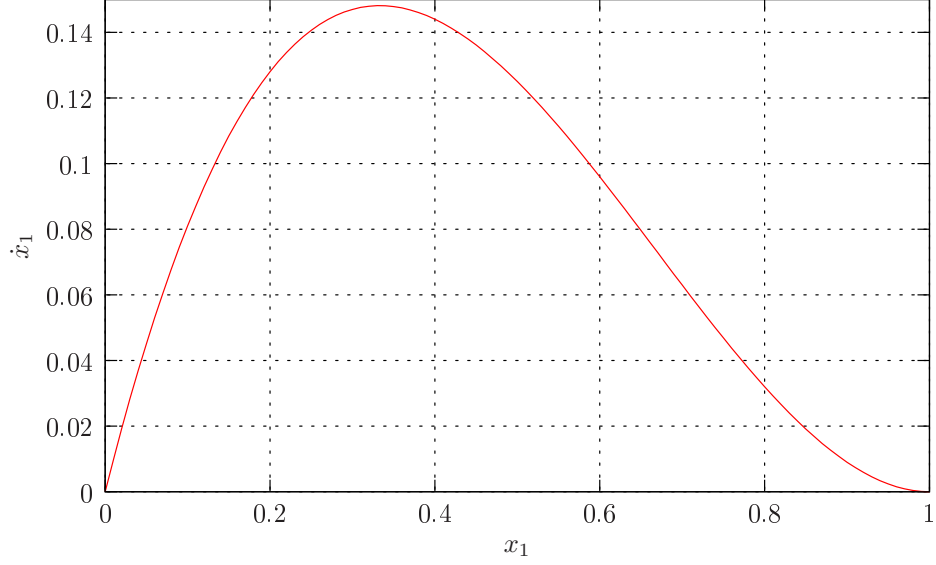


Figure 4.11: Class-7 Differential Equation

turn be used as a target for selection at the protocell level, assuming the protocell dynamics can be coupled to the molecular dynamics appropriately.

Class-7:

$$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\dot{x}_1 = x_1^2(1 - x_1)$$

One species is a self-replicase; the other is an obligate parasite of it. However, unlike class-4, the parasite also functions as a host for the self-replicase. This means that the (parasitic) replication service provided by the self-replicase to the self-inert species is offset by the (parasitic) replication service provided by the self-inert species back to the self-replicase and that eventually, the self-replicase will displace the obligate parasite, as seen in Figure 4.12.

Class-8:

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\dot{x}_1 = x_1(1 - x_1)^2 - x_1^2(1 - x_1)$$

Neither species is a self-replicase, but each can act as a replicase (obligate host) for the other (obligate parasite). Setting $\dot{x}_1 = 0$, and neglecting the case of $x_1 = 0$, there is a single fixed point at: $x_1 = 0.5$. Of course, we also therefore have $x_2 = 1 - x_1 = 0.5$ in this state (which would follow from the symmetry of the situation

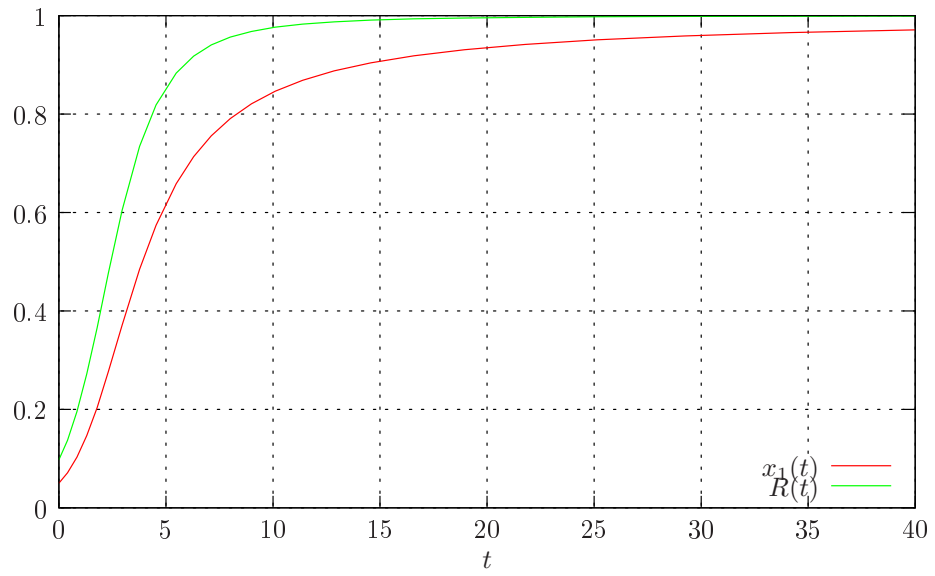


Figure 4.12: Class-7 Trajectory: Growth of self-replicase [$x_1(0) = 0.05$]

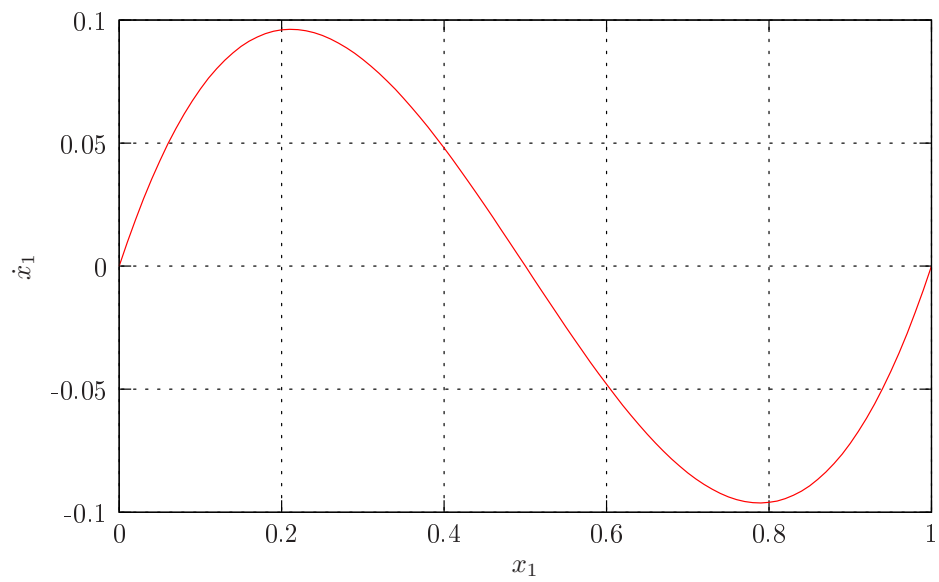


Figure 4.13: Class-8 Differential Equation

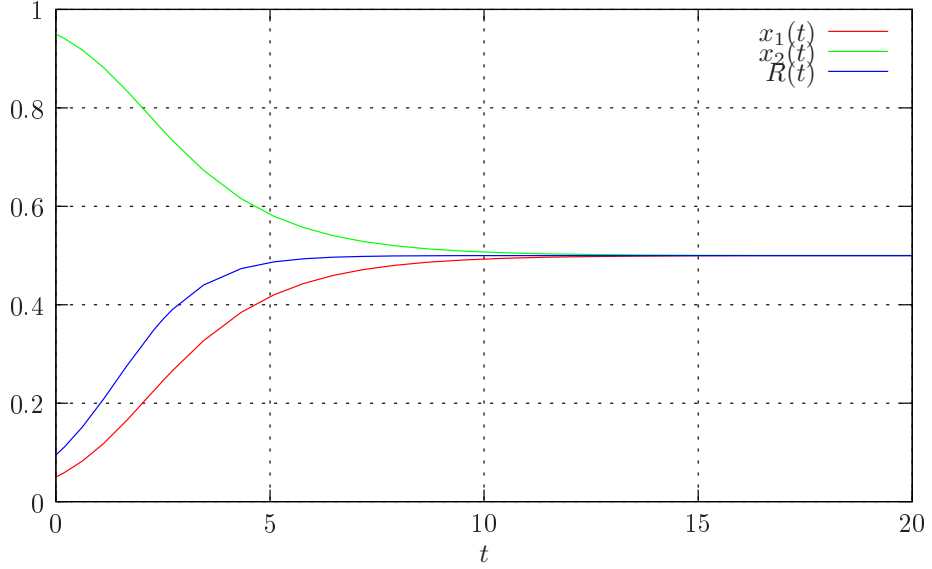


Figure 4.14: Class-8 Trajectory: Stable Co-existence/Hypercycle
[$x_1(0)=0.05$]

in any case). As shown in Figure 4.13, for $x_1 < 0.5$ we have $\dot{x}_1 > 0$ and for $x_1 > 0.5$ we have $\dot{x}_1 < 0$. Accordingly this is a stable fixed point. The reactor will return to it, even under essentially arbitrarily large perturbation. Figure 4.14 shows an example trajectory, initialised from $(x_1, x_2) = (0.05, 0.95)$. This essentially means class-8 networks are formal examples of two-component hypercycles in the sense of Eigen and Schuster (1977). Neither species can replicate itself but each can catalyse the replication of the other. The concentration of each species will therefore be maintained at exactly equal levels. In the presence of mutation, however, these two-component hypercycles are subject to the same weaknesses as larger, $n > 2$ hypercycles: should a parasite emerge such that it receives benefit from one part of the cycle and does not pass it on to the next member of the cycle, the hypercycle breaks down and the parasite wins.

Class-9:

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\dot{x}_1 = x_1^2(1 - x_1) - x_1(1 - x_1)^2$$

Both molecular species are independent self-replicases. In a certain sense this is exactly complementary to class-8. The expression for \dot{x}_1 is precisely the negation of that in class-8. Accordingly, the dynamics have exactly the same fixed points: the two states of $x_1 = 0$, $x_1 = 1$ where one or the other species is no longer present, and the state $x_1 = x_2 = 0.5$. However, the latter state is now unstable

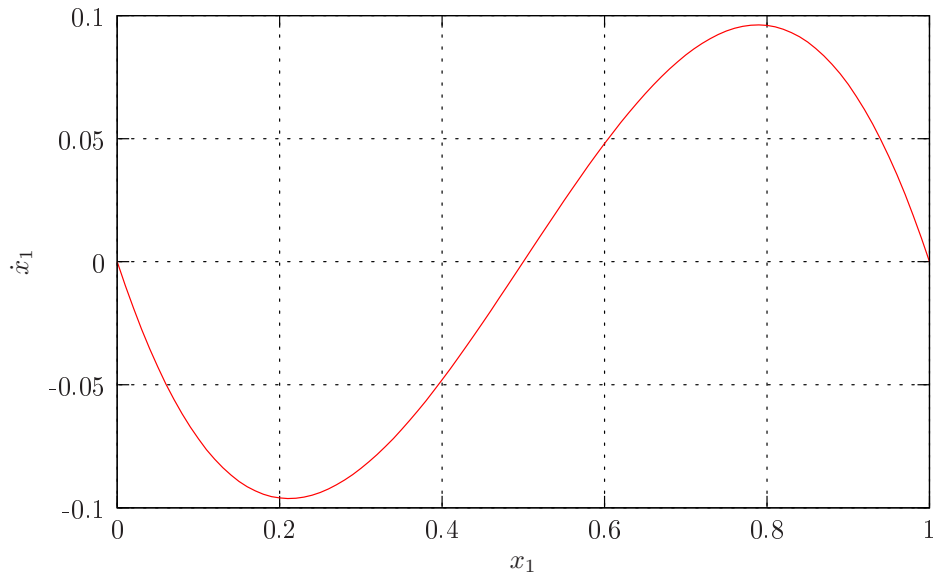


Figure 4.15: Class-9 Differential Equation

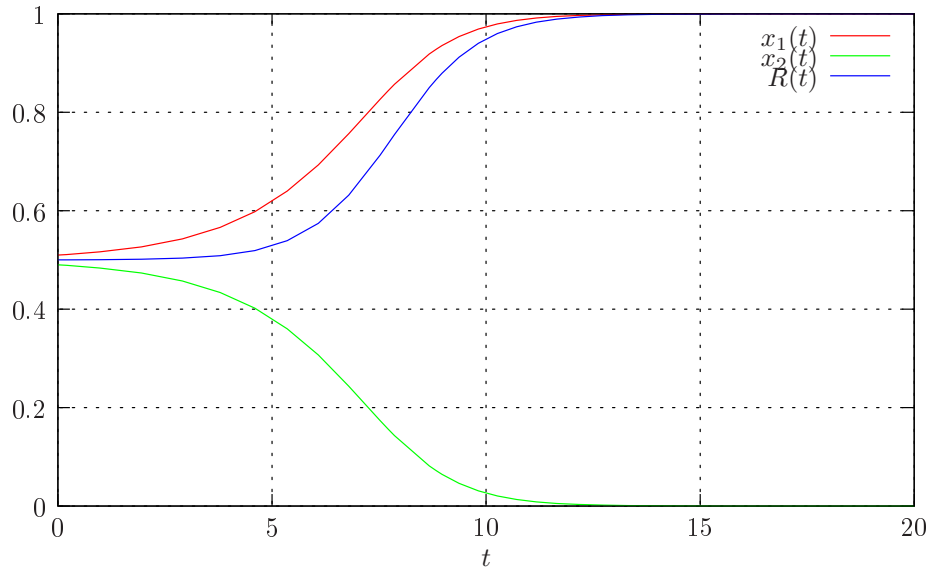


Figure 4.16: Class-9 Trajectory: Survival of the Common
 $[(x_1(0), x_2(0)) = (0.51, 0.49)]$

rather than stable. The effect is that any perturbation means the system rapidly collapses into a state where one species displaces the other. Figure 4.16 shows an example trajectory initialised from $(x_1, x_2) = (0.51, 0.49)$. In effect, there is positive frequency-dependent selection, so that whichever species first achieves a higher concentration benefits from a positive feedback effect which further amplifies its concentration until it takes over the complete reactor. The effect is well known, and has been termed “the survival of the common” by Szathmáry and Maynard Smith (1997). It is characteristic of any system of replicators undergoing hyperbolic rather than exponential growth.

4.5.4.2 Limitation of $n = 2$ analysis

As described at the beginning of this discussion (Section 4.5.1), the exhaustive analysis of $n > 2$ systems is infeasible due to the number of potential networks that arises. Since the majority of the interesting experiments that will be undertaken with the MCS will involve significant mutation rates, we will necessarily be dealing with $n > 2$ systems however. This highlights a significant limitation of this analysis. We will take the approach of assuming that the complicated dynamics of $n > 2$ systems can be approximately described by the superposition of multiple instances of $n = 2$ systems concurrently present in the same reactor.

4.6 Conclusion

An agent based artificial chemistry, MCS, was described in this chapter. MCS was designed to be an artificial life system which supported hierarchical selection. The hierarchy of MCS is based on two interacting sub-systems, molecular chemistry and protocell chemistry. The molecular chemistry has been designed to act as an idealised template replicator world, and the protocell chemistry has been designed to act as an idealised lipid membrane chemistry. These ideal systems should make it possible to carry out experiments in the later chapters which are focused on the complex evolutionary dynamics that emerge from the interaction between these systems, rather than the low level chemical makeup of the individual systems in isolation. In the upcoming experimental chapters, the MCS will undergo several modifications as these systemic evolutionary dynamics are uncovered. The experimental work culminates with a demonstration of ongoing protocell evolution, driven by the molecular dynamics of the lower level in the hierarchy.

Part III

Experimental Results

Chapter 5

MCS-0

This chapter presents the results of the MCS-0 experiment set. The MCS-0 experiments explore the basic replicator dynamics of the MCS platform and form the foundation for the work presented in the later experimental chapters (Chapter 6 and Appendix B). The aim of this chapter is thus to explore the properties and evolutionary potential of MCS-0. This exploration takes the form of making some predictions followed by experimental verification of those predictions. The focus of the experiments in this chapter is to explain the low-level molecular interaction dynamics of MCS-0 with a view to identifying possible molecular-level traits that may be used to apply selectional pressure at the protocell level. The core prediction (and experimental verification) explored in this chapter is that in the basic case of complete molecular recognition (MCS-0), class-6 facultative parasite displacement is the only expected outcome. These class-6 parasitic takeovers coincide with a length increase in the molecular bit-strings due to the mechanism of parasitism involved. When hierarchical selection is applied here via artificial protocells, we see the somewhat counter-intuitive result of selectional stalemate, as the selectional pressures at the molecular level and the protocell level of the hierarchy act in direct opposition.

5.1 MCS-0 Specification

MCS-0 is characterised by a molecular binding scheme which enforces molecular binding rules comparable to the “complete molecular recognition” concept presented by Altmeyer et al. (2004). Specifically, two molecules may bind if the catalyst is a substring of the substrate, and subsequently, a copy of the substrate molecule is produced and added to the molecular population. Some examples of molecular matching are shown in Table 5.1.

With respect to the binary replicase interaction classes presented in Section 4.5.4, the molecular binding scheme of MCS-0 will only allow system dynamics of class-

Substrate	00101000	Substrate	00101000	Substrate	00101000
Catalyst	1010	Catalyst	1111	Catalyst	0010100
Match?	yes	Match?	no	Match?	yes

Table 5.1: Examples of Molecular Matching (MCS-0)

9 and class-6. Qualitatively, since the MCS-0 binding scheme is straight-forward sub-string matching:

- all MCS-0 molecular species are self-replicases (all strings are sub-strings of themselves)
- all proper superstrings of any given species will be facultative parasites of it (by definition, the host is a substring of the parasite)
- no other relationships are possible in MCS-0 for self- and binary-interactions.

The remainder of this section discusses in detail some of the phenomena that are expected to heavily influence the outcome of the experiments. This discussion is then built upon to make predictions based on the corresponding ODEs from Section 4.5.4 and subsequently undertake a programme of experiments to test the validity of these predictions.

5.1.1 Molecular Selection

In the case of a single, fixed-size tank reactor (Section 4.4.2), selection applies at the level of individual molecules, or lineages of particular molecular species. Since every MCS-0 molecule is a self-replicase by definition, we would expect there to be “quasi-stable” periods of dominance by a particular replicator species which are punctuated by occasional displacement events by fitter replicators. This is due to the fact that we are dealing with hyperbolic growth which gives rise to the phenomenon of the “survival of the common” (Szathmary and Maynard Smith, 1997). Hyperbolic growth is characteristic of systems, like MCS-0, which rely on an



production rule where two molecules react to produce a third. In contrast, exponential growth arises from the



production rule, where a single molecule undergoes some transformation resulting in two copies of the original. During these predicted quasi-stable periods of single

species dominance, the system dynamics can be approximately described as a pairwise class-9 dynamic. Of course, in the presence of mutation, there will be many more than two species present in a reactor at a given time, so the true system dynamics can be better described as a superposition of multiple pairwise class-9 systems—each pairing between the dominant species and one of the distinct mutant species giving rise to a separate parallel or concurrent class-9 system in the reactor. In fact, all the *other* pairings would also normally constitute class-9 systems, but as the concentrations involved remain very small, we neglect all but the pairings involving the dominant species.

5.1.1.1 Prediction: Quasi-Steady-State Dominant Species Concentration

During time-periods where there is clear dominance by a particular molecular species, we can make some quantitative predictions about the steady-state concentration of that molecular species in the reactor.

Firstly, take x as the concentration of the dominant molecular species, X , and y to be the combined concentration of all other molecules in the reactor which we shall denote Y . In other words:

$$x + y = 1 \tag{5.3}$$

We will assume that Y is sufficiently diverse that the individual concentrations of each molecular species are low enough that we can neglect replications between identical molecules in Y . Furthermore, we assume that there is no species present which is parasitic relative to the dominant species. Of course, such an assumption depends critically on the choice of mutation rate, discussed further in Section 5.1.3, though for the purposes of this discussion, we assume that this mutation rate is set to an appropriate value to uphold our earlier assumption. So, while the reactor is actually a superposition of multiple parallel class-9 systems, the assumptions just described allow us to idealise the mutant population of the reactor and describe it in terms of a superposition of multiple class-2 pairwise systems with respect to the dominant species. While the assumption of a class-2 dynamics may seem contradictory to the specification given in Section 5.1, such an assumption is justified given that class-9 dynamics give rise to the phenomenon of survival of the common, and in such cases we can be sure that none of the individual mutant species will achieve significant concentrations. A further idealisation that arises from these assumptions is that, among the mutant species, we can assume a superposition of parallel class-0 systems. In this way, we can treat the combined effects of the mutant spectrum as a single “equivalent” species, with class-2 dynamics relative to the dominant species,

Catalyst	Substrate	Product	Replaces	Net Effect	Probability
X	X	X	X	-	$x.x.(1 - V).y$ $x.x.V.x$
X	X	X	Y	increase x	
X	X	Y	X	decrease x	
X	X	Y	Y	-	
X	Y	-	-	-	
Y	Y	-	-	-	
Y	X	-	-	-	

Table 5.2: Interactions in a Fixed-Size Reactor

and use this equivalence to make an approximate quantitative analysis of the ongoing mutational load in the reactor.

As described earlier in Section 4.4.5.2, replication of molecules in MCS-0 is subject to a per-bit mutation rate, v . So, for a molecule of bit-length l , we can compute the probability of producing a faithful copy of that molecule as:

$$(1 - v)^l \quad (5.4)$$

And consequently, compute the probability of producing a mutant copy of the molecule, V , as:

$$V = 1 - (1 - v)^l \quad (5.5)$$

Since we are dealing with a fixed size reactor (Section 4.4.2), the general behaviour is that as new molecules are added to the reactor as a result of successful replication events, they replace randomly chosen molecules from the reactor to keep the overall size constant. The possible interactions and their net effects on the fixed-size reactor are enumerated in Table 5.2. An approximate differential equation for the change in concentration, \dot{x} , can be derived from this table by summarising the outcomes as follows:

$$\begin{aligned}
\dot{x} &= (x^2(1 - V)y) - (x^3V) \\
&= x^2((1 - V)y - Vx)
\end{aligned} \quad (5.6)$$

Given that $x + y = 1$, this yields

$$\dot{x} = x^2(1 - x - V)$$

By setting \dot{x} to zero, we can determine the steady-state. Neglecting the case

where $x = 0$, this leaves one fixed point with:

$$\begin{aligned} 1 - x - V &= 0 \\ x &= 1 - V \end{aligned} \tag{5.7}$$

V here is then a lower-bound on the standing mutational load in the reactor; as this estimate of the mutational load neglects all replications among the mutant species, the actual concentration of mutant molecules is expected to be somewhat higher. If the individual concentrations of particular molecular species within the mutational cloud¹ increase, they may begin to undergo successful self-replication events, albeit at a slower rate than that of the dominant species. These “off-sequence” replication events will necessarily contribute to the global reaction rate. In theory, it should be possible to revisit the analysis to incorporate the approximate contribution of the off-sequence reactions within the mutational cloud, but it is not clear how to avoid either over-estimating or under-estimating the extent of this contribution. Nevertheless, these species will be present in at least low relative concentrations for non-zero mutation rates, and their presence would give rise to significant statistical fluctuation from any ODE estimate.

By assuming that individual species concentrations in the mutational cloud are low enough to give rise to pairwise class-0 dynamics, we can predict the replication rate for the dominant species, and therefore the approximate *overall* reaction rate as:

$$(1 - V)^2 = (1 - v)^{2l} \tag{5.8}$$

One consequence of this is that the steady-state overall reaction rate will fall as l , the bit-string length of the molecules, increases.

5.1.1.2 Prediction: Selective Displacement by Facultative Parasites

The quasi-stable periods of dominance predicted in Section 5.1.1.1 are expected to be punctuated by occasional displacement events by species which are facultative parasites relative to the incumbent dominant species. Such events are characterised by the class-6 system dynamics presented in Section 4.5.4, though in practice, the stochastic dynamics of the MCS make the occurrence of these events quite unpredictable. In particular, if the absolute number of parasite molecules is small, then it is not at all certain that such a displacement event will follow. In other words,

¹This is not strictly a “quasi-species” in the sense of Eigen (1989), since there are no intrinsic fitness differences between the molecular species in question.

Substrate	0000	Substrate	0000
Catalyst	0000	Catalyst	00001
Match?	yes	Match?	no
Substrate	00001	Substrate	00001
Catalyst	0000	Catalyst	00001
Match?	yes	Match?	yes

Table 5.3: Example of Molecular Parasitism in MCS-0

the presence of a single instance of a parasite species is necessary but not sufficient to cause a parasitic displacement event.

5.1.1.3 Prediction: Elongation “Ratcheting”

As a direct consequence of the substring molecular binding rules outlined in Section 5.1, we can also predict that a monotonic increase in average molecule length will accompany the displacement events presented in the previous section. To understand this, note the sample interactions of Table 5.3—the parasite must be *at least* one bit longer than its host, and this is precisely because of the mechanism of molecular binding. In this example, the parasite is identical to its host with the exception of a one bit prefix or suffix. Since it is impossible for any species which is either the same length or shorter than the incumbent dominant to take over the reactor, this length increasing phenomenon is perhaps described better as elongation “ratcheting”. The cumulative effect of this ratcheting will be that the global reaction rate for any particular reactor will quasi-deterministically decrease over time because longer molecules give rise to a lower overall reaction rate due to higher mutation rates (Equation 5.8). Higher mutation rates mean that for each replication event, there will be a progressively lower chance of producing faithful copies of the parent molecule. As this phenomenon manifests, the quasi-steady-state concentration of the current dominant molecular species will become lower and lower. This will have a direct effect on the global reaction rate of the reactor as it becomes decreasingly likely that two compatible copies of a replicase are selected for reaction.

5.1.2 Parasitic Mutation Network

In Section 5.1.1.1, it was shown that the quasi-steady state of a reactor dominated by a self-replicase was such that the concentration of the dominant species was $(1 - V)$ and the combined concentration of the mutant cloud was V . Here, we define the “parasitic mutation network” as the subset of molecular species in the mutational cloud which are parasites of the dominant species. The topology of this network is that of a directed graph, where the nodes correspond to distinct molecular species,

and the edges correspond to mutations that occur with some non-zero rate. Apart from the node which represents the dominant species—the central node—all nodes represent molecular species which are class-6 facultative parasites of the species represented by the central node. The edges carry weights which reflect the actual *rate* at which that mutation occurs. Directionality applies because, in general, the rate at which a species A mutates to B need not be equal to the rate at which B mutates to A. In fact, for the MCS system, and dealing strictly with parasitic mutations, not only need these rates not be equal, if species A can mutate to species B then species B definitely *cannot* mutate to species A. In summary, the “parasitic mutation network” is highly non-uniform in the senses that:

1. some species give rise, with significant probability, to more or fewer parasitic mutants than others,
2. the overall rates of production of parasitic mutants vary significantly between species.

5.1.3 Further Characteristics of Molecular Mutation

In the previous section, the highly non-uniform nature of the “parasitic mutation network” was described. By that description, a reactor currently dominated by a species with many parasitic nodes and whose edges had relatively high weights would therefore be expected to undergo selective displacement sooner than a reactor dominated by a molecular species which had a lower tendency to produce parasitic mutants. Given the implementation of the mutational process, the probability of a particular mutation occurring falls rapidly with increasing Levenshtein² distance. We can therefore make the problem more tractable if we consider only mutations that are within one unit of levenshtein distance of some “master” sequence. On this basis we can calculate a lower bound on the number of facultative parasites that *any* given molecule would have. To do this, recall that for a molecule to bind to another molecule, the structure of one of the catalyst molecule must be a substring of the substrate structure (see Table 5.1). So, taking an arbitrary molecule as our starting point, we are guaranteed to produce facultative parasite mutants if we add a single bit (0 or 1) to the beginning or the end of the starting molecule without modifying any other bits. Given a particular set of mutation rates then we can compute how often these “target” mutations would arise, giving us a lower bound on the rate of production of facultative parasites. As described in Section 4.4.5.2,

²Levenshtein distance is a metric which represents the number of modifications that need to be made to sequence A to change it into sequence B—the so-called *edit* distance.

Sequence	P(parasite)
1000000001	4.02×10^{-4}
0000000000	9.05×10^{-4}

Table 5.4: P(parasite) from 10^6 replication events

mutations in MCS-0 can be either bit-insertions, bit-deletions or bit-flips. These mutations happen at different rates to one another.

A general expression representing the lower bound probability of generating a parasite from a given molecular replication event can then be produced as follows. Taking l as the length of the molecular sequence which is being replicated, v as the per-bit mutation rate, β as the proportion of mutations that are bit-flip mutations— $(1 - \beta)$ is therefore the proportion of length-changing mutations, and since we are particularly interested in length *increasing* mutations, $\frac{(1-\beta)}{2}$ is the proportion of mutation events that result in a length increasing mutation. The probability of a single bit-insertion either at the beginning or the end of a sequence, with no other mutations, can therefore be expressed as:

$$[(1 - v)^{(l-1)}](\frac{v(1 - \beta)}{2})$$

Since either mutation (beginning insertion or end insertion) will result in a parasite, we can multiply this expression by 2 to give a general expression for the lower bound probability of generating a parasite from the replication of an arbitrary molecular sequence:

$$(1 - v)^{(l-1)}(v(1 - \beta)) \tag{5.9}$$

Taking $v = 0.05$, $l = 10$ and $\beta = 0.99$, values typical for the simulation runs presented in this chapter, we can compute this lower bound as 3.15×10^{-4} . We are now in a position to revisit the above claim that “some molecular species have a higher propensity to produce mutants that are facultative parasites than others”. Specifically, there are two distinct mechanisms by which these differences in parasite production arise. Firstly, and most obviously, molecular species of different lengths will have different lower bound parasite production rates, by definition. Secondly, a potentially larger effect is that the precise bit-pattern of the string may allow more or fewer parasites to be generated as a result of mutations at positions *other* than the beginning or the end. This effect is not captured at all by the lower bound analysis, but *can* usefully be investigated for specific example strings using a monte carlo approach.

Table 5.4 shows the results of a monte carlo style investigation which took the

given molecular sequences and made 1×10^6 replications using the same mutation rate parameters as will be used in subsequent MCS-0 experiments. Of these, the number which resulted in the production of a parasite molecular species were noted. For example, the sequence ‘10000000001’ is highly sensitive to mutations within the sequence and therefore the value obtained for it (4.02×10^{-4}) approaches the lower bound. Nonetheless, it comes in slightly above the lower bound because the structure of the bit-string of this molecule *does* permit a certain level of internal mutation in the case where multiple mutations may cancel each other out—a possibility which is explicitly ignored by the lower-bound analysis. By contrast, the sequence ‘00000000000’ is far less sensitive to internal mutations, since the insertion of a ‘0’ bit anywhere in its structure will result in a parasitic molecule. It is not surprising then that the lower bound significantly under-counts the probability of producing parasites here (9.05×10^{-4}), by a factor of roughly three.

The implication of this variability in mutational characteristic is that the length of time between the successive parasitic displacements predicted in Section 5.1.1.2 will vary significantly depending on the particular molecular species that is “dominant” in the reactor.

5.1.4 Cellular Selection

The exploration of the effects of hierarchical selection is one of the main objectives of this thesis. Each of the Major Evolutionary Transitions, as presented by Maynard Smith and Szathmáry (1997), describe the emergence of a new level of Darwinian selection that is an aggregation of “individuals” from before the transition. The containers that will be provided to MCS-0 replicators in the hierarchical experiments presented here take the form of artificial protocells, similar in concept to lipid vesicles or modern cellular membranes.

Since these protocell containers will grow and divide at a rate which depends on the growth rate of the contained molecular population (Section 4.4.4), we can expect to see competition between protocells over a fitness landscape which is based on heritable characteristics of the composition of the replicator populations inside the cells, assuming there *are* any heritable differences (Section 5.1.4.3). For the purposes of the experimental systems presented here and in Chapter 6, cells which have molecular populations that exhibit a high reaction rate will be fitter than cells whose interior molecules have a lower overall reaction rate³. By applying population constraints at this higher level (either in terms of absolute numbers of molecules

³In Appendix B, modifications will be made to the cellular dynamics which will modify this behaviour somewhat.

allowed, or in terms of absolute numbers of protocells allowed (Section 4.4.2) we can set up a process of Darwinian evolution at the protocell level. Again, having some traits which are heritable at the protocell level is absolutely crucial to this process of evolution.

At the moment, in this highly idealised hierarchical model, the only coupling between the molecular level and the protocell level of selection will be the reaction rate of the molecules, which, by design, governs the growth rate of the protocell. Earlier, we saw that the predicted outcome for fixed-size populations of MCS-0 molecules was that of quasi-stability punctuated by displacements by facultative parasites (Sections 5.1.1.1 and 5.1.1.2). This causes a molecular length ratcheting scenario—where progressively fitter molecules take over the reactor with a direct negative effect on overall molecular reaction rates (Section 5.1.1.3). Since we know that protocell growth rates are proportional to molecular reaction rates in this system (Section 4.4.4.2), we can predict that the protocell level of selection should counteract the parasitic takeovers at the molecular level and in doing so, prevent the length ratcheting that was inevitable before the application of hierarchical selection. Implicit in this assumption is that such protocells are capable of producing offspring sufficiently similar to themselves so as to meet the criteria for natural selection (Section 5.1.4.3).

5.1.4.1 Steady-State Concentrations in Artificial Protocells

The analysis presented earlier in Section 5.1.1.1 can be further applied to variable sized reactors—artificial protocells—due to the equivalence with tank reactors as presented in Section 4.4.2.3. The key difference with this analysis is that the size of the reactor, measured in terms of the number of molecules within it, grows until the protocell splits. In Section 4.4.2.3, it was shown that an appropriate dilution term is necessary for the ODE to account for the continuously growing molecular population within a protocell. Assuming that *all* molecular reproduction is due to the activity of species X , then the molecular population grows at a rate of x^2 . Taking the pseudo species Y as the aggregate of all mutant species in the protocell, then the total species concentration of the protocell is given by $(x + y) = 1$, regardless of the instantaneous number of molecules in the protocell. However, as the absolute numbers of species X and species Y increase, the concentrations x and y must be re-scaled appropriately for the instantaneous growth rate of the protocell, x^2 . Each concentration needs to be scaled in proportion to its instantaneous concentration. In this way, the dilution term for x becomes $x.x^2 = x^3$, and the generalised ODE for a growing molecular population becomes:

$$\dot{x} = (x^2(1 - V)) - x^3$$

and solving for ($\dot{x} = 0$), we get:

$$\begin{aligned} 0 &= (x^2(1 - V)) - x^3 \\ x^3 &= x^2(1 - V) \\ x &= 1 - V \end{aligned}$$

This of course is the same as Equation 5.7 due to the equivalence presented in Section 4.4.2.3. This means that the quasi-steady-state dynamics for protocells are the same as the dynamics for a fixed-size reactor, and that the predictions for parasitic takeover of fixed-size reactors also apply to protocells under the appropriate dilution term.

5.1.4.2 Protocell Gestation Times

Using the “steady-state” analysis for protocells derived in the previous section, we can further derive an estimate of the gestation times of protocells—in other words, the amount of time taken for a protocell to grow large enough to divide. This analysis assumes that an appropriate division threshold has been set, and that at $t = 0$, the protocell in question has just arisen by binary fission from its parent—in other words, its size is half that of a fully grown protocell. We know that the approximate molecular replication rate is x^2 , or $(1 - V)^2$. Taking $(X + Y) = N$ as the number of molecules inside a protocell then, we can model the growth of a cell as:

$$\frac{dN}{dt} \approx kN$$

where $k \approx (1 - V)^2$ for which the solution is:

$$N = N(0).e^{kt}$$

from which we can derive the gestation, or “doubling” time, τ , as follows:

$$\begin{aligned} N(\tau) &= 2N(0) \\ N(0).e^{k\tau} &= 2N(0) \\ e^{k\tau} &= 2 \\ k\tau &= \ln(2) \\ \tau &= \frac{\ln(2)}{k} \end{aligned}$$

which is approximately equivalent to:

$$\tau = \frac{\ln(2)}{(1 - V)^2}$$

at this point, recall that $V = 1 - (1 - v)^l$, so reintroducing this definition yields:

$$\tau = \frac{\ln(2)}{(1 - v)^{2l}}$$

From this derivation, we see that protocell gestation time, τ , shares a relationship with the average length of the molecules that are contained within it. Furthermore, in the presence of a constant, but greater than zero per-bit mutation rate, v , protocell gestation time increases as molecular length, l , increases. Of course, all molecules within a protocell may not be of the same molecular length. In this case, taking l as representative of the length of the dominant molecule in the protocell should suffice. In terms of protocell level fitness then, we can see that protocells which have a higher average molecular length have higher gestation times than protocells with lower average lengths. Under limited resources then, protocells which reproduce quicker (those with lower gestation times) will tend to displace those protocells with longer gestation times.

5.1.4.3 Protocell Inheritance Mechanism

In Sections 4.4.2.3 and 5.1.4.1, we saw that the quasi-steady-state concentrations of molecules in artificial protocells would be equivalent to the steady state concentrations that might be observed for a fixed-size reactor, as long as the number of molecules is sufficiently large enough to justify the ODE analysis. Based on this equivalency, we can make the further assumption that quasi-stable artificial protocells can also be approximated to a superposition of pairwise class-9 molecular species. From Section 4.4.4.3, we know that molecules are distributed by random assignment between daughter cells upon reaching maturity and undergoing binary fission. Given that molecular populations are subject to hyperbolic selection, ensuring survival of the common, it is reasonable to expect that the molecular populations of the daughter cells will be dominated by the same molecular species as that which dominated the parent. We can therefore use the molecular bit-string of the dominant molecule in a particular protocell to label that protocell, and furthermore, this label is a heritable trait and can be used to classify strains of protocells.

5.1.4.4 Protocell Mutations

The protocell inheritance mechanism just described breaks down in the situation where the mutational cloud contains an instance of a molecular species which is a parasite relative to the dominant species of the protocell (Section 5.1.1.2). When this happens, the protocell can no longer be reliably treated as a superposition of pairwise class-9 molecular species, and the “survival of the common” dynamic no longer applies. In Section 5.1.1.2, it was argued that the presence of a *single* parasite was necessary, though not sufficient, to cause a displacement event at the molecular level. The variability of outcome here is due to statistical fluctuation at low absolute numbers of parasites. Using “sufficiently” high numbers of molecules, however, reduces this case to an approximately continuous process where the class-6 parasite ODE analysis presented in Section 4.5.4 applies. Experimental data are presented later as Figure 5.7 which explores this variability for different, relatively low absolute numbers of parasite molecules.

From the previous section (Section 5.1.4.3) we know that the molecular species of the dominant molecule is a reliably heritable trait for protocells, in the absence of molecular level parasites. When parasitic mutants of the dominant molecule *are* present, the fate of that particular lineage of protocells depends on whether or not the parasite takes over. We shall call protocells who have a dominant molecular species which has no parasitic mutants present in the population “pure” cells, and classify them by the sequence of their dominant molecule. By this definition, protocells which have even a single instance of a parasitic molecule are not “pure” cells, and will therefore be labeled as “mixed” cells. This classification will take place at protocell birth, meaning that there will be cases where a daughter cell will have a different classification to its parent, and potentially even a different classification from its sibling. Furthermore, by a mechanism similar to the “stochastic corrector” model proposed by Szathmáry (1986), it might indeed be possible for a cell which was born as mixed cell to give rise to a daughter cell which is of the original pure line. In the simplest case, this can happen as follows:

- the original mixed cell produces a single instance of a parasite of the dominant molecule,
- throughout the lifetime of the cell, no further parasites were produced
- upon protocell division, one of the daughter cells possesses no parasite molecules, and the other contains the parasite molecule.

The foundation of a new protocell strain therefore occurs in two distinct phases:

1. one or more parasite molecules emerge in the parent cell, causing one or both of the daughter cells to be labeled as mixed cells upon binary fission,
2. the concentration of parasites in a mixed cell continues to grow, perhaps spanning multiple generations, until eventually, the parasite becomes the dominant molecule in the protocell, and a new pure line is founded.

In this way, we can completely isolate the behaviour of pure-line cells and clearly delimit the emergence of new “pure” protocell strains. Given these new protocell labels, “pure” and “mixed”, we can define some new terminology to describe the key molecular species within such protocells:

- The **Principal** species is that which has the largest single concentration, measured instantaneously—the majority species.
- The **Dominant** species is such that, against the background of all currently present species, and of its own non-parasitic mutants, it is likely to become the **Principal** species (if it is not already so, and modulo statistical fluctuation at low absolute population levels) and to stably remain as such (in the absence of parasites of it).

5.1.4.5 Protocell Mutation Rates

In the previous section, it was proposed that a protocell strain undergoes mutation into another new protocell strain by a two step process. Firstly, the presence of parasitic molecules of the dominant molecular species in a pure-line protocell would cause the offspring of that protocell to be relabeled as a mixed cell if they had been assigned some of these parasitic molecules during the fission event (Section 4.4.4.3). These protocells, and all of their progeny, would then be labeled as mixed cells until such time as the lineage matures to a pure-line once again. There is also the possibility for “back-mutation” at the protocell level. The concept of back-mutation in MCS protocells may seem counter-intuitive—unlike the fixed-size reactor, individual molecules are never explicitly removed. Consider a pure-line protocell which undergoes $\frac{S_{max}}{2}$ molecular replications during which a single class-6 parasite of the principal molecular species is produced. Following the fission event, the daughter protocell which was assigned with the parasite will be labeled “mixed”, and the other daughter protocell will be labeled “pure”. If that mixed protocell now undergoes $\frac{S_{max}}{2}$ molecular replications, and this time no further parasitic molecules are produced, then the daughter protocells will once again consist of one mixed protocell and one pure-line protocell. In this way, a mixed protocell has “back-mutated” into the original pure line.

Given a particular molecular species, we can estimate the rate at which it would be expected to produce parasite molecules either as a lower bound or by direct measurement for a particular molecular species by monte carlo trials (Section 5.1.3). Taking r as this parasite production rate, we can then derive an estimate of the rate of production of “mixed” protocells—the protocell mutation rate—as:

$$R = 1 - (1 - r)^{\frac{S_{max}}{2}} \quad (5.10)$$

This calculation significantly overestimates the *effective* protocell mutation rate, however, due to the fact that it does not account for back-mutation from mixed cells to pure-line cells.

5.2 Experiments: MCS-0

Following on from the predictions outlined above, a series of four experiments was carried out to test their validity. The first set of presented experiments (Section 5.2.1) is concerned with verifying the class-9 “Survival of the Common” dynamic. The second set of experiments (Section 5.2.2) is concerned with demonstrating the displacement of a dominant species by a facultative parasite, as predicted by the ODE analysis of Section 4.5.4. In the third set of experiments (Section 5.2.3), we will see that a direct consequence of these selectional displacements is a length ratcheting effect that manifests, in this case, with the average molecular length increasing monotonically which causes a direct reduction in global reaction rate for the system. Finally, the fourth experiment set (Section 5.2.4), applies hierarchical selection as a mechanism to explore the evolutionary potential for populations of MCS molecules.

5.2.1 Survival of the Common

As mentioned earlier (Section 5.1.1), systems like MCS-0 which use a production rule like:



result in a situation of hyperbolic growth. This means that any molecular species which is already at a higher concentration than other species will continue to increase in concentration, whereas those which are at lower concentrations will be driven lower still—all other things being equal. In MCS-0, there are no intrinsic fitness differences between molecular species of the same length, and only relatively small differences between molecules of similar length. It is important therefore to

verify that the predicted “survival of the common” phenomenon does indeed hold, especially in the case of same-length molecular species which have no intrinsic fitness differences. Any alternative outcome would imply that displacement of one molecular species by another is necessarily a selective displacement via a class 6 “parasitic takeover” dynamic, and *not* a result of random drift.

The “survival of the common” dynamic predicts that if a reactor is seeded with equal concentrations of two distinct molecular species, then whichever first gains the upper hand will displace the other. To verify this dynamic, a reactor was set up under the following conditions:

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with equal concentrations (0.5) of two molecular species: $x_1 = 0000000000$ and $x_2 = 1111111111$,
- the per-bit mutation rate, v , was set to 0.

From 100 unique trials, species x_1 (0000000000) displaced species x_2 (1111111111) 53 times. Total displacement occurred across all runs at a mean time of 14.7, and standard deviation of 2.3, where time has been scaled by dividing by the total number of molecules in the population, 1×10^3 .

Since the “survival of the common” dynamic is due to the class-9 pairwise relationship as described in Section 4.5.4, we can plot the solution to the ODE for class-9 against the results of a typical experimental set to verify the phenomenon. Figure 5.1 contains such a plot, where the theoretical and experimental curves have been aligned along the time axis based on the mid-point of the transitions ($x_1Sim = x_1 = 0.75$).

The ODE solution for class-9 can also be used to predict the global reaction rate for these experimental runs, and the plot shown in Figure 5.2 shows that the global reaction rate of the experiment depicted in Figure 5.1 does indeed correspond with the predicted rate. The experimental and ODE plots have been aligned at a concentration value of 0.75.

The “survival of the common” dynamic should also apply to systems where the number of unique species is greater than two. A further demonstration of the survival of the common dynamic is shown in Figure 5.3. In this case, a reactor was set up as follows:

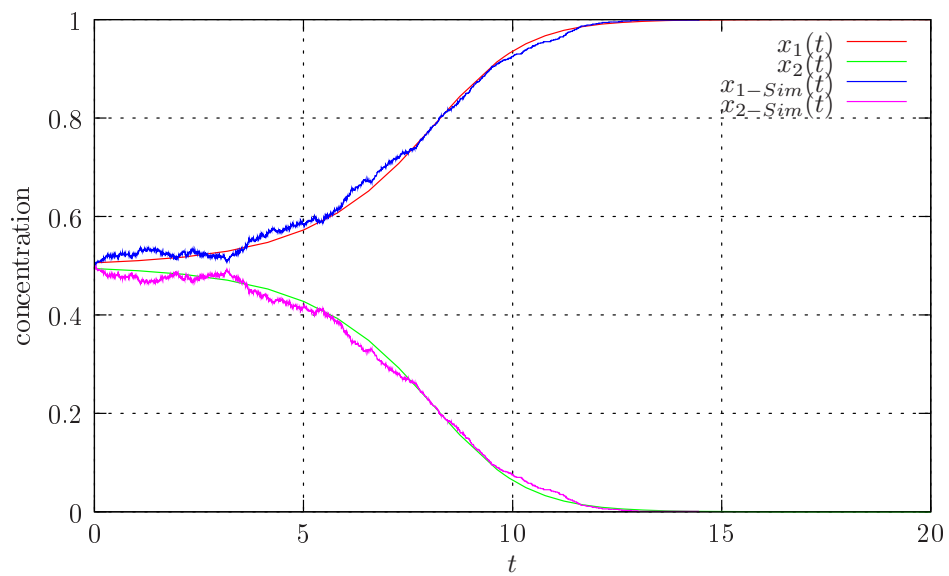


Figure 5.1: Survival of the Common

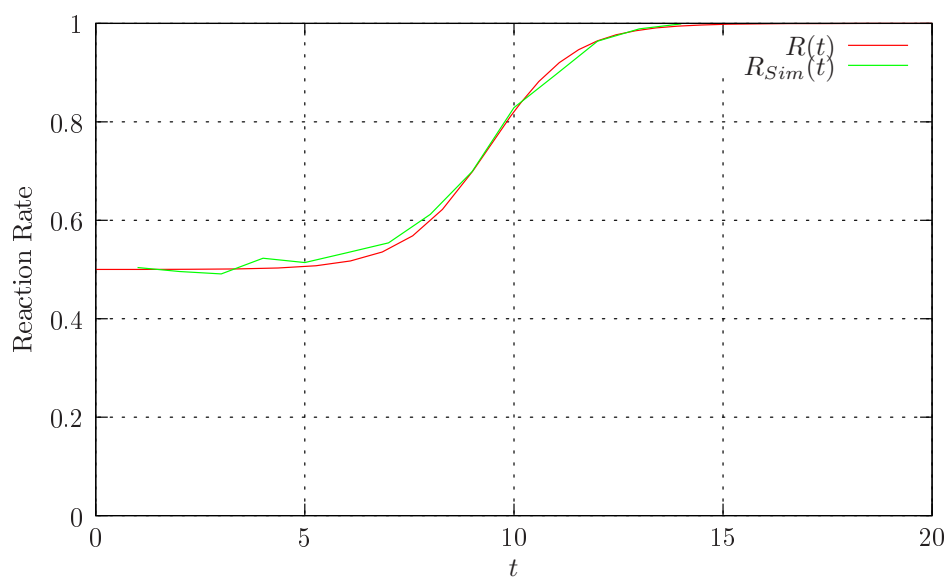


Figure 5.2: Survival of the Common Reaction Rate

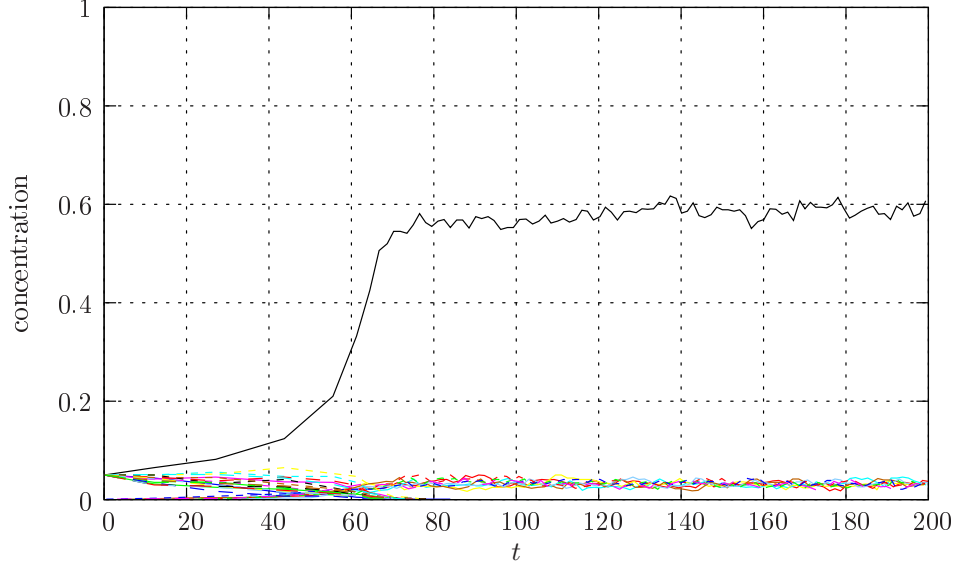


Figure 5.3: Survival Of The Common - 20 Species

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with equal concentrations (0.05) of twenty distinct molecular species of equal length $l = 10$,
- the per-bit mutation rate, v , was set to 0.05 and β , the proportion of mutations which are bit-flipping mutations (Section 4.4.5.2), set to 1.0.

Figure 5.3 shows that by $t = 20$, there was already one species whose concentration had begun to benefit from hyperbolic selection. This transition marks the beginning of a new epoch in the reactor, where there is a dominant (and by definition, *principal*) species with expected steady-state concentration of $(1 - V)$ (Section 5.1.1.1). Given that $v = 0.05$ in this experiment, and *all* molecules are of length $l = 10$ (since β , the proportion of mutations which are bit-flipping mutations (Section 4.4.5.2), was set to 1.0) :

$$\begin{aligned} (1 - V) &= (1 - v)^l \\ &\approx 0.6 \end{aligned}$$

This value of 0.6 is indeed the approximate steady state concentration of the dominant species during this epoch, with some superimposed stochastic fluctuation. The mutational load, V , consists *mostly* of the easily accessible mutations from the dominant species. This experiment was carried out multiple times using the same initial molecular population, but with different pseudo-random number generator seeds. In each case, the qualitative dynamics were identical, but the specific

molecular species which becomes dominant is unpredictable at the outset.

5.2.2 Takeover by Facultative Parasite

As we have seen in Section 5.2.1, the survival of the common effect is extremely robust. Given a population of diverse molecular species, which have no fitness differences between them, whichever species happens to first fluctuate to a concentration which is sufficiently higher than the others to be statistically effective, will then rapidly dominate the reactor; and, in the absence of mutation/replication error, will completely displace the other species. If, however, any of these other species could instantiate the class-6, facultative parasite system dynamics with the currently dominant species, then we would expect this parasitic species to take over the reactor. At low absolute numbers of molecules we expect that the early stochastic dynamics will significantly reduce the predictability of this result. At higher absolute numbers of molecules we expect that the ODE analysis presented in Section 4.5.4 will apply. To test this hypothesis, we devised the following set of experiments:

- M , the total number of molecules allowed in the reactor was set to 5×10^4 ,
- the reactor was seeded with a “host” molecule, species $x_1(0000000000)$, to concentration 9.99×10^{-1} (49,950 molecules) and a parasite molecule, species $x_2(00000000001)$ to concentration 1×10^{-3} (50 molecules).
- the per-bit mutation rate, v , was set to 0.

This run has been set up to verify that, at relatively high absolute numbers of molecules, the ODE analysis presented for class-6 dynamics in Section 4.5.4 applies. Figure 5.4 shows the outcome of this run with the curves aligned at $x_1 = x_1Sim = 0.5$. It is clear that the simulation results closely mirror the ODE predictions confirming that at relatively high absolute numbers of molecules, the ODE analysis can account for the behaviour of the system, even from relatively low initial parasite *concentration*. However, the reality of parasitic emergence and subsequent reactor takeover sees parasites emerge at *low* absolute numbers of molecules—in fact, the parasite lineage is *always* founded by a single mutant molecule. As such, we would expect that the early stochastic dynamics of these parasitic events mean that the ODE does not apply in all cases. The next set of experiments was setup to test the extent to which these early stochastic dynamics effect the behaviour.

Species	Mean	StdDev	Min	Max
x_1	3.12	6.94	0.001	107.4
x_2	53.9	16.6	27.7	156.6

Table 5.5: Time to Dominance Statistics

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with a “host” molecule, species x_1 (0000000000), to concentration 9.99×10^{-1} (999 molecules) and a parasite molecule, species x_2 (00000000001) to concentration 1×10^{-3} (now just a single molecule).
- the per-bit mutation rate, v , was set to 0.

The simulation was set to run until such time as one or other of the initialiser species dominated the reactor. From a total of 1×10^4 runs, the parasitic molecule reached dominance in just 262 (2.6%) of the runs. Table 5.5 shows the observed statistics of the time to dominance, separated according to which species was finally dominant, across these runs.

Figure 5.5 shows the outcome of a typical run in which parasitic takeover was observed overlayed with the ODE solution of class-6 facultative parasitism from Section 4.5.4. The graphs in this plot have been aligned at $x_1 = x_1Sim = 0.5$. The ODE solution showing global reaction rate is plotted alongside the measured global reaction rate for that run in Figure 5.6, which shows the graphs aligned at $t = 60$. These two figures generally align well with the ODE analysis. In the absence of mutation, however, no new molecular species can emerge, and a concentration of zero therefore becomes an absorbing state. Between $t = 0$ and $t = 20$ in Figure 5.5, it is clear that the parasitic species came close to extinction more than once. By contrast, the relatively high absolute number of molecules used for the simulation which produced Figure 5.4 meant that the simulation results were approximately continuous and deterministic even at the comparably low initial parasite concentration.

Figure 5.7 illustrates the results of a series of tests to determine the probability that a reactor will end up dominated by a parasite under the same conditions as above, by varying the starting number of molecules of the parasite. The plot shows three curves, representing varying values for M —the total number of molecules in the reactor. Each data-point in this graph is calculated from 1×10^3 runs of the model. Table 5.6 shows a more detailed view of the situation where there is a single instance of a parasite. Combining the data from this experiment with the lower-bound probability of producing a parasite on a single replication event presented in Section 5.1.3, we can estimate an approximate “waiting-time”, t_{wait} , for parasitic

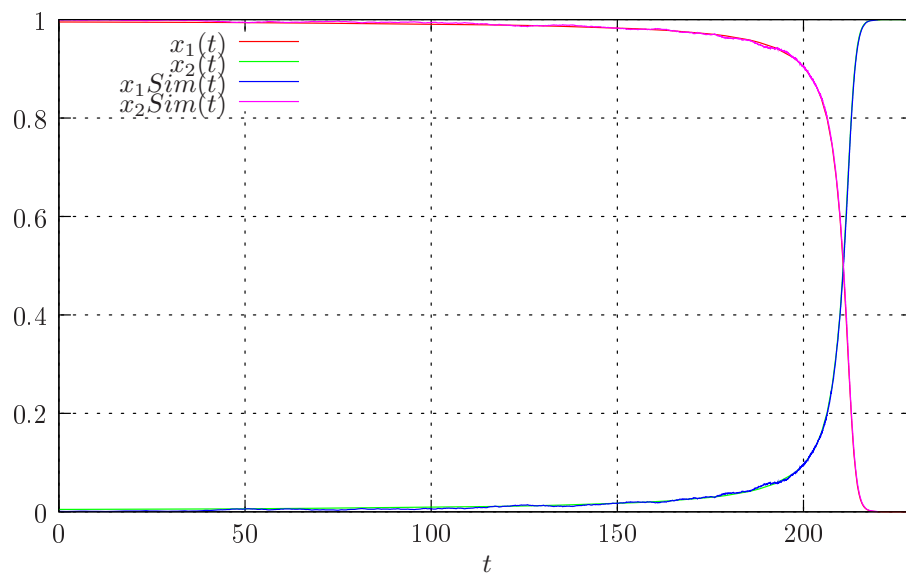


Figure 5.4: Parasitic Takeover—High Absolute Number of Molecules

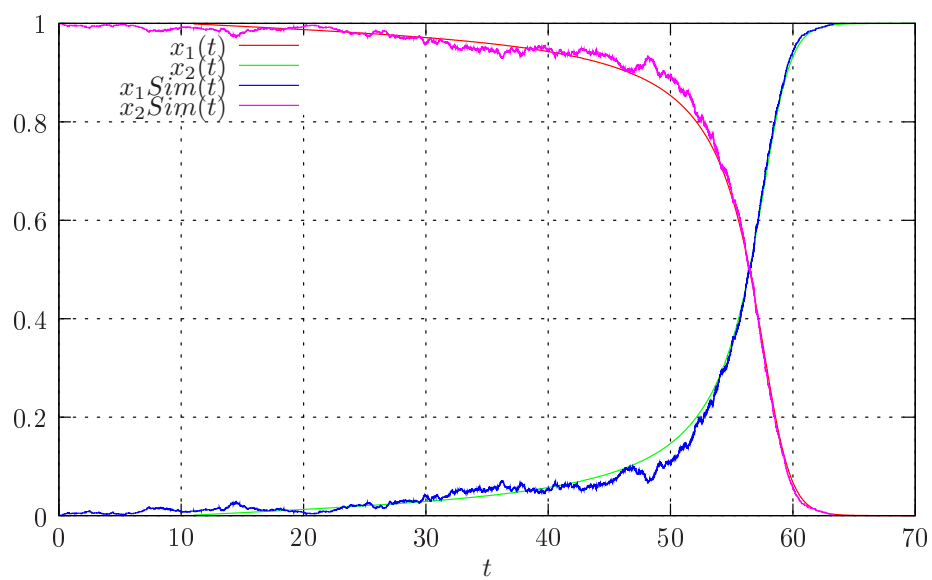


Figure 5.5: Parasitic Takeover—Low Absolute Number of Molecules

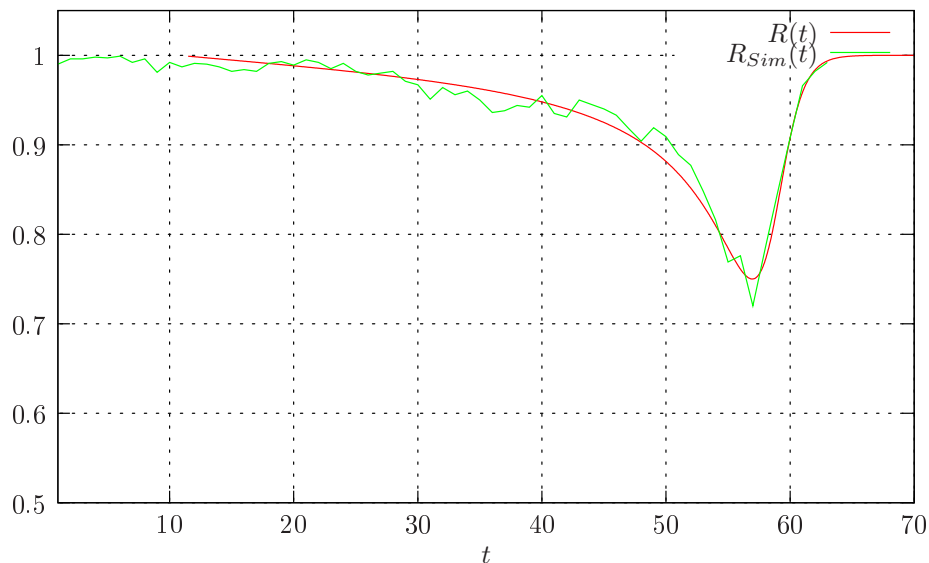


Figure 5.6: Takeover Reaction Rates—Low Absolute Number of Molecules

M	Parasite Concentration	$P(\text{takeover})$	t_{wait}
500	0.002	0.038	167
1000	0.001	0.0262	121
2000	0.0005	0.02	80

Table 5.6: $P(\text{takeoverBySingleParasite})$ (t_{wait} calculated for $v = 0.05$, $l = 10$ and $\beta = 0.99$ using equation 5.9)

takeovers in MCS-0:

$$t_{\text{wait}} = \frac{1}{[P(\text{parasiteOccurs}) \times P(\text{singleParasiteTakesOver})]} \quad (5.11)$$

5.2.3 Elongation Ratcheting

As we have seen so far, MCS-0 reactors will typically be dominated by a single self-replicating species until such time as a suitable class-6 parasite occurs within the population and displaces the previously dominant species. An important point to remember about the class-6 parasitic relationship in MCS-0 is that the parasite is necessarily *longer* than its host by at least one bit due to the substring matching rules that are applied (Section 5.1). It is due to this fact that we can reliably expect to see the monotonic increase of average molecular length in an MCS-0 reactor leading to a gradual reduction of reaction rate in the reactor. We should be able to plot this against both time and average molecular length to see that increasing molecular length does indeed correspond to decreasing reaction rate. The following

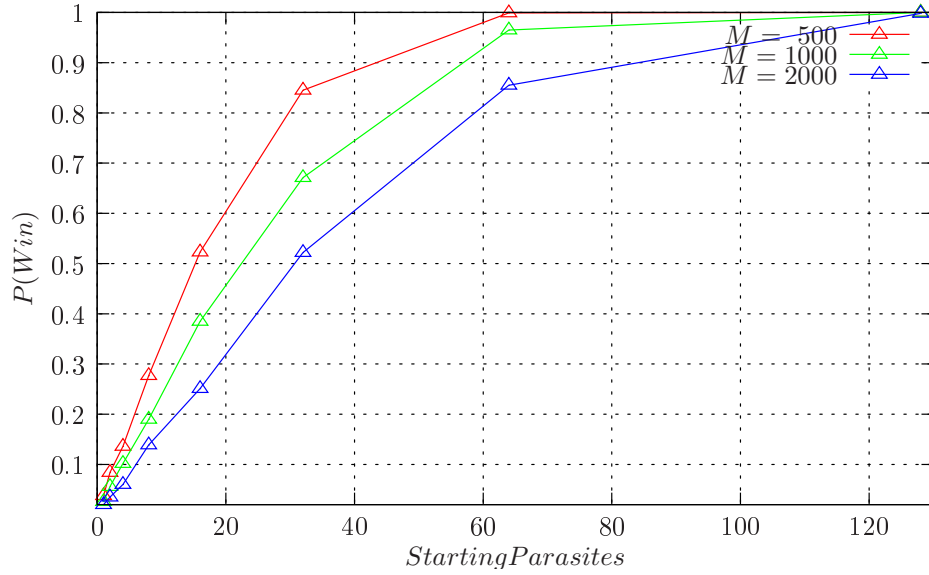


Figure 5.7: Probability of Parasite Win

experimental set-up was used to explore the phenomenon of elongation ratcheting:

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the per-bit mutation rate, v , was set to 0.05, and β , the proportion of mutations which are bit-flips (Section 4.4.5.2) set to 0.99.
- the reactor was seeded to capacity by making error prone copies of a seed molecule, species $x_1(00000)$ which results in a reactor initialised to a “normal” state for the appropriate per-bit mutation rate, v .

The simulation was run until such time as the length of the dominant molecule had reached 10 bits, double that of the seed species x_1 . Figure 5.8 shows the results of 10 runs of this simulation. Each of these 10 runs terminated within 500 timesteps, with an instantaneous reaction rate at termination less than half of the reaction rate at initialisation, demonstrating that the reaction rate does decrease as predicted. Figure 5.9 shows a more detailed plot of one of these runs. This time, the run was permitted to complete 1200 timesteps. The plot shows the concentrations of every molecular species that ever became the principal species in the reactor. It is clear from the plot that the length of these principal molecules is indeed increasing monotonically, and furthermore, it is clear that this length ratcheting coincides with a monotonic decrease in global reaction rate.

Equation 5.11 presented an approximation for average t_{wait} , the waiting time for parasitic takeover. Table 5.7 uses this equation to show the estimated average t_{wait} for each molecular species that became the dominant species in the experimental

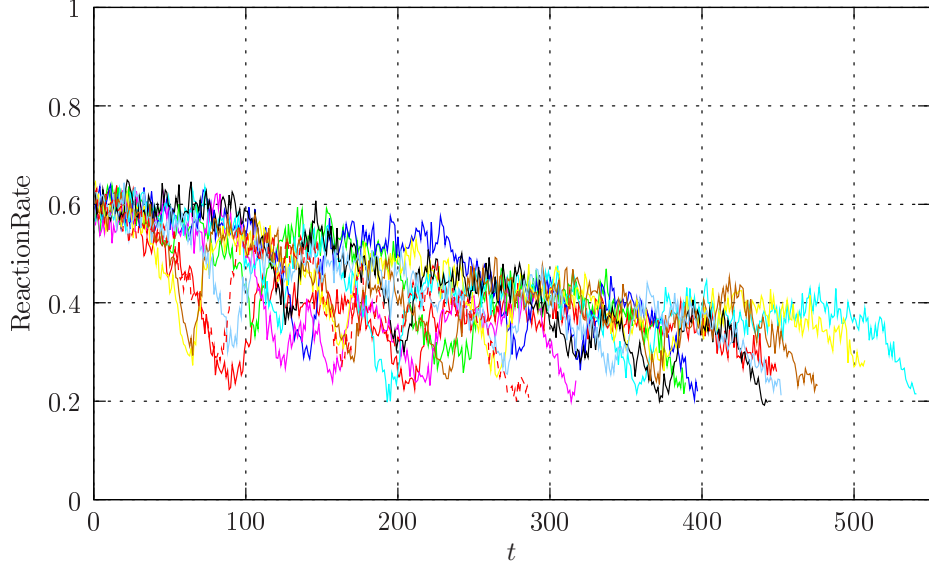


Figure 5.8: 10 Runs Overlaid Reaction Rates

Molecule	$P(\text{parasite})$	t_{wait}
00000	6.13×10^{-4}	62
000000	7.1×10^{-4}	53
0000000	7.22×10^{-4}	52
000000001	9.13×10^{-4}	41
0000000001	9.57×10^{-4}	39
00000000011	9.64×10^{-4}	39
000000000011	1.017×10^{-3}	37
1000000000011	4.16×10^{-4}	91
11000000000011	4.19×10^{-4}	91

Table 5.7: Approximate t_{wait} for Figure 5.9

run shown in Figure 5.9. The values for $P(\text{parasite})$ in each case were extracted by monte carlo simulation of 1×10^6 replication events in the same way as described earlier in Section 5.1.3. Of course, t_{wait} is essentially a random variable, and so using a sample size of one as we do here means that direct comparison of predictions and experimental results is impossible. Though we do not expect any precise quantitative correlation between the estimated average t_{wait} and the experimental results, we can make the qualitative correlation that the two longest lived species in Figure 5.9 also have the highest estimated t_{wait} values in Table 5.7.

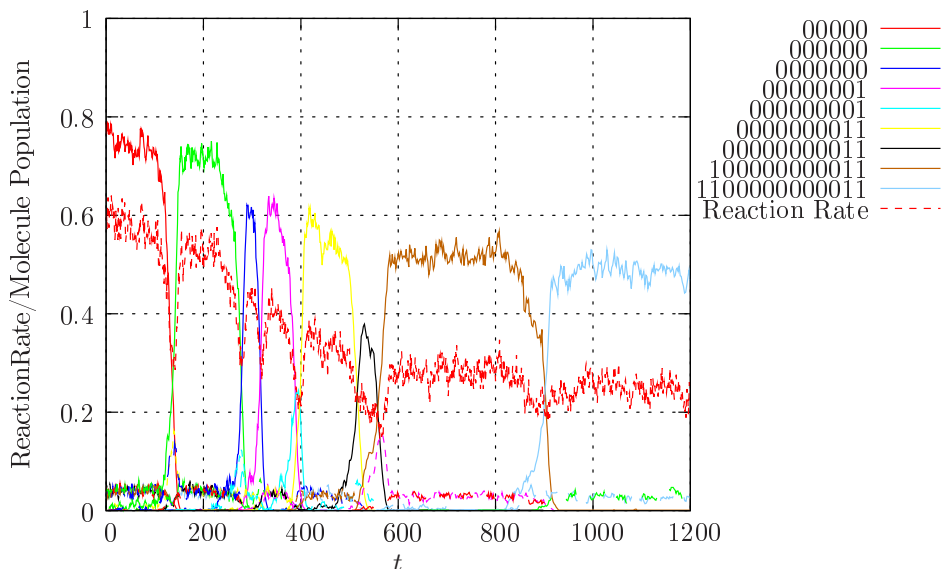


Figure 5.9: Length Ratcheting

5.2.4 Protocell Evolution in MCS-0 Systems

Now that we have quantified the effect of molecular parasitism on the global reaction rate, we shall explore the possible effects that hierarchical selection might have upon this. By constraining sub-populations of MCS-0 molecules to multiple distinct protocells, and by limiting resources at this new protocell level, we can initiate a process of natural selection between the protocells. Those protocell strains dominated by longer molecular species will have lower reaction rates, as seen in Section 5.2.3, which in turn will cause the protocell strains which contain molecular sub-populations with relatively longer molecular lengths to have lower reaction rates than those with relatively shorter average molecular lengths. Since we have limited resources at the protocell level, protocell strains with a relatively lower molecular reaction rate will be displaced by protocell strains which have higher molecular reaction rates. The experiments in Section 5.2.3 showed that molecular reaction rate—and consequently protocell growth rate—will be the key contributing factor in protocell gestation time.

In Section 5.1.4.4, a mechanism for classifying protocell strains based upon their principal molecular species was described. At that time, a pseudo-strain was defined for protocell strains which had molecular parasites of their principal molecular species—“mixed” protocells. The experimental results presented in this section will apply that classification algorithm to identify protocell strains.

The following experimental parameters were varied across the results presented below:

- $maxCells$ —the total number of protocells permitted in the simulation environment.
- S_{max} —Protocells split once their internal molecular population reaches S_{max} (Section 4.4.4.3)

Figure 5.10 shows a typical protocell run, set up as follows:

- $maxCells = 50$
- $S_{max} = 1500$

This example run shows that at these parameter settings the cellular population is almost saturated by “mixed” cells. This indicates that the protocell mutation rate is too high to allow coherent lineages to form—i.e., there is effectively an error catastrophe at the protocell level (cf. Maynard Smith, 1989, Chap. 2). Qualitatively, the ability of a pure lineage of protocells to maintain itself crucially depends on the production rate of molecular parasites. If even a single parasite is generated during the life time of a pure strain cell, then it is guaranteed that at least one of this cell’s daughter cells will be not be of the pure strain. The size at which protocells split, S_{max} , determines the number of molecular replications that take place during the lifetime of that protocell, and therefore has an effect on the protocell mutation rate. Clearly then, the overall rates of protocell level mutation can be substantially reduced by reducing the number of molecular replications before each fission event—i.e., by reducing the S_{max} parameter. It was found that reducing S_{max} to 150 (and correspondingly increasing $maxCells$ to 500 to maintain the same overall number of molecules) resolved this problem (Figure 5.11). No finer grained exploration was carried out on this particular issue.

Now that the strain classification parameters have been decided, we can begin to explore the true effects of hierarchical selection in MCS-0 protocells. This section began with the hypothesis that “protocells which suffer from a reduced reaction rate will be displaced by protocells which have higher reaction rates”, reaction rates being variable due to the length of the dominant molecular species in the protocell strain. The experiment to test this hypothesis was initialised as follows:

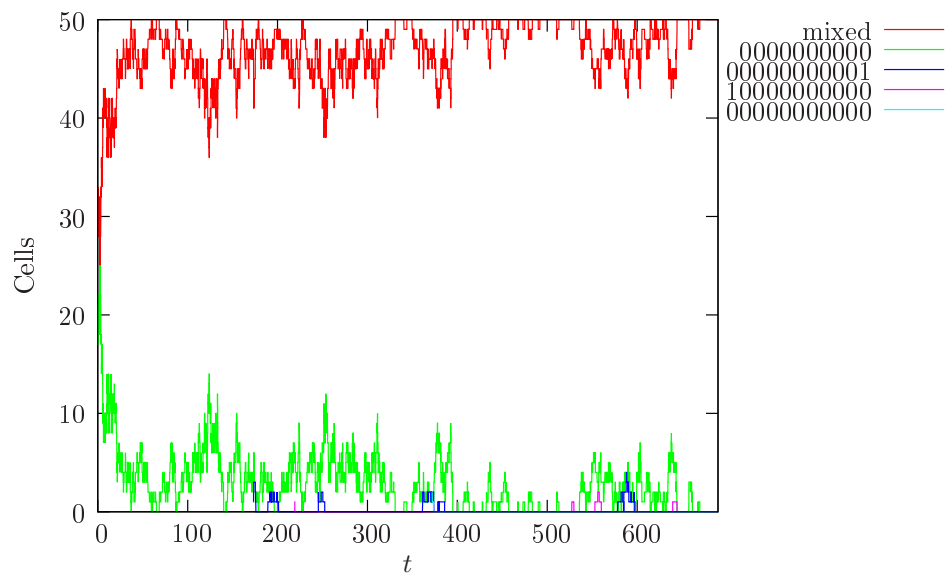


Figure 5.10: $S_{max} = 1500; maxCells = 50$

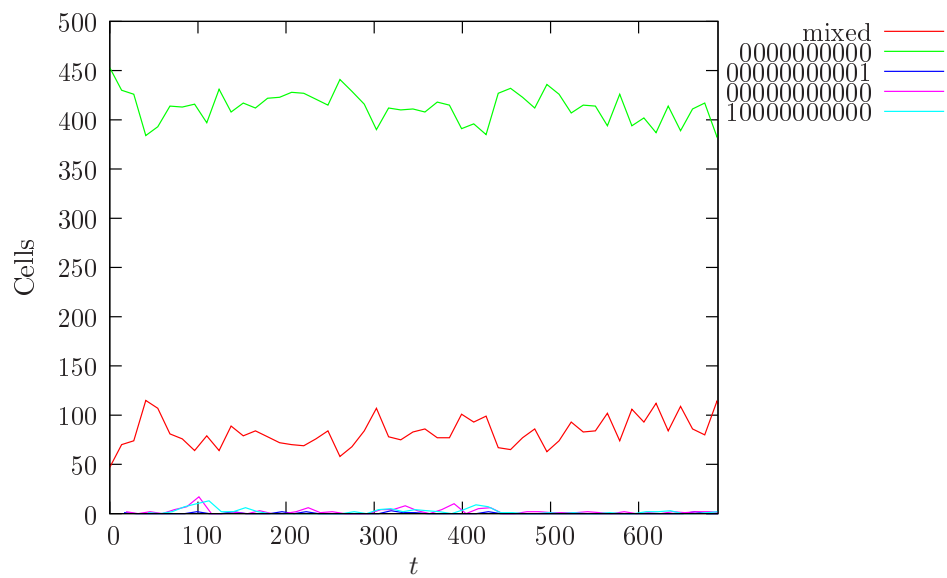


Figure 5.11: $S_{max} = 150; maxCells = 500$

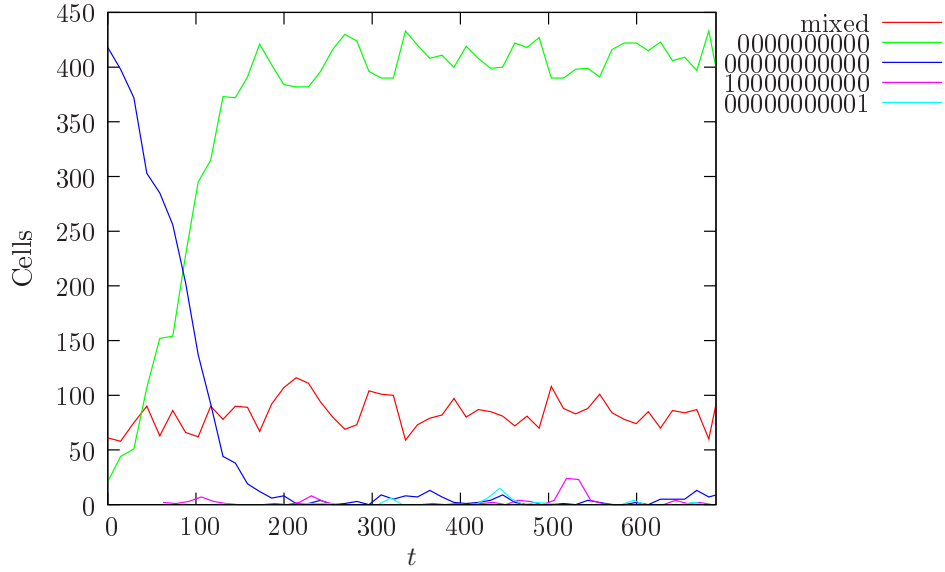


Figure 5.12: $S_{max} = 150; maxCells = 500$

- $v = 0.05; \beta = 0.99$
- $maxCells = 500; S_{max} = 150$
- Two initialiser molecular species of length l and $l + 1$ respectively are chosen.
- $(0.9)maxCells = 450$ are initialised with the $l + 1$ initialiser
- $(0.1)maxCells = 50$ are initialised with the l initialiser.
- The process for initialising individual protocells ensures that at $t = 0$, the protocell population is in a normalised state. This means that protocell sizes, measured as total number of molecules, will be uniformly random, and that within protocells, a suitable mutational load is present.
- Repeated molecular interactions, as described in Section 4.4.5, are then carried out.

Figure 5.12 shows the results of a typical run of this experiment. The protocell strain seeded with the l initialiser molecule rapidly displaces the strain seeded with the $l + 1$ initialiser. Periodically throughout the experiment, new cell strains are founded due to parasitism at the molecular level within the main strain of protocells, but these new strains can *never* get an established toe-hold in the protocell population as their principal molecular species is *longer* than that of the predominant cell strain (the cells initialised with the molecular species of length l). None of these

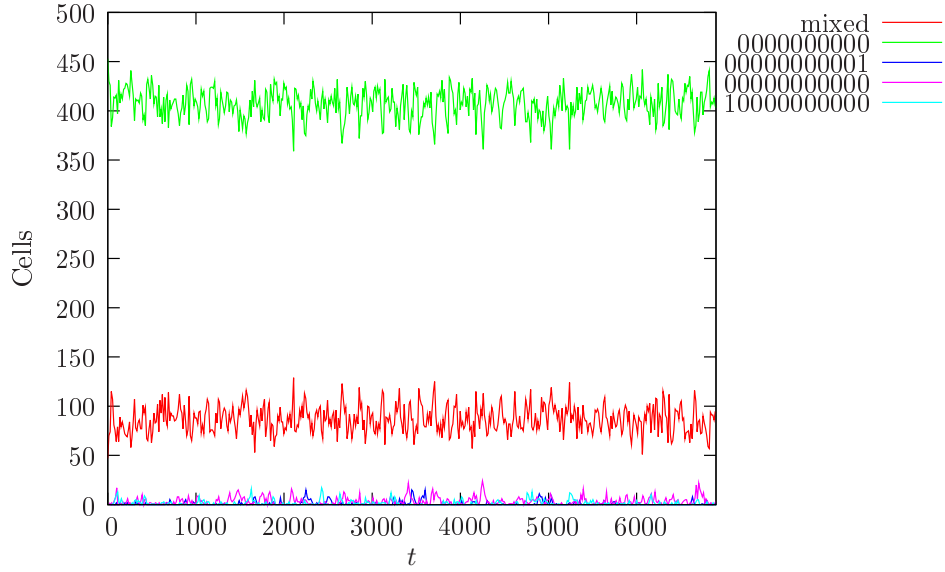


Figure 5.13: $S_{max} = 150; maxCells = 500$

new strains can ever increase significantly in the protocell population as they are constantly selected against. They persist, however, due to continuing mutational flow at the protocellular level to these mutant strains.

Regardless of the choice of l in this experiment, the l strain will *always* displace the $l+1$ strain, since shorter molecules have a lower per-string mutation rate implying a higher protocellular growth rate. However, no matter how long the experiment will run for, it will never be possible for a new strain of l , or shorter to be founded. This is because:

1. The founding of a new protocell strain corresponds to a parasitic takeover at the molecular level within a protocell, and,
2. it is *impossible* for a molecule to be a parasite of a molecule that is longer than it ($lengthOfParasite < l$) or even the same length as it ($lengthOfParasite = l$).

It will be possible, however, to witness new strains, of $l+1$ or greater, being founded continuously. The combination of these factors gives rise to the phenomenon of “selectional stalemate” between selectional forces at the molecular and protocellular level acting in opposite directions, essentially balancing each other out. Figure 5.13 shows the outcome of a much longer run, demonstrating that the “stalemate” phenomenon is indeed robust over extended periods of time.

5.3 Conclusion

It is clear from the experiments presented in this chapter that selection at both the molecular level and the protocell level can be coupled in such a way as to allow the selection of a characteristic protocellular trait (the species of the dominant molecule). This hierarchical selection, in combination with the specific MCS-0 molecular interaction rules used gave rise to the distinctive, and original phenomenon of “selectional stalemate” that arose in the later experiments of this chapter. This effect clearly demonstrates the potential for distinctive evolutionary phenomena to arise through the interaction between different levels of selection.

In the next chapter, I will describe a revised version of MCS-0, MCS-1, which includes modifications to the molecular interaction scheme aimed at opening up a richer variety of molecular dynamics. In this way, it should be possible for MCS protocells to display ongoing evolution.

Chapter 6

MCS-1

This chapter presents the results of the experiment set: MCS-1. In the previous chapter, we saw that the proof of concept MCS-0 system provided the basis for an artificial chemistry in which to carry out hierarchical selection experiments. However, the potential for evolutionary growth of complexity in the MCS-0 system was constrained by the (novel) phenomenon of “selectional stalemate”. It was concluded that this minimal MCS-0 system would need some carefully focused enrichment to now break selectional stalemate while at the same time retaining *some* characteristics of the molecular level chemistry that could be used as a target for selection in protocell level experiments.

The work presented in this chapter aims to expand upon the MCS-0 by modifying some of the implementation details of the molecular chemistry, essentially enriching the reaction rules. To use terminology introduced by Altmeyer et al. (2004) (Section 5.1), MCS-0 could be described as a system of “complete molecular recognition”, and the MCS-1 that will be presented here will implement “partial molecular recognition”. As we will see, these enrichments do indeed disrupt the selectional stalemate as desired.

6.1 Introduction

The MCS as presented so far has displayed some interesting phenomena both in terms of molecular selectional dynamics and hierarchical evolution. We have seen that it is possible to couple two layers of selection (molecular selection and protocell selection) in a way that permits the protocell level of selection to have an indirect effect on the molecular level. In other words, we have been able to select for preferable molecular sub-populations without explicitly examining these sub-populations in detail. This hierarchical selection along with the particular molecular interaction scheme that has been used so far has given rise to a robust “selectional stalemate”

phenomenon due to the opposition of the two levels of selection. Essentially, the MCS-0 molecular binding rules enforced a scheme which resembles the “complete molecular recognition” described by Altmeyer et al.. In this chapter, MCS-1 is presented as a revised version of MCS-0 which aims to open up new evolutionary pathways for the MCS platform to explore. As we will see, however, the initial modifications presented here will have a dramatic effect on the underlying molecular dynamics to the point that individual replicator species in the MCS-1 system lose the ability to maintain a level of dominance in a population. The early analysis presented in Chapter 5 relied foremost on there being a quasi-stable dominant replicator according to whose interactions we can approximate the molecular dynamics in a fixed size reactor, or in the case of the hierarchical experiments, a growing cell lineage. If we lose the ability to approximate that behaviour, we lose the protocell inheritance mechanism that was described in Section 5.1.4.3; and if there is no mechanism for protocell inheritance, then there can be no basis for evolution at the protocell level. This chapter focuses then on firstly introducing some modifications to MCS-0—to yield the modified system called MCS-1—and subsequently on characterising the problems that arise and proposing a remedy for them. Particular attention is given to maintaining the protocell inheritance mechanism that was derived in Chapter 5. One of the key contributions of this chapter is an analysis of protocell “fitness” in terms of gestation time and protocell mutation rate.

6.2 MCS-1: The Modifications

In MCS-0, there was little distinction between informational and enzymatic functions of the molecules, other than in name. Recall that the enzymatic function of MCS-0 molecules was to bind with molecules that were a bitwise *super*-string of themselves (Section 5.1.1.1) and then to make a bit-wise, error-prone copy of the bound molecule (Section 4.4.5.2). This simple dynamic severely limited the potential for interaction between any two molecules. In Section 4.5.4, a full characterisation of all potential “binary replicase interaction systems” was presented. Depending on the choice of molecular binding rules, some or all of those possible binary replicase interaction systems may be impossible to instantiate. In other words, particular pairs of molecules may or may not be able to interact with each other for different choices of molecular binding rules. In Section 4.5.3, “self”-replication was identified as a special case, since there is only one species of replicase present in a “self”-replication reaction. Depending on the molecular binding rules used, it might not be possible to instantiate some of the generic set of binary replicase interaction systems presented in Section 4.5.4. Table 6.1 shows a summary of the systems that can be

Class	Possible?
0	×
1	×
2	×
3	×
4	×
5	×
6	✓
7	×
8	×
9	✓

Table 6.1: Binary Replicase Interaction Systems Possible in MCS-0

instantiated for the binding rules used in MCS-0. From the analysis presented in Section 4.5.4, it can be seen that the class-9 behavior is responsible for the “survival of the common” phenomenon seen in the experiments of Section 5.2.1 and also the periods of dominance by a single self-replicase in the experiments of Sections 5.2.2 and 5.2.3. The potential for facultative parasitism as seen in the experiment in Section 5.2.2 is also encapsulated in Table 6.1. Since each and every molecule is a *self*-replicase, any of the reaction classes which involve molecules which are not capable of self-replication can be immediately excluded. It is clear from Table 6.1 that the MCS-0 molecular interaction ruleset is quite restrictive since eight of the ten binary replicase interaction networks cannot be instantiated.

With MCS-1, a more complicated scheme for transforming a molecule from its informational structure into an enzymatically active structure will be introduced. This transformation mechanism will take the form of molecular folding, and will change MCS from a system of *almost* “complete molecular recognition” to a system more like that of “partial molecular recognition”. The transformation mechanism has been chosen arbitrarily and takes inspiration from the ribozymes of the RNA-World hypothesis which have both informational and enzymatic roles, and which undergo spontaneous folding into characteristic hair-pin loops (Gilbert, 1986). MCS molecules will be said to be in their enzymatically active form when they are folded, and in their informational form when they are not folded, as with MCS-0. The following two sub-sections will describe this folding mechanism, but as we will see in Section 6.3, there is a need for some further adjustment of the MCS-1 system. These further modifications will be presented in Section 6.3.1.1.

6.2.1 Molecular Folding

In order to establish a system of “partial molecular recognition”, we need to introduce a transformation mechanism which will be used to convert the molecules into their enzymatically active form. This transformation mechanism is termed “folding” from this point onwards—an analogy inspired by the ribozymes of the RNA world (Kruger et al., 1982; Gilbert, 1986; Joyce, 1991) and protein folding into active conformations. Molecules in their un-folded, informational form are said to be represented by their primary structure whilst molecules in their enzymatic, “folded” form are said to be represented by their secondary structure. We use the idea of folding here to highlight the fact that there is some indirect relationship between primary and secondary strings which can be algorithmically generated.

Folded MCS molecules are represented as strings over the alphabet $[L, H]$. The folding mechanism is arbitrary and follows a pre-defined folding scheme (for example, see Table 6.2). This scheme determines the various transformations between one or more monomers of the primary structure and their corresponding tokens from the secondary alphabet. The primary structure of a molecule is processed from left to right and the transformations indicated by the folding scheme are applied. In the event that there remain bits which cannot be processed fully, these bits are disregarded for the purposes of folding, but they are not removed from the primary structure. In other words, they have no effect on the folding process.

Molecular interactions in MCS are always between two *individual* molecules, though these molecules may of course be of the same *species* (Section 4.5.3). One of the molecules in any given reaction will be represented by its primary structure, and the other will be represented by its secondary structure. In the next sub-section, we will see how to determine whether two molecules may “bind” or not—the updated molecular binding rules. It is the interaction between the folding scheme and the molecular binding rules that gives rise to a particular set of possible molecular dynamics.

As previously mentioned, the particular folding scheme used must be pre-defined. Under the conditions of:

1. a two-bit codon length, $[00, 01, 10, 11]$,
2. a two letter secondary alphabet, $[L, H]$, where each letter must appear in the folding scheme at least once,

there are $(2^4 - 2) = 14$ possible folding-schemes for MCS-1. Tables 6.2, 6.3 and 6.4 show some example folding schemes. In the absence of any theoretical guidance for the choice of folding scheme, the folding scheme presented in Table 6.2 will be the

one selected for all experimentation. We conjecture that the specific phenomena presented in this thesis do not depend critically on this choice of folding scheme¹.

Unfolded (primary) String	Folded (secondary) String
00	L
01	L
10	H
11	H

Table 6.2: Example MCS Folding Scheme #1

Unfolded (primary) String	Folded (secondary) String
00	L
01	H
10	H
11	H

Table 6.3: Example MCS Folding Scheme #2

Unfolded (primary) String	Folded (secondary) String
00	L
01	L
10	L
11	H

Table 6.4: Example MCS Folding Scheme #3

6.2.2 Molecular Binding

When two molecules collide (Section 4.4.5), they may or may not be capable of binding to one another. As with MCS-0, this binding process begins by treating

¹If the secondary alphabet were expanded however, then significant further research would be required to understand the molecular dynamics fully. For example, Holland’s Classifier Systems use a further two symbols: $[\#,:]$ to represent a wild-card and syntactic separator respectively (Holland and Reitman, 1977). Some preliminary investigations into using wild-card operators in MCS-1 revealed complex behaviours that were difficult to relate to the MCS-0 dynamics. Accordingly it was decided to focus here on the simpler case involving only the literal symbols of the classifier formalism.

Secondary String Character	Matches
0	0
1	1

Table 6.5: MCS-0 Pattern Matching Rules

Secondary String Character	Matches
L	0
H	1

Table 6.6: Pattern Matching Rules

one of the molecules as a substrate, and the other as a catalyst, or enzyme. MCS-1 differs by firstly “folding” the enzyme molecule into its “active” form (secondary structure per Section 6.2.1). This “activated enzyme” is then processed together with the substrate molecule to determine if a binding has occurred.

In MCS-0, the possibility of binding between a catalyst and a substrate was determined by a process of bitwise sub-string pattern matching. Table 6.5 represents the rules used for that pattern matching algorithm. The molecule which acts as the enzyme is assumed to have been folded into its secondary structure, which for MCS-0 was merely an identity transformation. If the string represented by the folded enzymatic structure of the catalyst molecule is found within the primary structure of the substrate, then the molecules are deemed to be able to bind. This binding is a short-lived one however and will spontaneously decay once the reaction processing (which is outlined in Sections 4.4.5.2 and 4.4.4) has completed.

In MCS-0, the secondary string alphabet was $[0, 1]$ whilst in MCS-1 we have changed the alphabet to $[L, H]$. Accordingly, the molecular binding rules will be updated for MCS-1 as described in Table 6.6. An updated version of Table 5.1, which presented an example of the molecular binding rules for MCS-0 is shown in Table 6.7. This updated version represents the molecular folding described in Section 6.2.1 via the transformations laid out in Table 6.2.

6.2.2.1 Self-Replicases in MCS-1

In MCS-0, every molecular species was capable of acting as a self-replicase—every string is, by definition, a substring of itself. In MCS-1 however, the new molecular matching rules introduce the possibility that the secondary structure of an instance of a particular molecular species may not be capable of binding to the primary

Catalyst (<i>primary</i>)	11001001	Catalyst (<i>primary</i>)	11111010
Catalyst (<i>secondary</i>)	HLHL	Catalyst (<i>secondary</i>)	HHHH
Substrate	00 10 1000	Substrate	00101000
Match?	yes	Match?	no

Catalyst (<i>primary</i>)	00101000
Catalyst (<i>secondary</i>)	LHHL
Substrate	00101000
Match?	no

Table 6.7: Examples of Molecular Matching (MCS-1)

structure of another instance of that same species. In other words, not all MCS-1 species are capable of acting as *self*-replicases. In the experiments carried out for MCS-1, it was necessary to have a method of producing random self-replicase species to act as seeds for the initialisation of reactors. The algorithm used to produce these random self-replicases was a brute force one. Essentially, random molecules of a given bit-string length were produced and then checked to see if they were self-replicases. This process was repeated until a self-replicase was found.

6.2.3 Expanded Matching Capabilities

The modifications to MCS that have been presented in Sections 6.2.1 and 6.2.2 will, at the very least, allow for a larger repertoire of potentially successful molecular interactions as summarized in Table 6.8. In particular, it should be noted that the new possibility of partial molecular recognition, which has been enabled by the new folding mechanism, will allow a species to act as a catalyst for other species that are shorter than itself in primary form and conversely, allow species to parasitise other species that are longer than themselves—a situation which was impossible in MCS-0. The experiments presented in the remainder of this chapter will show that this modification to MCS does indeed make a significant difference in the potential evolutionary outcomes for both the molecular level and the protocell level of selection, most notably in the elimination of the selectional stalemate phenomenon that was characteristic of MCS-0.

6.3 MCS-1: The Experiments

In light of the modifications that have been presented in Section 6.2, it is necessary to revisit some of the early experiments of Chapter 5 in order to determine precisely where the effects of these changes will manifest. It is also important to remember at this stage that while we have only modified the interactions of the molecules and have not explicitly modified the implementation of the protocell level dynamics of

Class	Possible?
0	✓
1	✓
2	✓
3	✓
4	✓
5	✓
6	✓
7	✓
8	✓
9	✓

Table 6.8: Binary Replicase Interaction Systems Possible in MCS-1

the model, there may well be noticeable changes at the protocell level nonetheless. In fact, the changes to the protocell level evolutionary dynamics are what we are most interested in seeing, since it is hierarchical selection that is the main subject matter of this thesis. The qualitative behaviours of MCS-0 that we want to maintain in MCS-1 are the following:

1. “stable” periods of domination by a single self-replicase species
2. occasional displacement events by class-6 facultative parasites

since these essentially make up the protocell inheritance mechanism that we will use for hierarchical experiments. At the same time, the hope is that the modifications presented in the early sections of this chapter will be able to eliminate the selectional stalemate which was present in the hierarchical experiments of MCS-0 and blocked evolutionary progress in the system.

6.3.1 Collapse of Seed Self-Replicase Species

The first experiment required is one that will verify the molecular level dynamics, in the absence of the higher protocell level of selection. We seek to reaffirm our primary assumption from Chapter 5—that reactors will be dominated by a single self-replicase until such time as another, fitter, self-replicase displaces it and takes over the reactor. In essence, what is called for at this point is to re-run a combination of the experiments presented in sub-sections 5.2.1 and 5.2.2 and take note of the differences, if any, between the MCS-0 results and the MCS-1 results.

This experiment will be initialized with appropriate parameters to make it comparable to the MCS-0 experiment presented in Section 5.2.2. The reactor will be seeded with 1×10^3 molecules, representing a quasi-steady-state² sub-population of

²In other words, there is an adjusted mutational load, as described in Section 5.2.2

a particular self-replicase. As with the earlier MCS-0 experiments, this experiment uses a fixed size reactor with a maximum capacity of 1×10^3 molecules (Section 4.4.2). The key difference between this experiment and the one presented in Section 5.2.2 is that the molecular interactions respect the modifications detailed above in Section 6.2.

The experiment was initialised with the following conditions:

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with a normalised population of molecules centered on the self-replicase species: $x_1 = 0001001100$,
- the per-bit mutation rate, v , was set to 0.05, and β , the proportion of mutations which are bit-flips (Section 4.4.5.2) set to 0.99.

Figure 6.1 shows the fate of the concentration of the seed self-replicase, x_1 , during a single run of this experiment. It is clear that the seed species, far from maintaining a period of dominance, rapidly declines in concentration. Of course, this *might* be due to the chance early presence of an appropriate class-6 facultative parasite. If this were indeed true, then we would expect that re-initializing the experiment with different random number generator seeds and a different seed molecular species might result in a set of experiments which support both of the conditions in the MCS-0/MCS-1 continuity list presented above—periods of relative stability punctuated by parasitic displacements.

Figure 6.2 shows a composite of 10 runs which were set up as follows:

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with a normalised population of molecules centered on a randomly chosen self-replicase species of length $l = 10$ which was capable of self-replication (derived from the algorithm described in Section 6.2.2.1),
- the per-bit mutation rate, v , was set to 0.05, and β , the proportion of mutations which are bit-flips (Section 4.4.5.2) set to 0.99.

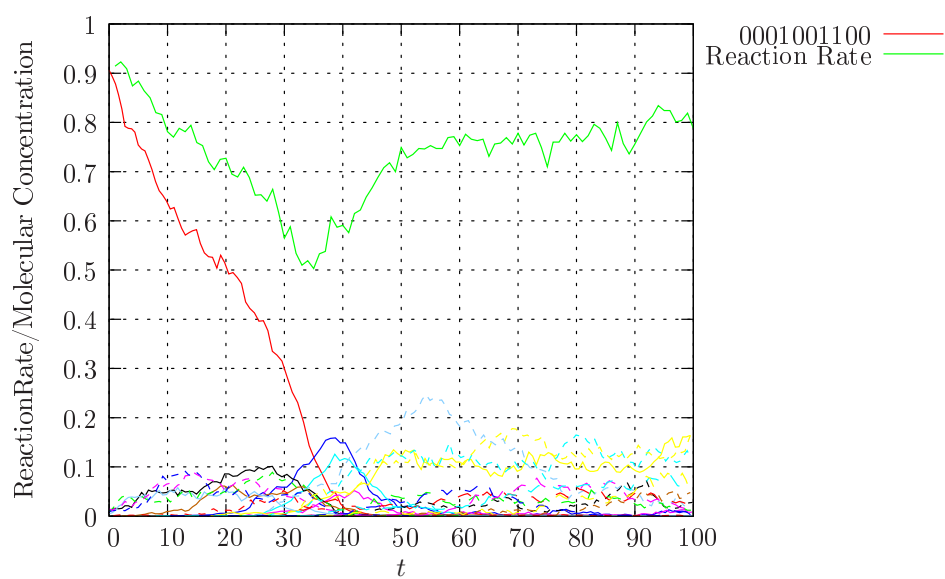


Figure 6.1: Collapse of Seed Self-Replicase

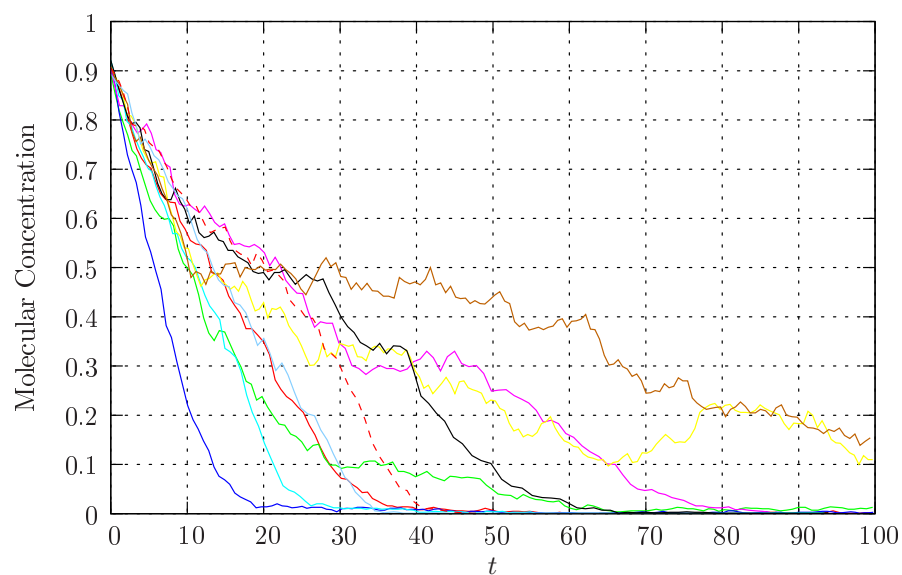


Figure 6.2: Collapse of Seed Self-Replicase (10 Runs) Molecular Concentrations

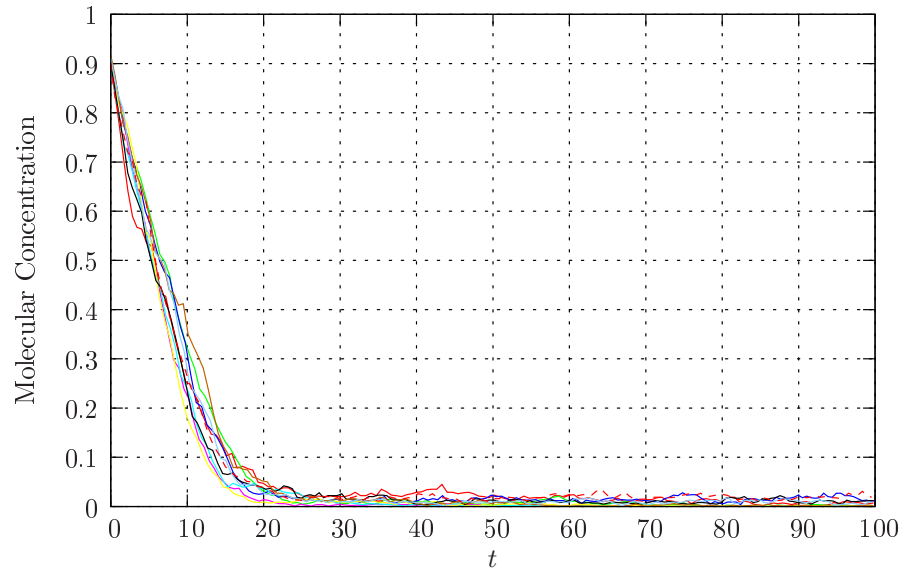


Figure 6.3: Collapse of Seed Self-Replicase (10 Runs) Molecular Concentrations (left outlier)

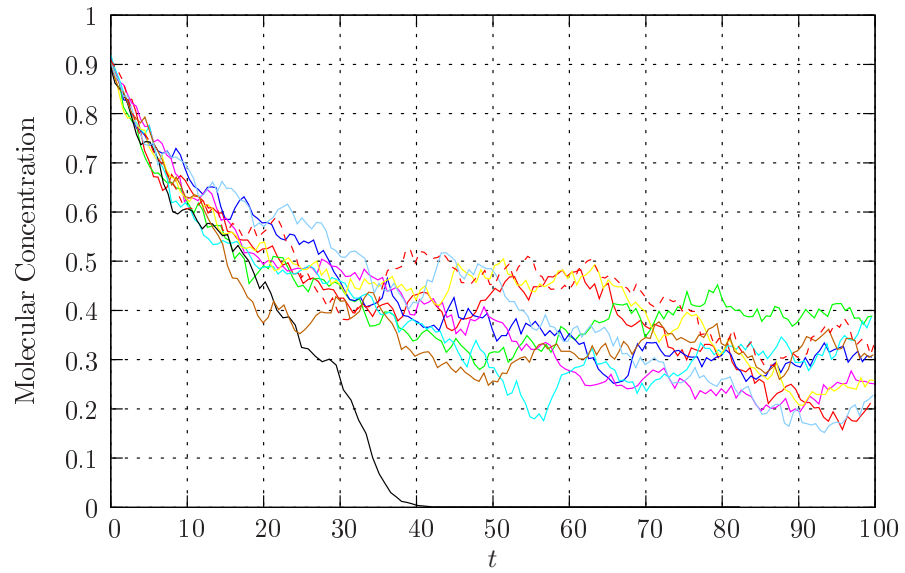


Figure 6.4: Collapse of Seed Self-Replicase (10 Runs) Molecular Concentrations (right outlier)

The plot in Figure 6.2 shows that in all cases, the concentration of the seed self-replicase decays rapidly suggesting that MCS-1 is *not* behaving as predicted, since it is unlikely that 10 independent runs with different seed molecules should show take over by parasites in approximately the same amount of time. Figures 6.3 and 6.4 show the results of a further 10 runs each for the species which collapsed most quickly and most slowly, respectively in Figure 6.2. These results show consistency for individual initialiser species, suggesting a possible systemic reason for the collapse of the seed self-replicase.

At this stage, let us refer back to the set of binary replicase interaction systems which were presented in Section 4.5.4 which can be used to characterise and label the interaction dynamics between two molecular species in MCS. It is possible to analyse any run of the MCS system and examine the contents of the reactor by considering the pairwise “class” of each molecular species and the principal species. Such an analysis should show the presence of parasites of the principal species, and if not, should at least shed some light on the reason for the collapse of the population of the principal, seed, self-replicase. If MCS-1 behaves in a qualitatively similar way to MCS-0, then we should expect to see a sharp increase in class-6 facultative parasites mirroring the decline in the concentration of the seed self-replicase.

Applying this diagnostic technique to the MCS-1 experiment shown in Figure 6.1, however, reveals that a fundamentally different process underlies MCS-1 as seen in Figure 6.5. Rather than the expected prevalence of species which are pairwise class-6 to the principal species, we can see a steady increase in species which are either pairwise class-1 (facultative mutualists) or pairwise class-4 (obligate parasites) to the principal species. It is important to note that the rise in pairwise class-6 molecules from circa $t = 25$ onwards is not a significant contributing factor in the decline of the initial seed. By the time this sub-population of species which are class-6 relative to the principal species begins to rise rapidly, the population of the principal species has already been steadily decaying for circa 30 timesteps.

This result was not anticipated in advance. The original analysis of the interaction classes which was presented in Section 4.5.4 suggested that the only mechanism for selective takeover of a reactor is that of class-6 facultative parasitism, and that under no circumstances would one expect to see the systematic displacement of a self-replicase by molecular species which shared any other pairwise relationship to it. Nevertheless, the results presented in Figure 6.5 are typical of MCS-1.

Closer inspection of Figure 6.1 reveals another unanticipated fact about these MCS-1 experiments. Apart from the initial relaxation period, no single molecular

species attains a concentration higher than about 0.23. This implies that the displacement events we see are not due to the presence of a single invading species but rather due to the collective action of a cloud of distinct molecular species. However, given the detailed analytical results presented in Section 4.5.4, the hypothesis that a dominant molecular species can be displaced by anything other than the selective pressure due to class-6 facultative parasitism requires a deeper explanation and re-examination.

The explanation proposed here is as follows:

1. Mutational flows are generally *asymmetric* in systems of modular-replicators such as MCS. In other words, there are numerous ways of producing a mutant *from* a given master sequence, but there is only one pathway (the reverse of the “inbound” path) to directly return from one of these mutants *to* that master sequence. If mutational flows were the only relevant dynamic then this asymmetry would have a diluting effect on the concentration of the master sequence of a given reactor since the expected rate of back-mutation is negligible. However, asymmetric mutation flows alone cannot explain the observed behaviour, as this was also true in MCS-0, but was there outweighed by the underlying, quasi-deterministic, survival-of-the-common dynamic.
2. Given the changes to the reaction rules presented in Section 6.2, we can revisit the analysis presented in Section 4.5.4:
 - (a) The system dynamics due to the pairwise classes 0, 3 and 8 will not have a significant influence on a reactor since these classes only arise when neither species is a self-replicase.
 - (b) Instantiations of classes 2, 5, 7 and 9 will be actively selected against relative to the principal species, though due to the constant presence of mutation, such species will never be completely eliminated from the reactor (Section 5.1.1.1)
 - (c) Facultative parasites, species which are pairwise class-6, will be actively selected for, as before.
 - (d) Mutations into pairwise class-1 species are selectively neutral at the molecular level, which should give rise to a “family” of mutually pairwise class-1 species. That is, per point 1 above, a master species would indeed tend to diversify across this family of class-1 mutant species. Such a family could be seen as a single “quasi-species” which might still collectively behave in the same way as a single dominant species. An increase in class-1 mutants is indeed clearly visible in the early phase of Figure 6.5. However, in the more likely case that at least some of the family members

are not “fully connected” (in the sense of class-1 pairwise dynamics) to *all* other members, then the overall reaction rate will begin to fall and the equivalence to a single dominant species system would begin to fade.

- (e) As in the previous point, mutations into pairwise class-4 species are also selectively neutral at the molecular level. However, there is one key difference. As the principal species diversifies across these pairwise class-4 species, the overall reaction rate will unconditionally decrease as the concentration of class-4 mutants increases. An increase in class-4 mutants is also clearly visible in the early phase of Figure 6.5.

Given that the goal of this research is to explore the effects of hierarchical selection in artificial protocell systems, these molecular level dynamics will need to be changed before we can reliably use the replicase properties of the molecular level as an inheritance mechanism for the artificial protocells which we want to build. In the following section, one such approach to addressing these issues is presented, though it is worth clarifying at this point that the MCS is so abstract in comparison to realistic natural systems that the problems identified here are primarily relevant to the design of *computational* evolutionary systems rather than necessarily being systemic problems for a properly biological “origin of life” hypothesis, or intrinsic barriers to the evolution of artificial, physical, protocells.

6.3.1.1 Further Modifications

The results of the previous experiments show that the adoption of the new folding scheme gives rise to some pairwise reaction classes which can have a drastically detrimental effect on the overall molecular dynamics—so much so that we no longer have an effective mechanism for protocell level heritability. At this point, there are at least two distinct pathways that might be followed to solve this problem:

1. implement some mechanism giving rise to an intrinsic fitness difference between molecular species, thus allowing conventional dominance of a fitter species over its spectrum of immediate mutants, or
2. impose restrictions on the nature of the reactions that are allowed to occur, in order to limit dilution of a principal self-replicase species by its immediate class-1 and class-4 mutants.

An important design principle for the MCS has been to keep the molecular model as simple as possible. The changes to MCS so far have introduced significant complexity at the molecular level. Rather than making the molecular mechanism more complex still (via option one—by introducing some basis for variation in intrinsic fitness) we

Class	Unmodified	Modified
0	✓	✓
1	✓	×
2	✓	✓
3	✓	×
4	✓	×
5	✓	×
6	✓	✓
7	✓	×
8	✓	×
9	✓	✓

Table 6.9: Binary Replicase Interaction Systems Possible in MCS-1 (Modified)

have taken the approach of “re-simplifying” MCS-1 (via option two: by limiting the realisable molecular level dynamics). It is clear from the analysis of the problem observed in Figure 6.5 presented at the end of the previous section that at the very least, class-1 and class-4 pairwise interaction systems must be inhibited, while the specific treatment of classes 2, 5, 7 and 9 is not so important³. In this way, we still retain the core idea of having an *indirect* relationship between primary structure and binding pattern, but we limit the realisable molecular level dynamics to just quasi-stable dominance and parasitic takeover which underpin heritability and mutation respectively at the protocell level.

The further modifications to the MCS system to implement this change are straightforward and are applied at the molecular binding stage of the reaction algorithm (Section 5.1). A list of allowable reaction types will be maintained and consulted upon each interaction attempt. This reaction filtering makes it possible to prevent certain pairwise interaction dynamics that are known to cause the breakdown of our simplified, $n = 2$, pairwise assumptions and analysis. Essentially, the modifications described here involve modifying Table 6.8. That table has been reproduced here as Table 6.9 and includes an extra column to show the proposed changes.

³A significant limitation of this assumption is that it neglects interactions between mutant molecules, which may not always have a negligible effect on the dynamics. This is an artifact of the “pairwise” ($n = 2$) analysis approach, since as n increases, the number of possible classes of interaction system rises to 2^{n^2} , and the validity of viewing the dynamics as a simple superposition of pairwise dynamics becomes much less clear.

6.3.2 MCS-1 Molecular Evolution

Considering the further modifications presented in Section 6.3.1.1, it might appear that MCS-1 offers no substantial change from MCS-0. However, the key difference between the restricted MCS-1 and MCS-0 is that MCS-1 allows parasites to be shorter than, or the same length as their hosts, in addition to the MCS-0 “longer-only” parasites case. This is significant because it provides the mechanism to break the evolutionary stalemate that was observed in the experiments of Section 6.3.3. The following experiment was carried out to determine the extent to which these changes affect the behavior of MCS-1.

- M , the total number of molecules allowed in the reactor was set to 1×10^3 ,
- the reactor was seeded with a normalised population of molecules centered on the randomly chosen self-replicase species: $x_1 = 0000000100$,
- the per-bit mutation rate, v , was set to 0.01. Although this represents a change in mutation rate from the value of 0.05 that has been applied previously, the MCS-0 and MCS-1 systems are not directly comparable in terms of their mutational dynamics with respect to parasitism. This alternative rate was chosen to enable the presentation of a qualitatively similar result to those presented in Chapter 5, even though the experiments cannot be compared *quantitatively*
- the proportion of mutations which are bit-flips, β , set to 0.99 (Section 4.4.5.2).
- the reaction restrictions outlined in Table 6.9 were applied at the molecular binding stage of the reaction algorithm (Section 5.1).

Figure 6.6 shows the outcome of one such experiment. It is clear from Figure 6.6 that:

1. quasi-stable dominance by a principal molecular species has been restored, and,
2. the “longer-only” parasitic displacements, characteristic of MCS-0, are no longer the only possibility; indeed, all displacement events observed were by equal length parasites. Table 6.10 shows an analysis for each displacement event.

Though Figure 6.6 does not show any displacements by shorter species, there is no doubt that such displacements are at least possible—Table 6.11 shows an

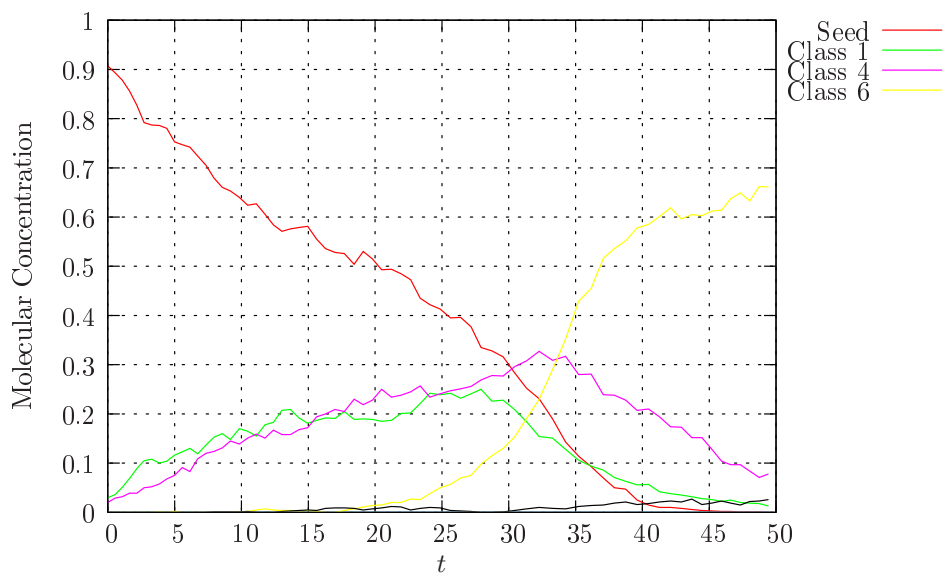


Figure 6.5: Classwise examination of MCS-1 decay of seed self-replicase

Catalyst (<i>primary</i>)	0000000100	Catalyst (<i>primary</i>)	1000000100
Catalyst (<i>secondary</i>)	LLLLL	Catalyst (<i>secondary</i>)	HLLLL
Substrate	1000000100	Substrate	1110000101
Match?	yes	Match?	yes

Table 6.10: Examples of Molecular Matching (MCS-1)

Catalyst (<i>primary</i>)	0000000100	Catalyst (<i>primary</i>)	1000000100
Catalyst (<i>secondary</i>)	LLLLL	Catalyst (<i>secondary</i>)	HLLLL
Substrate	100000010	Substrate	111000010
Match?	yes	Match?	yes

Table 6.11: Examples of Shorter Parasites in MCS-1

example parasite which is shorter than its host for each of the first two molecular species observed in Figure 6.6.

6.3.2.1 MCS-1 Mutational Neighbourhoods

With this formulation of MCS-1, we can now revisit the mutational neighbourhood question that was first addressed in Section 5.1.3. The modifications implemented in MCS-1 introduce some new complications that alter the long-term dynamics of MCS-1. Recall that in MCS-0, all species

- were self-replicases
- were guaranteed to have class-6 parasites at most one mutational operation from their own structure.

In Section 5.1.3 we also calculated the lower bound on the probability of generating a class-6 parasite for an MCS-0 molecule, using parameters $v = 0.05$, $l = 10$ and $\beta = 0.99$, as 3.15×10^{-4} . MCS-1 unfortunately does not afford such a simple analysis. Firstly, not all species in MCS-1 are self-replicases which immediately has the secondary effect of limiting the number of class-6 parasites, since such a parasite must also be capable of self-replication, by definition. Secondly, a further limitation on the analysis of class-6 parasites is that there is no direct algorithm to create a parasite from a given molecule as there was in MCS-0—by adding a random bit to the beginning or the end of a molecule. This means that there is a potential that a particular species might arise as the dominant species in a reactor such that it has no “nearby” class-6 parasitic mutants available to displace it. It is an unproven conjecture that given sufficient replications of a master sequence at some non-zero mutation rate, a class-6 parasite will inevitably arise but there is no straightforward analytical method to calculate even a lower-bound probability for the occurrence of a class-6 parasite for a given replication event. It is, however, possible to carry out a similar monte-carlo analysis to that presented towards the end of Section 5.1.3, given a particular starting molecular sequence. This monte-carlo approach to analysis will be adopted in the next section to understand the molecular selection events that arise in MCS-1.

6.3.3 Protocell Evolution in MCS-1 Systems

Assuming that the above modifications have successfully recovered our ability to use dominant molecular species as a protocell inheritance mechanism, we can once again apply hierarchical selection in the form of artificial protocells and examine the effects that this higher level of selection has on the outcome of the evolutionary runs. Furthermore, if we have indeed broken the evolutionary stalemate that existed in MCS-0 protocells, then we would also expect to see ongoing variability—due to mutation—in the traits that are heritable at the protocell level. The key difference however is that we can safely expect to observe that there is no longer a tendency for the molecular chemistry to favor a trend of increasing molecular length, since molecular parasites are no longer restricted to being longer than their hosts (Section 6.3.2).

6.3.3.1 Protocell Population Dynamics in MCS-1

Given that we expect the modifications made to the MCS in this chapter will disrupt the selectional stalemate phenomenon, we therefore also expect to see occasional protocell level population displacements. The population dynamics at the protocell level are that of straightforward exponential growth under the influence of variable reproduction and mutation rates, and equal death rates. Depending on the particular protocell strains involved, population dynamics could be driven by any combination of the following three forces:

1. Drift: protocell strains dominated by molecular species of the same lengths will have approximately comparable gestation times and reproduction rates. Depending on the relative structures of the dominant molecular species of a pair of protocells, they may also have comparable protocell mutation rates. These situations of comparable Darwinian fitness may lead to a situation of population drift with eventual fixation.
2. Reproduction rate: In Section 5.1.4.2, we saw that given the length, l , of the dominant molecular species in a protocell, we could derive a value for the average gestation time for that protocell strain. We can thus infer a value for the average reproduction rate for a given protocell strain as $\frac{1}{\text{gestaionTime}}$. Since all protocells have the same constant death rate, then the variability of birth rates, B , provides another axis for protocell level fitness:

$$\begin{aligned}
B &= \frac{1}{\text{gesta}t\text{i}o\text{n}T\text{i}m\text{e}} \\
&= \frac{1}{\frac{\ln(2)}{(1-v)^{2l}}} \\
&= \frac{(1-v)^{2l}}{\ln(2)}
\end{aligned}$$

3. Protocell mutation rate: In Section 5.1.4.5, a method was presented for estimating the average protocell mutation rate for a particular strain of protocell (Equation 5.10). Table 5.6 showed that for a reactor of size 500, the probability that a single parasite will take over the reactor is approximately 0.038. In the protocell experiments presented in this chapter, S_{max} has been set to 150. Based on the results presented in Table 5.6 then, we propose to use a multiplier of 0.05 to discount the protocell mutation rates for protocells with $S_{max} = 150$. This “discounting” is intended to allow—very approximately—for the effect of back-mutation of mixed protocells while the absolute number of contained parasite molecules is still relatively small. In other words, the “effective” mutation rate is significantly less than the “raw” mutation rate (to mixed cells, which neglects the effect of back-mutation).

$$R = 0.05(1 - (1 - r)^{\frac{S_{max}}{2}}) \quad (6.1)$$

It is important to note that the absolute values of the calculations above do not have any intrinsic application. Only in comparison to the values obtained for another protocell strain do these calculations hold significance. In particular, the values obtained for a protocell strain can be compared to the values for various mutant strains of the original, and this comparison might be used to understand displacements at the protocell level.

6.3.3.2 MCS-1 Protocell Experiments

The experiments to explore protocell evolution in MCS-1 were set up as follows:

- $v = 0.025$; $\beta = 0.99$
- $maxCells = 500$; $S_{max} = 150$
- An initialiser molecular species of length $n = 20$ is randomly chosen, with the added condition that this initialiser molecule is capable of self-replication.
- $maxCells = 500$ are initialised with the initialiser species.
- The process for initialising individual protocells ensures that at $t = 0$, the protocell population is in a quasi steady-state.
- Repeated molecular interactions, as described in Section 4.4.5, are then carried out, while also respecting the further interaction restrictions presented in Table 6.9.

The results of three independent runs of this experiment are presented below. Once again, an appropriate per-bit mutation rate, this time 0.025, has been applied to enable the presentation of results that can easily be qualitatively contrasted with those results presented in Chapter 5. Each run was initialised with a different initialiser molecular species and demonstrated qualitatively distinct behaviours.

1. Figure 6.7 is superficially reminiscent of the molecular stalemate result achieved with MCS-0 (compare Figure 6.7 with Figure 5.13). Monte Carlo analysis of the initialiser species, however, reveals that after 2×10^7 replication events, only 273 individual instances of class-6 parasites were produced—in other words, the probability of producing a parasite molecule on a single replication event, r , for this initialiser, is 1.37×10^{-5} . According to equation 6.1, this lineage therefore has an effective protocell mutation rate of approximately 5.1×10^{-5} . In other words, it is estimated that one protocell for every 2×10^4 reproductions will eventually mature to a new pure-line strain. For the timescale of this experiment, we can estimate that 75 *new* instances of pure-line cells emerged from the original pure strain during the run. To understand why this initialiser species maintained dominance of the protocell population, we can compare its effective protocell mutation rate with those of the mutant cells to which it gave rise. Table 6.12 shows the effective protocell mutation rates for initialiser and the three most common mutant protocell strains in Run 1. As before, these rates are estimates based on 2×10^6 Monte Carlo trials. Each of the mutant strains presented in this table have effective protocell mutation rates at least an order of magnitude higher than the initialiser strain. Since

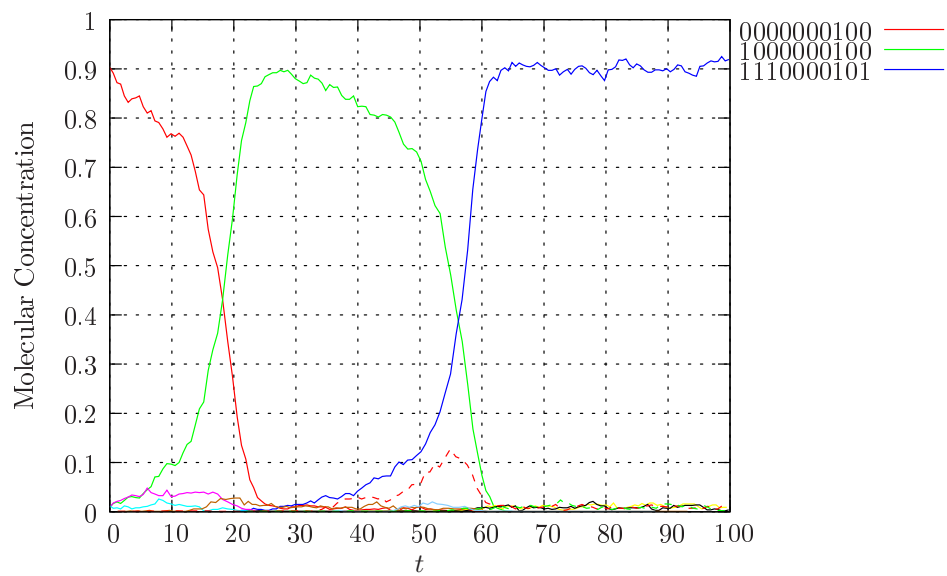


Figure 6.6: Molecular Evolution in MCS-1

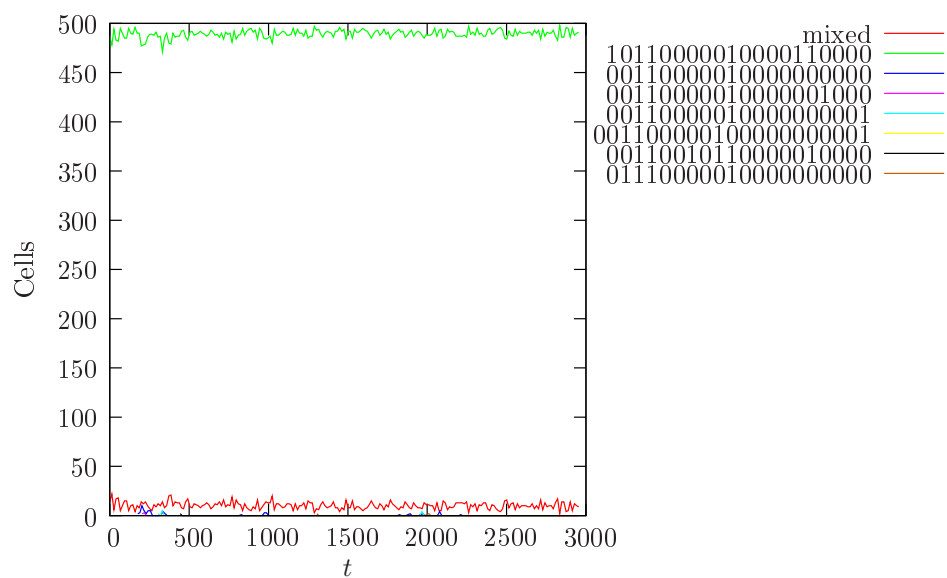


Figure 6.7: MCS-1 Cells

Dominant Molecular Species	Molecular Parasite Production Rate (r)	Estimated Effective Protocol Mutation Rate
101100000100000110000	1.365×10^{-5}	5.1×10^{-5}
001100000100000000000	2.29×10^{-3}	6.2×10^{-4}
001100000100000001000	3.06×10^{-3}	8.4×10^{-4}
001100000100000000001	3.6×10^{-3}	1.0×10^{-3}

Table 6.12: Estimated Effective Protocol Mutation Rates for Figure 6.7

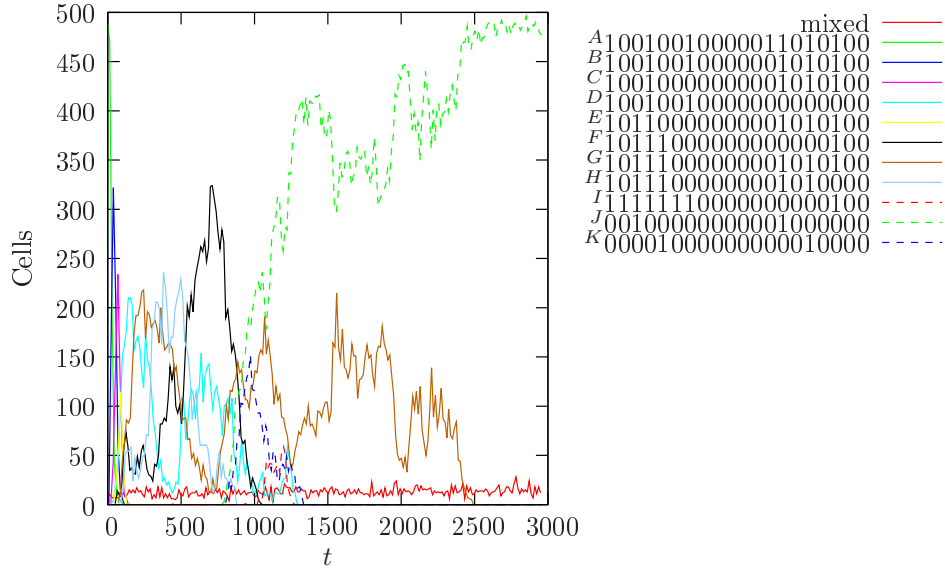


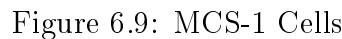
Figure 6.8: MCS-1 Cells

all of the protocell lineages examined are dominated by molecular species of the same length, their protocell reproduction rates are equal, and therefore the difference in mutation rates is the only factor affecting relative fitness in this case.

- Figure 6.8 shows a run which is initialised using a molecular species which turns out to be extremely unstable with respect to parasitic mutation. Figure 6.8 was produced by disregarding any pure cell lineage which did not reach an overall concentration of 0.1 (50 cells out of 500). Monte carlo analysis of this species reveals that the probability of generating a parasitic mutant is 1.695×10^{-2} , three orders of magnitude greater than the initialiser used for Figure 6.7. The effective protocell mutation rate for the lineage dominated by this molecule is estimated at 3.6×10^{-2} according to equation 6.1—in other words, around 18 new pure-line protocells are produced *every generation* by this protocell strain. It is unsurprising, then, that the initial pure line of protocells in Figure 6.8 is rapidly replaced by a new protocell strain. In fact, the dynamics of the beginning of this experiment appear to show many such relatively unstable protocell strains. As the experiment progresses, however, new protocell strains emerge which are more stable due to lower parasite production rates for their dominant molecular species. Table 6.13 shows the results of Monte Carlo analysis for the dominant molecular species of each of the pure cell lines that was observed during the course of this run, as shown in Figure 6.8. Protocell

Dominant Molecular Species	Molecular Parasite Production Rate (r)	Estimated Effective Protocol Mutation Rate
A 100100100000011010100	1.69×10^{-2}	3.6×10^{-2}
B 10010010000001010100	1.74×10^{-2}	3.6×10^{-2}
C 10010000000001010100	1.82×10^{-2}	3.7×10^{-2}
D 10010010000000000000	4.4×10^{-4}	1.6×10^{-3}
E 10110000000001010100	1.72×10^{-2}	3.6×10^{-2}
F 101110000000000000100	4.6×10^{-4}	1.7×10^{-3}
G 10111000000001010100	8.5×10^{-6}	3.2×10^{-5}
H 10111000000001010000	1×10^{-5}	3.8×10^{-5}
I 11111110000000000100	4.5×10^{-7}	1.7×10^{-6}
J 00100000000001000000	2.32×10^{-5}	8.7×10^{-5}
K 00001000000000010000	1.52×10^{-5}	5.7×10^{-5}

Table 6.13: Estimated Effective Protocol Mutation Rates for Figure 6.8



131

3. The behaviour presented in Figure 6.9 is qualitatively different from the previous two runs. Figure 6.9 was produced by disregarding any pure cell lineage which did not reach an overall concentration of 0.1 (50 cells out of 500). In it, there are two clear displacement events, both by cell lines “shorter” than the previously dominant line. Table 6.14 shows the results of monte carlo analysis for each of the dominant molecular species of the protocell strains observed during this run. The symbols A , B and G are used to highlight the three cell lineages which reached a cellular concentration of over 0.9. The displacement of A by B occurs due to alignment of the selectional forces presented in Section 6.3.3.1. The dominant molecular species of strain B is shorter than that of strain A . Strain B therefore has a shorter gestation time than strain A . Furthermore, strain B has an estimated effective cellular mutation rate less than half that of A . The displacement of B by G is interesting because, while B has a lower estimated effective cellular mutation rate than G : G is two bits shorter than B , which means that G protocells have significantly shorter gestation times than B cells, and can therefore reproduce faster than B cells. Further inspection of the results reveals that during the “reign” of the B strain, there are three other protocell strains whose dominant molecular species is one bit shorter than that of the B strain. These strains have a lower gestation time, and therefore a higher overall reproduction rate, due to their shorter dominant molecules. However, strain B has an estimated effective cellular mutation rate which is at least an order of magnitude less than the three “shorter” cell lines (D , E and F) which tried to compete with it.

6.3.3.3 Summary Remarks

These results verify the complicated interplay between selectional forces at the protocell level and demonstrate that there are at least *some* cases where having a shorter dominant molecule does not automatically imply an effective selective advantage, as was the case with MCS-0. The corollary here is that a situation may arise where a *longer* cell line displaces a shorter one, on the condition that the longer cell line has a significantly lower cellular mutation rate, though it is unclear as to how significant this difference needs to be to offset the benefit of being shorter.

Dominant Molecular Species	Molecular Parasite Production Rate (r)	Estimated Effective Protocol Mutation Rate
A 010111010001000000010	1.21×10^{-4}	4.5×10^{-4}
B 00011010001000000010	3.5×10^{-5}	1.3×10^{-4}
C 0011101000110100010	1.15×10^{-5}	4.3×10^{-5}
D 0011101000100000001	1.42×10^{-3}	5.1×10^{-3}
E 0011010010000000001	4.78×10^{-4}	1.8×10^{-3}
F 001101000000010000	1.45×10^{-3}	5.2×10^{-3}
G 001101001000000001	4.91×10^{-4}	1.8×10^{-3}
H 00100100100000101	9×10^{-4}	3.3×10^{-3}

Table 6.14: Estimated Effective Protocol Mutation Rates for Figure 6.9

6.4 Summary of Experimental Results

In Chapter 5, we saw that MCS-0 was extremely limited in its potential for cell-level evolution most notably due to the “selectional stalemate” phenomenon. In this chapter, MCS was enriched by incorporating the concept of partial molecular recognition through a simple molecular folding mechanism. We have seen in the early experiments of this chapter, however, that the reduced level of molecular recognition caused the breakdown of the ability of the molecular dynamics to support both heritability and occasional mutation at the cell level—and thus of any continuing selectional evolutionary process at the cell level. This problem was isolated and corrected by placing some further restrictions on the pairwise interaction classes that were permitted in the system. In fact, the restrictions put in place had the net effect of limiting the allowable pairwise interaction classes to precisely the set that occurred naturally in MCS-0. However, as the later experiments of this chapter show (Section 6.3.3), the phenomenon of “selectional stalemate” can be successfully disrupted using the modifications to MCS-0 presented here.

A consequence of the more complex binding function adopted in MCS-1 is that it is no longer possible to analytically examine a particular molecular sequence to determine any upper or lower bounds on the probability of producing parasites during replication. In the case of MCS-0, *every* molecular species was a self-replicase, and the molecular binding rules applied made the production of parasitic molecules trivial. The binding rules used for MCS-1, however, introduce the possibility that some molecules might not be capable of self-replication. The lower bound probability of producing a parasitic molecule on a single replication in MCS-1 then is zero. In the absence of any usable analytical estimate, monte carlo trials were used to examine molecular species that had been observed in the simulation and parasite production rates could be determined on a per species basis, allowing us, in turn, to roughly estimate a consequent effective cellular mutation rate. The combination of cellular mutation rate and cell gestation time then determined relative fitness between protocell species, and the population dynamics showed examples of both selection and drift.

6.5 Conclusion and Wider Impact of Results

The main aim of this thesis has been to design and build a minimal artificial life system incorporating agents operating at at least two interacting hierarchical levels of selection, exhibiting unlimited heredity at both levels, and to test and characterise the resulting evolutionary dynamics. The particular hierarchical levels chosen for modeling were the molecular level and the cellular level of a novel protocell architec-

ture. Careful attention has been given to understanding the evolutionary dynamics of each of the hierarchical levels individually, as the concept of hierarchical selection holds that the evolutionary dynamics of the system as a whole will depend on the lower level interactions that take place. In the framework of the major evolutionary transitions, the MCS platform could be seen as an abstraction of the first major transition—from individual replicating molecules to populations of such molecules constrained to coexist inside a rudimentary cellular membrane. The approach taken has been to significantly abstract away from the reality of the bio-chemistry in which life as we know it must operate (Sections 2.3 and 4.1) and to focus on the evolutionary dynamics that arise in such hierarchical systems. This approach might also lend itself to the study of the dynamics of the other major transitions, for example those involving social insects or primate societies.

The major evolutionary transitions are thought to have been critical to the evolutionary growth of complexity in life as we know it (Maynard Smith and Szathmáry, 1997), though the concept of hierarchical selection might also apply outside the fields of artificial life and the origin of life. The MCS platform was never intended to be a concrete model of life as we know it, but rather an abstraction which could be used to explore the evolutionary dynamics in hierarchical systems, and as such, the results achieved in this thesis might have an impact outside the specific area of artificial protocell research. Potentially, any hierarchical system in which the survival of the whole depends on the actions of agents at lower levels of the hierarchy might be subject to hierarchical selection. In business for example, large organisations are typically organised as a hierarchy, with clear-cut division of labour between different levels. The success of the enterprise as a whole depends upon the combined actions of the lower level business units, who themselves, in some cases, could not exist as entities in their own right. The drastic abstraction approach adopted by this thesis might lend itself to the modeling of complicated systems such as these, and the emphasis on understanding the dynamics at all levels of the hierarchy ensures that meaningful results can be extracted from such models and perhaps applied in the real world.

More specifically, the PACE project (Programmable Artificial Cell Evolution) (EU-FP6-IST-FET-002035) provided funding for the majority of the research presented in this thesis. PACE received EU funding on the basis of carrying out research in the field of Information and Communication Technologies (ICT). The objective of the PACE project was to investigate artificial protocell fabrication, and the potential applications of such artificial protocells. The PACE project benefited from the co-operation of a large number of partners across multiple disciplines, and involved the modeling of protocells using many different approaches and over a range of scales.

The work presented in this thesis has contributed to one of these project branches. Specifically, a significant objective of the PACE project was to investigate whether protocells (Section 3.3) have potential as novel computational devices. Appendix B describes a series of modifications that were made to MCS-1 with a view to carrying out such an exploration using artificial protocells. The generic concept is that certain molecular transformations might be used to represent computations, and if so, populations of such molecules constrained to co-exist within protocells might be programmed to carry out an aggregate computation at the protocell level. For the purposes of this thesis, and the PACE project in general, programming is expected to be achieved via directed evolution. In Section 2.3, the DNA based translation mechanism used by life as we know it was described. Protocells, by definition, are not endowed with such a complex translation mechanism, which raises the question of whether it is even feasible to apply selection at the protocell level in order to evolve appropriate sub-populations of molecules. MCS-2, presented in Appendix B, provides a proof of principle, at an extreme level of abstraction, that hierarchical selection can be used in this way to evolve molecular sub-populations appropriate for carrying out a very basic computational operation.

Part IV

Concluding Remarks

Chapter 7

Discussion

Life, as we know it, is a fine example of growth of complexity by evolutionary means. Pioneers like Darwin, Mendel, Watson and Crick described the low-level mechanisms by which evolution by natural selection produced the diversity of life that we see around us from a humble, most likely shared origin. Further research in this field suggests that hierarchical selection has had a major contribution to the evolution of complexity in life as we know it. The almost mechanical nature of living systems has led many researchers to attempt to reproduce them both physically in the laboratory and virtually, using computers. From this “artificial life” perspective, many key problems for origin of life theories have been identified, and some have had plausible solutions proposed. To date, however, no artificial life system has demonstrated an evolutionary growth of complexity comparable to that observable in the evolutionary history of life as we know it. This main contribution of this thesis has been the design and implementation of a model of hierarchical selection which demonstrated interesting evolutionary behaviours which were due to the interaction between selectional pressures at different levels of the hierarchy. This proof of concept might make a significant contribution to the realisation of the evolutionary growth of complexity in future artificial life models.

As Dawkins (2009) remarks: “Right up to the middle of the twentieth century, life was thought to be qualitatively beyond physics and chemistry. No longer. The difference between life and non-life is a matter not of substance but of *information*.” Many research groups have approached the problem of creating life ex-nihilo from many different directions. In anticipation of *some* future success in building artificial life-forms, the work presented in the experimental chapters of this thesis set out to explore the kinds of evolution that we might reasonably expect to observe in such systems, and to provide a framework for the incorporation of hierarchical selection in artificial life models. In the appendix material, a further demonstration of the MCS system is presented and aims to explore the potential for using MCS style

protocells as novel computational devices.

7.1 Recap of Key Achievements

Throughout this thesis, **four key achievements** were presented.

1. A novel artificial chemistry, **Molecular Classifier System (MCS)**, was developed as a framework for the investigation of the effect of hierarchical selection in artificial life systems. MCS was built to support hierarchical selection and incorporated a rudimentary membrane chemistry to enable the construction of artificial protocells.
2. A **generic ODE analysis** was developed to facilitate an understanding of the dynamics of the molecular interactions in MCS. Given that these catalytic interactions are *trans-*acting, the resulting population growth is hyperbolic, rather than exponential. The pairwise analysis presented in Section 4.5.4 provides a useful classification for protocell species in the more advanced versions of MCS, but more importantly, it enables the low level analysis of the molecular dynamics of the protocell experiments presented in Section 6.3.3 (and later in Section B.5).
3. During the software testing phase of the work, a previously unseen phenomenon, **selectional stalemate**, was demonstrated. This phenomenon arose from the impoverished molecular binding mechanism used in MCS-0 (Chapter 5), and its identification is the third key achievement of this work. Having identified this selectional stalemate, the focus then shifted towards potential mechanisms to disrupt it, while at the same time retaining *some* heritable trait at the protocell level.
4. In Chapter 6, the concept of fitness at the molecular level was expanded to include not only the sequence length of molecules but also the richness of the mutational neighbourhood of the molecules in terms of available class-6 parasites. This is in contrast to typical approaches to replicator fitness in artificial life systems which incorporate intrinsic differences based only on the efficiency of the replication process across molecular species. This chapter also presented a mechanism to expand the evolutionary potential of the MCS—a rudimentary molecular folding mechanism which enforced a more “partial molecular recognition” than the “complete molecular recognition” that went before it. The resulting demonstration of ***ongoing* protocell evolution** is the fourth and final major contribution presented in the thesis.

7.2 Potential for Further Work

One of the key motivations of the PACE project, which partially funded this research, was the investigation of the utility of artificial protocells as novel computational devices. Appendix B presents a modified MCS which demonstrates the **directed evolution of protocell embedded molecular computation**, and is presented as a supplementary contribution of this thesis. This work sees the inclusion of an externally imposed molecular activity factor to apply a fitness differential based on protocell size. This extended version of MCS stands as a proof-of-concept that an artificial chemistry can support protocell embedded molecular computation. There also exist potential pathways that could be taken to expand upon the core MCS as it has been presented in the main body of the thesis, not to mention the potential for investigating the applicability of the results presented here for incorporation in a programme of wet-lab experimentation. In terms of a road-map for incremental further development of the MCS system *in silico*, there are at least two routes that might be followed:

- a deeper investigation of the **evolutionary phenomena** that were identified, at both levels of the hierarchy,
- implementing extended **computational functionality**,

The following sub-sections describe each of these pathways in further detail.

7.2.1 Further Investigation of Evolutionary Phenomena

During the work presented in the experimental chapters of this thesis (Chapters 5 and 6), it became apparent that in spite of being an idealised, supposedly simplified, artificial chemistry, it nevertheless demonstrated interesting evolutionary behaviours.

7.2.1.1 Molecular Phenomena

Typically, replicator systems such as MCS will measure replicator fitness as a function of the speed at which it can replicate something—*intrinsic* fitness. Even though molecules in MCS had no intrinsic fitness differences by this definition, various methods were identified to describe the effective fitness of molecules. This effective fitness could be derived from a number of sources. Molecular length, l , and the composition of the mutational neighbourhood were both identified as contributing components to molecular “fitness”. While these techniques assisted in the explanation of the complicated evolutionary behaviour observed during the experiments, most notably

the detail presented in Appendix B, our understanding of the low level evolutionary phenomena remains poor. That being said, these results highlight the fact that full dynamical systems approach becomes impractical as the number of interacting species increases. Our approximation that the dominant molecular species was responsible for the all protocell growth is based on the fact that we have disabled most of the **binary replicase interaction systems** that were introduced in Section 4.5.4. Interesting results might be achieved by selectively re-enabling some of those interaction systems, especially since the analytical method used during this thesis relied completely on the assumption that the reactor could be approximated as a super-position of multiple independent pairwise ($n = 2$) systems. Increasing n however will necessarily void this analytical approach, though frameworks like **Chemical Organisation Theory** (Dittrich & Speroni, 2007) could prove fruitful, at least at the molecular level.

7.2.1.2 Protocell Phenomena

The coupling between molecular level selection and the fate of entire protocells is currently based on a combination of molecular level dynamics and the application of an externally imposed molecular activity factor. The modifications discussed in Section 7.2.1.1 would be expected to have at least *some* effect on the protocell level dynamics, but it is clear that there are also incremental changes that could be applied at the level of protocell dynamics. For example F , the **protocell fissioning factor**, might be *decreased* by certain reaction types. A more involved modification might be to add a function similar to the “**quorum sensing**” that was described in Section A.3.3, so that molecules might have access to the precise instantaneous concentration of a particular molecular species, or even the concentration of the fissioning factor itself. Such modifications, though implemented at the molecular level, would presumably give rise to phenomena that are only appreciable at the protocell level.

7.2.2 Extension of Computational Capabilities

MCS was inspired by Holland’s **Classifier Systems**, as introduced in Section 3.2.2.2. Conceptually, MCS has rules and messages, though these are instantiated inside one and the same physical object, the molecule, or to be more specific, the molecule will act as a rule or as a message, depending on the its context for a given reaction (**substrate** or **catalyst**). A precursor to Holland’s learning classifier systems was what he called the “Broadcast Language” (Holland, 1975). Decraene et al. (2006) presented a modified version of MCS which aimed to model artificial cell

signaling networks, and which used the Broadcast Language to represent molecules. The version of MCS presented in this thesis was always intended to be a sample of a precursor system to that presented by Decraene. Incremental modifications could be made to the *Comp_{func}* implementation (Section B.3.1.2) in order to extend the enzymatic computational capabilities of MCS molecules, most likely inspired by the Broadcast Language.

7.3 Limitations and Potential Criticisms

The MCS as presented in this thesis is a highly abstract platform upon which a model universe was built and used to investigate the effects of hierarchical selection on the virtual agents contained within it. MCS is not intended to be a model of any particular target biological system, but rather a novel artificial universe which is inspired by aspects of biological life which have until now received little attention from the artificial life community. With that in mind, it is fair to say that MCS is indeed a highly abstract toy-model, and while it has been a successful demonstration of hierarchical selection in an artificial life system, the immediate applications of the model in its current form are difficult to predict. Some of the behaviours exhibited by the system, most notably the phenomenon of selectional stalemate presented in Chapter 5, are arguably due to the configuration of the model itself, rather than indicative of the kinds of behaviour that might be expected in other artificial life systems which implemented this kind of hierarchical selection. Another potential criticism of the model as presented is that the idealisations which have been applied have tended to be in areas which have received considerable attention in the artificial life literature, most notably, the molecular replication mechanism. Typically, artificial life modelers will focus much of their attention on the mechanics of the replication mechanism, as this is seen as a crucial component. In defense of the idealised approach adopted by MCS, consider the complex molecular interaction dynamics observed throughout Chapters 5 and 6, the understanding of which would not have been possible without the generalised ODE model described in Section 4.5. In spite of the simplified replication mechanism, complex dynamical behaviours were nonetheless observed. Were it not for the idealised nature of the molecular structure and replication mechanism, these behaviours might not have been so straightforward to explain. As discussed in Section 4.5, this model can be generalised to any $n = 2$ system of *trans*-replicators, which would suggest that regardless of the details of the molecular replication mechanism, qualitatively similar behaviours would be expected, adding weight to my claim that an idealised approach is what is required to enable a more complete understanding of the behaviour of the system. The various “idealisation” steps that

were taken during the construction of the artificial protocell system presented in this thesis highlighted the kinds of problems that *in vitro* artificial protocell systems might also be faced with. On the other hand, the core properties of the MCS as presented in this thesis are still far beyond the capabilities of any wet-lab systems to date, suggesting that success in creating artificial life forms *in vitro* will be the *beginning* of a journey, rather than the end.

7.4 Conclusion

I have achieved many personal goals during the programme of research leading to the production of this thesis beyond the concrete area of artificial life. It has been a journey which involved learning many new transferable skills. For example:

- research methods
- time management
- project planning and organisation
- academic writing and presentation

Also, since artificial life is such a broad interdisciplinary field, I was required to develop significant knowledge outside my previous academic spectrum. At undergraduate level, I studied Software Engineering, and I completed a taught Masters degree in Digital Security and Forensic Computing prior to beginning this research. Thanks to the fact that the majority of my research was funded by a large EU project—the FP6 funded PACE project (FP6-IST-FET-002035; McCaskill et al. (2008))—I was given the opportunity to work with researchers outside my undergraduate fields. This gave me exposure to members of the community who would not normally attend artificial life conferences per se, but whose work was so closely related to my own that it was possible to have many productive conversations that otherwise would not have been possible.

Part V

Appendices

Appendix A

Living Computation

In Chapters 2 and 3, the distinctions between “life as we know it” and “life-as-it-*could-be*” were highlighted. The motivations for drawing these distinctions have been, so far, philosophical in nature, and have essentially focused on:

- the issues surrounding the definition of life,
- identifying the common sub-systems which underpin life
- simulating living processes inside computational devices

In Section 2.5, it was shown that the challenges facing life can be essentially split into two categories:

1. physico-chemical challenges
2. systemic, organisational challenges

The first category of challenges are specific to the details of life as we know it, and are therefore most suited to being addressed by the sciences of Chemistry, Biology and Physics. The second category however is a more abstract set of challenges which are specific to the general phenomenon of “life”, rather than any particular incarnation of it. In the final section of Chapter 3, some of the current approaches to building artificial protocells were described. If such artificial protocells are realised, that is to say, the physico-chemical challenges can be overcome, then these life-forms too will presumably be faced with a similar set of systemic, organisational problems to those faced by life as we know it. If we think of the physico-chemical makeup of the living system as defining the possible interactions that may happen, then the key to solving problems is by re-organising and modularising, rather than requiring new primitives. In this chapter, we shall narrow our focus to some of the interesting problems that living systems have succeeded in solving. So, the purpose of this chapter then is to identify the kinds of task that “life as we know it” has proved to

be very good at dealing with, and then examine the possible origins of these abilities so that the emergence of similar capabilities in artificial lifeforms such as protocells can be encouraged and, more importantly, *recognised* in their earliest stages to give a better understanding of the emergent pathways that exist.

A.1 Introduction

In order to identify the kinds of things that we would expect artificially created life-forms to be able to do, it will be instructive to at least recognise the key abilities of life as we know it. In the earlier chapters of this thesis, we were particularly interested in the interaction between the various low-level systems that give rise to living things—membrane chemistry, metabolic chemistry, information chemistry. However, given a living system as a starting point, a new set of potential problems arises—problems which relate to sustaining life. It is obvious that the problems faced by a life-form in relation to maintenance of life are specific to that life-form and also potentially the environment in which that life-form exists. The focus of this chapter will be on a small subset of the challenges facing various forms of *uni*-cellular life, and examining the elegant solutions that have evolved to overcome these problems. The problems facing *multi*-cellular life-forms, for example humans, are the subject of many other works of science and art, though at the nano-scale, interactions both within and between cells are responsible for the high-level emergent behaviour that we see. By focussing on the molecular level organisation of single celled life, it should be possible to examine the simplest mechanisms that life has devised to overcome these problems. We shall see that the mechanism by which these problems are faced and overcome is analogous to “information processing”, which in turn implies “living computation”. While it is not difficult to imagine the mammalian brain as a kind of “living computer”, extending this abstraction to cover the day to day activities of single cells requires some deeper thought. The following sections will first present a more detailed description of the general case for “living computation”, before describing some examples of this kind of computation at work in contemporary biological cells.

A.2 What is Living Computation?

The following definitions will be used to support the subsequent argument of this section:

1. *Computation* is used here as a general term for any type of information processing and is not to be confused with the more technical concept of *computability*

related to the works of Alonzo Church and Alan Turing (Church, 1936; Turing, 1937).

2. *Information* is a term used to describe the configuration of a system of interest. Once again, this should not be confused with the formal definition of information proposed by Shannon (Shannon, 1948).
3. *Processing* typically describes the act of taking something through an established and usually routine set of procedures to convert it from one form to another.
4. A *System* is a set of interacting or interdependent entities forming an integrated whole.

Like the satisfactory definition of life presented in Section 2.1, the above definitions do not place restrictions on the mechanisms or representations involved. Synthesising these definitions and applying them to the day to day activities of life as we know it, we might say that a living system processes information about its environment¹ every time it interacts with it. This synthesis is clearly related to the concept of “systemic computation” proposed by Bentley (2007). Let us define the environment as a representation of some information about what things are in the environment, and how they interact. Then, if a “living system” were to exist in that environment, the interactions between this living system and the environment can be seen as an exercise in Information Processing. The living system gathers information about its environment and interacts with the environment, and indeed *becomes* a part of the environment for whatever *other* living systems exist. For “life as we know it”, this information processing takes place most obviously in the brains and nervous systems of living systems which possess those faculties. For living things which do not have brains or nervous systems, particularly single-celled organisms, one must look a little deeper to see the level of information processing that takes place, for example:

1. *Vibrio fischeri* (*V. fischeri*) are a quorum-sensing bacteria which bioluminesce when they detect the presence of “enough” nearby individuals of their species,
2. *Escherichia coli* (*E. coli*) are a species of bacterium which is commonly found in the lower intestine of warm blooded animals. *E. coli* can detect the presence of a chemical food source and then move towards that food so that it may feed upon it, a process known as chemotaxis.

¹The environment could also be described as a system in its own right.

Given a pre-existing artificial living system such as the proposed protocell system described in this thesis, this information processing analogy might then be applied to it also. It is clear that having an influence over the *kinds* of information processing that these artificial systems engage in would be very powerful. On some level then, such systems might be described as “living computers”. If this is the case, then the question arises as to how we might program such a computer, and more importantly, how might we insert and extract information from it. Furthermore, if building artificial protocells is indeed an exercise in building “living computational systems”, it will be beneficial to take a closer look at the kinds of computation that life as we know it is already capable of. By identifying sufficient examples of “living computation” as it currently exists, we might then be able to reproduce something similar in artificial protocell systems, and at the very least we can narrow our focus to the key information processing capabilities of living systems.

A.2.1 Programming Living Computers

All life as we know it has been shaped through evolution by natural selection. Given the central position of the concept of evolution in the satisfactory definition of life presented in Section 2.1, it is clear that any new life-forms created will also be subject to the powerful force of evolution, suggesting that evolution is one, if not the only, way to “program” the behaviour of artificial living systems. Living systems evolve to succeed in their environment. By coupling this survival pressure with a suitable environment, we might hope that the living systems evolve to cope with, and eventually exploit this environment. The solution to the programming problem then is to configure the environment in such a way that it encourages evolution toward a particular end point, perhaps in a similar way to the Avida system (Adami and Brown, 1994).

A.3 Examples of Living Computation

At this point, it is almost trivial to identify the “ability to stay alive” as a core information processing ability that all life can do, considering how closely related this is to the very definition of life. The ability to self-repair and self-maintain is seen as one of the most fundamental abilities of life (Varela et al. (1974); McMullin and Varela (1997); McMullin (1997)). However, self-repair and self-maintenance require raw material from the environment and therefore require other supporting functions in order to work properly. This suggests that a strong synergistic coupling between components of these life-forms is essential for the system as a whole to function correctly. Contemporary biological cells are magnificent examples of such tightly

coupled systems which carry out a range of different behaviours. In Section A.2, the ability for some single-celled organisms such as *E. coli* to seek food by the process of *chemotaxis* was highlighted. Higher level computation, such as that carried out by mammalian brains and nervous systems, which at their lowest levels rely on signal transduction between and within the cells that make up those systems, was also briefly mentioned. These computational functions are facilitated by complicated, emergent “cell signalling networks”. These networks enable communication both between and within cells (Palsson, 2006). CSNs also regulate the state of cells in the absence of signals at the membrane. In these cases, the “signals” that the CSN responds to are the internal states of the cell and other CSNs within it. Inside the nucleus of a eukaryotic cell, the process of expressing the genome to build the living thing, or components of it, that it encodes has already been identified as a highly modular task (Section 2.3). The order of events during the expression of the genome, and the ongoing referencing of the genome to maintain and repair the living creature depend on “gene regulatory networks” (GRNs), essentially switching on and off certain portions of the genome at specific times, in accordance with the “program” described by the genome (Palsson, 2006). GRNs have also inspired artificial systems which have demonstrated the ability to “compute”, notably Bentley’s fractal protein system (2004) which exhibited sequential counting as genes were switched on and off. The following sub-sections highlight some examples of “living computation” according to the synthesis of definitions presented in Section A.2.

A.3.1 Cell Cycle Control

Cell cycle control is the mechanism by which cells, both prokaryotic (cells with no nucleus) and eukaryotic (cells with a nucleus), manage their growth and division processes. Prokaryotes and eukaryotes divide by a process of binary fission. Eukaryotic binary fission is a more complicated process than the prokaryotic case and can be described in terms of two phases (Smith and Martin, 1973). The first phase, known as the interphase, sees the cell roughly doubling in size. This size doubling requires the cell to gather enough nutrients from the environment in order to ensure that the surrounding membrane is sufficient for division and that each individual internal sub-system inside the cell is also sufficiently prepared to be divided amongst the daughter cells. The following is a short list of some of the tasks that must be completed during the interphase:

1. all of the genetic information must be duplicated,
2. all of the organelles (cell machinery) must be duplicated,

3. the cell membrane must grow to accommodate the growing number of entities contained within it.

The interphase is followed by mitosis, where the cell divides into two separate daughter cells. In order for the cell-division to be successful, each of the components that have been doubled up must end up in one or other of the cells, as if any component is missing from either daughter cell, that cell will almost certainly be non-functional². This growth and division cycle places a significant demand on the individual sub-systems of the cell, since they must remain functional at each point in the process.

To increase the chances of successful cell division, various checkpoints have evolved to break the process up into stages. These checkpoints enable the cell to determine if everything is going well in the cycle. If a checkpoint fails, the cell will typically abort the normal division process and either revert back to a quiescent (dormant) state or self-destruct by a process known as apoptosis, depending on how much of the division process had taken place up to the point of failure. One such checkpoint, the spindle assembly checkpoint, is responsible for ensuring that the duplicated genetic material, now in the form of chromosomes, is correctly positioned for the cell splitting operation. Other checkpoints monitor the integrity of the reproduction of other cellular components and have the ability to halt the normal cell cycle until such time that the damage is either repaired or the cell self-destructs by apoptosis.

The checkpoints of the cell-cycle are controlled by a complicated interaction between the information in the genome which encodes them and the proteins which carry out the various tasks during cellular reproduction. Mutations to the portions of the genome responsible for these checkpoints may allow cells to grow and divide in spite of the breakdown of genomic integrity. These kinds of mutation are known to be the cause of some types of cancer—cells under the control of faulty machinery or genomic information multiply out of control and form tumours.

While contemporary examples of cell-cycle control networks are intricate in nature, early single-celled life would have possessed far simpler capabilities. The ability to control and synchronise the behaviour of the individual components of the cell, especially during cell growth and division, has clear selectional benefits. There is no doubt that any attempt by a cell lineage to improve the chances that one or both offspring are at least as fit as their parent will be favoured over a cell lineage in which no time or effort is spent on maintaining integrity. In the earliest cells, it is unclear to what extent this control might have been possible, but the ability to

²It might also be disastrous for a cell to have more than one instance of any of these sub-systems, for example, the erroneous inclusion of an extra chromosome during foetal development can cause birth abnormalities.

maintain control over the reproduction cycle would be a great benefit to those cells which evolved it, suggesting that this was an early achievement for life.

A.3.2 Cellular Chemotaxis

Cellular chemotaxis is the process by which cells can move, under their own direction, towards or away from a chemical signal source. Single celled chemotaxis is a particularly interesting form of locomotion because it does not involve a complicated nervous system or brain to control the system—the entire system exists within a single cell. When a sperm cell “swims” toward an egg to attempt to fertilise it, it does so by chemotaxis. *Escherichia coli* uses a similar mechanism to locate food in the intestines of its host (Adler and Tso, 1974). In *E. coli*, the physical component of locomotion is controlled by a number of flagella, or tails, which it uses to propel itself towards, or away from, a chemical signal. An *E. coli* bacterium, as far as chemotaxis is concerned, exists in one of two states:

1. random tumbling,
2. directed swimming.

The *E. coli* is built with the flagella, or tails, along its sides, towards the rear, and with a receptor area at the front. The receptor area at the front is sensitive to changes in chemical concentration, allowing the *E. coli* to determine the relative direction of a chemical signal (food or poison).

In the absence of an appropriate chemical signal, the internal state of the *E. coli* allows the flagella to perform their default behaviour: each independently and randomly rotating either clockwise or anti-clockwise, causing the *E. coli* to tumble around in a non-directed fashion. The presence of an appropriate chemical signal at the receptor initiates a cascade of chemical signals inside the bacterium which trigger all the flagella to rotate in the same direction, propelling the *E. coli* forwards towards the food-source. By alternating between these two states, the *E. coli* can navigate its way towards a food source. *E. coli* is a prokaryote, and divides into two daughter cells roughly once every 20 minutes meaning that the entire cell signalling network which enables this elegant locomotion solution must also be replicated once every 20 minutes, highlighting the importance of harmony between the various cellular control mechanisms.

A.3.3 Quorum Sensing

Another example of emergent computation is the so-called “quorum sensing” capability of populations of the bacterium such as *Vibrio fischeri* (Fuqua et al., 1994).

V. fischeri is capable of bioluminescence—the ability to produce and emit light. Luminescence is a biologically expensive task, and, given the physical size of a single individual bacterium, the amount of light emitted and the relative benefit received would be quite wasteful. *V. fischeri* have therefore evolved a mechanism to determine when the local concentration of bacteria is enough to make bioluminescence worthwhile. The bacteria constantly emit signalling molecules to their environment while simultaneously monitoring the local concentration of these signalling molecules. If this local concentration exceeds a particular threshold, the bacteria luminesce and increase the rate at which they emit the signalling molecule. Since the signalling molecules are subject to constant diffusion once they have been emitted, the chances of a particular bacterium sensing one of its own signalling molecules are slim. Free-living *V. fischeri* then do not bioluminesce as their concentration, and by association, the concentration of the signalling molecule, never reaches the appropriate threshold. When “enough” *V. fischeri* are concentrated in a small enough space, they all begin to sense each others signalling molecules and begin to luminesce. The Hawaiian Bobtail Squid, *Euprymna scolopes*, maintains a symbiotic relationship with *V. fischeri* and uses them to power a light-emitting organ which enables the squid to control the shadow that it casts on the sea-bed as a defence mechanism. *E. scolopes*, when born, begins to gather up *V. fischeri* from the surrounding environment. The squid matures once it has gathered enough *V. fischeri* to power its light-organ (Wei and Young, 1989).

A.4 Conclusion

In this chapter, we’ve seen some examples of how life as we know it is capable of computation. For this “computation” to emerge, it must bestow some sort of fitness benefit on the organism which possesses it—without the fitness benefit, then there can be no selective pressure applied by evolution to maintain and improve the trait. Cell-signalling networks and gene regulatory networks are examples of the kinds of computation that life as we know it is capable of. These communication channels and the corresponding information processing they enable necessarily evolved from somewhere. If the tape were replayed, it would be a fair to assume that *some* form of information processing would evolve, since living entities rely so heavily on this ability for their most basic functions (reproduction, energy acquisition). Once the environment becomes saturated with living things, it is obvious that any of these entities who develop any mechanism by which to exploit their environment more efficiently than its contemporaries will be favoured by natural selection. Given the ongoing process of “blind variation” that characterises evolution by Natural Selec-

tion, and given the infinite number of ways that could arise which might marginally improve the ability for an organism to gather and process information about its environment, the idea that some form of “living computation” might emerge in an artificially created living system does not seem unreasonable.

Appendix B provides a preliminary investigation of the potential of using a system like MCS to model some aspects of (wetlab) protocell computation by examining a minimal extension of the molecular enzymatic function beyond replication to incorporate a rudimentary computational function. The kind of computation that will be demonstrated is relatively arbitrary, and is not motivated by any specific candidate molecular replicator-enzymes that are being considered for use in physical wetlab protocells.

Appendix B

MCS-2

In Chapters 5 and 6, an artificial chemistry, MCS, was described and an artificial protocell system was built using that chemistry. These earlier chapters showed that the MCS platform could be used to simulate hierarchical selection by enclosing populations of MCS molecules inside artificial containers. Furthermore, an emergent protocell level trait was identified which could be used as a target of selection to enable the evolution of these protocells, though no explicit fitness function was applied. The hierarchical experiments carried out in Chapters 5 and 6 demonstrated the evolution of protocell lineages driven by the internal molecular dynamics of the individual protocells, but as yet, the coupling between the molecular level dynamics and the higher, protocell level dynamics has been somewhat arbitrary, in the sense that no explicit coupling was defined. In Appendix A, some examples of the kinds of computation that “life as we know it” can carry out were presented. The earliest computations demonstrated by these emergent control systems can be seen as a side-effect of the coupling of the various systems. One of the motivations of wet-lab artificial protocell development has been that it may be possible to impose artificial selection to evolve protocell lines having some desired functions. This chapter will introduce, implement and examine one possible mechanism for coupling the molecular level dynamics and the protocell level dynamics in the MCS system to explore the directed evolution of a “useful” protocell level trait. Given that MCS is an idealised system, if such directed evolution is unsuccessful here, some doubt might be cast over the possibility of achieving a similar goal in a real, and necessarily much more complicated, artificial protocell system.

B.1 Introduction

Artificial protocells if realised, in *any* medium, would open new technological possibilities in terms of computation and information processing. The potential com-

puting power of living cells is different to that of traditional computing devices in many important ways:

1. Parallelism is inherent
2. Living systems are self-maintaining and self-repairing
3. “Programming” a living computer more than likely takes place through a process of evolution
4. Living systems exist on the nano-scale by default

In light of these possibilities, any exploration of artificial protocell systems would be incomplete without attempting to demonstrate their computational capability. At the same time, exploring the origin and evolution of these computational functions might provide some insight into some of the key problems surrounding the origin of life. Biological life is an example of the coupling of interacting systems to perform some higher level task. The systems self-regulate and regulate each other through their interactions. In many cases, these systems even rely on the same physical components to carry out certain tasks (for example, informational molecules such as RNA may also catalyse some of the metabolic chemistry) thus further cementing the relationship between the coupled systems. In this chapter, the MCS is extended to support a rudimentary computational function, and evolutionary experiments are carried out to explore the evolutionary dynamics that arise.

In Chapters 5 and 6 we implemented and tested an artificial chemistry, MCS, which displayed quasi-predictable molecular dynamics and was suitable for incorporation within a hierarchical model of selection. We explored in considerable detail the extent of the predictability of MCS simulations and presented some issues which required significant modifications to the original MCS design in order to achieve conditions under which protocell evolution is feasible. Specifically, the molecular level dynamics give rise to protocell reproduction, with heritable traits which affect protocell-level fitness where there are at least some accessible mutations which yield protocell strains of different fitnesses. So far, protocell fitness in MCS has been determined by a combination of the average length of the molecular species inside the protocell and the composition of the mutational neighbourhoods of the molecular species which are inside the protocell. In this chapter, an additional, external, artificially imposed fitness function is introduced which is parametrised by the size of the protocell, measured as the number of informational molecules inside the protocell. Until this point, protocell membranes have been assumed to grow at a rate which maintains protocell stability. MCS is modified in this chapter to permit variation in average protocell size across protocell strains. This variation in

protocell size emerges from the molecular dynamics which arise from the addition of a rudimentary computational function. Size then functions as the direct target of the externally imposed fitness function.

B.2 Overview of modifications to MCS

The changes proposed in this chapter are motivated by a desire to observe some form of minimal computation in MCS. In the hierarchical experiments presented in Chapters 5 and 6, protocells were free to grow until a certain molecular population threshold, S_{max} , was reached. In this chapter, a notional protocell fissioning factor, F , is introduced, and it is this protocell fissioning factor which triggers the binary fission of a parent protocell. For MCS-0 and MCS-1, the instantaneous amount of protocell fissioning factor in a given protocell is implicitly equivalent to the number of informational molecules within that protocell. In other words, it is assumed that protocell fissioning factor is produced upon each molecular replication, as a side reaction, and protocells undergo binary fission when the protocell fissioning factor reaches an appropriate level. The proposed modification at this point is to break this implicit coupling between molecular population and protocell fissioning factor, and implement a mechanism where *only certain* molecular replications give rise to the side reaction which produces the fissioning factor. The exact details of this mechanism will follow later (Section B.3.1).

Since we are using a feature of the bit-sequence of the molecules involved to determine whether a particular interaction affects fissioning, this new feature will necessarily be mutable during the replication process. Given these modifications, it should be possible to evolve populations of protocells against a particular fitness function so that the internal molecular composition of these protocells evolves to form a certain organisational structure that is conducive to producing protocells with the highest fitness. The extended molecular function that will be implemented can be seen as the ability for molecules to support rudimentary computation. In particular, a basic counting function will be enabled so that the effects of counting at the molecular level can be propagated to the protocell level, via the protocell fissioning factor, F . This kind of counting behaviour is somewhat similar to the “sequential counting behaviours” observed in the simulated Gene Regulatory Networks presented by Bentley (2004) though in this case, it is considerably simplified. Applying artificial selection at the protocell level might then have an indirect effect on the organisation of the internal molecular dynamics of the protocells. Conversely, if even this highly idealised test were unsuccessful, then it might raise significant doubt on the short-term potential for wet-lab artificial protocell research to demonstrate

evolution of more general computation in protocells.

B.2.1 Relevance to Wet-Lab Protocells

The coupling mechanism that will be applied in MCS-2 is not based on any concrete mechanism that might be used in real wet-lab protocells, but the fact remains that *some* real mechanism of achieving such coupling will be necessary in real protocells if selection is to be applied at the protocell level. In this case, the mechanism has been chosen ad-hoc in order to *initiate* an investigation of whether, or how, the overall evolutionary dynamics will then work. In other words, to provide at least a proof of principle that such directed evolution of protocell-based molecular computation might be possible.

B.2.2 Variability of Fission Condition

The key difference between MCS-2 and the previous versions of MCS is the breaking of the direct coupling between the number of informational molecules and the fissioning of a protocell. This will be achieved using an identifying region within molecules which determines whether or not reactions involving that molecule will give rise to the side-reaction which increases the quantity of fissioning factor, F . The MCS modifications proposed here will enable this region to be in one of a number of states, one of which causes an increase in protocell fissioning factor when that molecule is created as the product in a replication event. The configuration or state of this region will be determined during molecular replication: specifically, this state never changes during the lifetime of an individual molecule.

This identifying marker is active when molecules are in the substrate role. It is therefore represented in the primary string of a molecule, and will of course be subject to mutations during replication. Furthermore, we shall modify the molecular replication process to allow the catalyst to also deterministically modify the identifying region of the molecule that is produced. When a new molecule is produced as a result of the replication of an already existing molecule, the identifying region of the newly produced molecule will be modified according to a pre-defined scheme, explained in more detail in Section B.3.1. This ability to explicitly modify the variable identifying region of a newly produced molecule is also described as the *computational function* of the catalyst. It is worth mentioning at this point that the modifications to the molecular chemistry are made such that they do not otherwise change the molecular level dynamics that have been described in Chapters 5 and 6.

Upon successful replication then, the newly produced molecule will have one of the possible identifying regions, the precise one being determined by a particular

substring of the molecule that catalysed the replication event. Given two molecules, C and X , where C is the catalyst and X is the substrate, if C binds X , then the result is to replicate X and modify the configuration of its variable region according to the computational function specified by C . Then, if the variable region of the newly produced molecule has a particular pre-defined value, the amount of protocell fissioning factor in its containing protocell will be increased. The available computational functions are equivalent to simple counting (with wrap around at a variable modulus or *base*), the details of which will be provided later (Section B.3.1.2).

B.3 MCS-2: The Modifications

In the previous section, an overview of the modifications to be applied was presented. The modifications described involve changes to both the molecular level, in terms of the states of identifying markers, and the protocell level, in terms of protocell fissioning factor, F . As discussed in Section B.2, the end result of these modifications will be that MCS molecules can undergo an enzymatically mediated transformations which might be interpreted as very basic computational operations (Holland, 1975).

The following sections present the implementation details of the modifications which add this rudimentary “counting” function and will be followed by a set of experiments which seek to demonstrate an initial proof of principle that the behaviour of this molecular counting function can be affected by an explicit, externally applied fitness function at the protocell level. If this can be done then it provides a basis for investigation of more complex molecular computation, embedded in protocells, in the future. Conversely, if it cannot be done, even for such “minimal” kinds of computation, then this might indicate a significant potential barrier to the long-term application of protocell-embedded molecular computation in realising general purpose information and communication technologies.

B.3.1 Modifications to Molecular level chemistry

As previously described, the modifications required to the implementation of the MCS molecules are two-fold. Firstly, we shall set aside a portion of the primary bit-string to represent the variable computational data region (state) of the molecule, when it is acting as a substrate. The second modification to the MCS molecular implementation is to bestow upon catalyst molecules the ability to modify the computational data region of molecules that are produced by catalytic replication of a substrate and to devise a scheme by which this modification takes place. We therefore set aside a portion of the folded sequence to encode the computational function implemented by the molecule when it is acting as a catalyst.

10110011110000	
<i>Comp_{state}</i>	<i>Binding Target Region</i>
10	110011110000

Table B.1: Modified Molecular Structure (substrate role)

10110011110000		
<i>Comp_{state}</i>	<i>Comp_{func}</i>	<i>Remainder</i>
10	1100	11110000
-	<i>HL</i>	<i>HHLL</i>

Table B.2: Modified Molecular Structure (catalyst role), *Comp_{state}* ignored

B.3.1.1 Representing Computational State

For the MCS as presented during the remainder of this chapter, the (primary) structure of all molecules will have a variable region which will be in one of four different configurations or states. Since the primary structure is binary, we need to reserve $\log_2(4) = 2$ bits of the bit-string in order to represent four different states. These two reserved bits will not be included in any of the molecular binding mechanisms described in Sections 4.4.5 and 6.2—they will form neither the matching string nor the target string for pattern matching and molecules which differ only in their variable state region will be taken as members of the same molecular species. Since these bits are reserved for representing a *computational state*, they will henceforth be known as the *Comp_{state}* region, where “Comp” is an abbreviation for “computational”, as shown in Table B.1.

B.3.1.2 Representing Computational Function

A further modification to the MCS is the reservation of a portion of the secondary, folded structure, which will serve to identify the particular computational function of that molecule. This function is only active when the molecule plays the role of “catalyst” in an interaction. Since these (secondary alphabet) bits are reserved to denote the manner in which this molecule may modify the state of a newly created molecule, they will henceforth be known as the *Comp_{func}* region, where “func” is an abbreviation for “function”, as shown in Table B.2.

B.3.1.3 Remaining Bits

Table B.2 shows an MCS molecule segmented into 3 components: *Comp_{state}*, *Comp_{func}* and *Remainder*. These *Remainder* (secondary alphabet) bits form the matching pattern used when this molecule is in the catalyst role in the same way as described

earlier in Chapters 5 and 6, where the matching takes place against the “Binding Region” of the substrate, as shown in Table B.1.

B.3.2 Modifications to Protocell level chemistry

Until now, protocell fissioning in MCS was triggered when the number of contained molecules reached a certain fixed threshold, S_{max} . In MCS-2, we need the ability to make fissioning (and thus average protocell size) vary depending on some feature of the underlying (computational) molecular dynamics. We do this by introducing a separate “fissioning factor” molecule whose instantaneous quantity is represented by F . Protocell fission will now take place when this fissioning factor, F , reaches a pre-defined threshold, F_{max} . MCS-0 and MCS-1 can then be retrospectively described as having implicitly assumed that the amount of protocell fissioning factor in a particular protocell will grow in tandem with the replication of molecules inside the protocell, $F = S$. If we take it that a side-reaction from molecular replication causes this increase in F , then in MCS-0 and MCS-1, this side-reaction could be said to accompany *every* replication, while in MCS-2 it will accompany only *certain* replications. Upon protocell fission, the total number of fissioning factor molecules of the parent cell will be distributed equally between the two daughter cells. The amount of protocell fissioning factor then varies between $\frac{F_{max}}{2}$ at birth and F_{max} at fission. These fissioning factor molecules are considered to be completely disjoint from the informational molecules in MCS—they are all identical and have no reaction dynamics in themselves except to trigger protocell fissioning when they reach some threshold number. For this reason, we do not have to represent or track individual fissioning factor molecules, but only the aggregate count of such molecules within each protocell. The actual change to the protocell implementation is therefore to only increase the number of fissioning factor molecules when an informational molecule is produced (by replication of some substrate molecule) which has a particular computational state, the details of which are presented later in Section B.3.3.3. In this way, different strains of protocells should grow and divide at different rates due to their varying computational functionality of their internal molecular compositions.

B.3.3 Modifications to the MCS Reaction Algorithm

In Sections B.3.1.1, B.3.1.2 and B.3.2, the precise *structural* details of the proposed modifications were described. What remains is to describe how these structural changes to the molecules and protocells will be incorporated into MCS, and to ensure that the modifications maintain compatibility with the experimental results

$Comp_{func}$	Value
LL	4_{10}
LH	1_{10}
HL	2_{10}
HH	3_{10}

Table B.3: Converting $Comp_{func}$ to a specific computational function (counting base)

presented in Chapters 5 and 6.

B.3.3.1 Updating Computational State

As was the case in MCS-0 and MCS-1, we continue to take two molecules to serve as the catalyst and the substrate respectively. Pattern matching proceeds as before with the exception of the previously mentioned reserved regions (Sections B.3.1.1 and B.3.1.2). If the catalyst can bind to the substrate, then replication proceeds with all primary structure bits, comprising $Comp_{state}$ and the target binding region (incorporating the encoded Res_{mod} region), being subject to a per-bit mutation rate as before. Once the new molecule is produced, we apply a mechanism so that the catalyst can deterministically update the variable computational state region, $Comp_{state}$, of the new molecule. This takes place by examining the $Comp_{func}$ secondary structure region of the catalyst molecule. The $Comp_{func}$ region is taken to represent a binary coded number. This number will be treated as specifying the computational function of that particular catalyst. Upon the replication of a substrate, the $Comp_{state}$ region of the newly produced molecule is also treated as a binary coded number. The computational state of the newly produced molecule is updated by incrementing the value that was previously contained within its $Comp_{state}$ region and then taking the remainder of this when divided by the value represented by the $Comp_{func}$ region of the catalyst molecule and derived from Table B.3. That is, the computational function is always simple modular counting, parametrised by a base specified by the catalyst molecule, and applied to the variable computational state region of the produced molecule.

B.3.3.2 Updating Computational State: A Worked Example

In the example interaction shown in Table B.4, the substrate molecule and the catalyst molecule are identical with the exception of their $Comp_{state}$ regions: the substrate has 00_2 as its $Comp_{state}$ region, and the catalyst has 11_2 as its $Comp_{state}$ region. These two molecules are therefore members of the same molecular species, since their primary structures differ only in the computational state region, $Comp_{state}$.

substrate : 0010101111		
<i>Comp_{state}</i>	<i>Binding Region</i>	
00	10101111	

catalyst : 1110101111		
<i>Comp_{state}</i>	<i>Comp_{func}</i>	<i>Remainder</i>
-	1010	1111
-	<i>HH</i>	<i>HH</i>

IntermediateMolecule: 0010101111		
<i>Comp_{state}</i>	<i>Binding Region</i>	
00	10101111	

Table B.4: Example Interaction, before State Update

To determine whether a binding may occur between these two individual molecules (and thus whether this molecular species is a self-replicase) we need to fold the catalyst into its secondary structure. For all experiments in this chapter, we use the same folding scheme as previously introduced in Table 6.2. When folded according to that scheme the catalyst becomes *HHHH*. Pattern matching then takes place on the substrate molecule having disregarded its *Comp_{state}* region. In this case, we attempt to match the “*Remainder*” region of the catalyst, *HH* (which is 11 by the transformations in Table 6.6) against the bit-string 10101111 (which is the substrate minus its *Comp_{state}* region). It is clear that pattern matching succeeds and therefore these two molecules can bind. Error-prone bit-wise molecular replication then takes place as described in Section 4.4.5.2, with all bits of the substrate (ie. *Comp_{state}* and *Binding Region*) being processed subject to the mutation parameter. In the first instance, let us assume that the substrate has been replicated with complete fidelity so that the newly produced molecule is that shown in Table B.4 and labeled “*IntermediateMolecule*”.

The next step in the molecular replication process is to update the *Comp_{state}* region of the new molecule according to the mechanism described earlier in this section before it is added to the protocell. In the situation where mutations occur during replication, the mutations are applied *before* the catalyst modifies the *Comp_{state}* region of the new molecule. To carry out this update, the *Comp_{func}* region of the catalyst is parsed from its secondary structure and its value is determined by Table B.3.

1. The value of the *Comp_{func}* region of the catalyst molecule is interpreted from Table B.3. For the example catalyst in Table B.4, we get a value of $HH = 3_{10}$.
2. The *Comp_{state}* region of the newly produced molecule is also taken as a binary coded number. In the example “*IntermediateMolecule*” presented in Table B.4, this is $00_2 = 0_{10}$. This value is incremented modulo the value of the *Comp_{func}*

NewMolecule: 0110101111	
<i>Comp_{state}</i>	<i>Binding Region</i>
01	10101111

Table B.5: Molecule produced from interaction in Table B.4, after molecular state update

region of the enzyme, 3_{10} .

$$\begin{aligned} 0 + 1 &= 1 \\ 1 \bmod 3 &= 1 \end{aligned}$$

So, the value that is assigned to the *Comp_{state}* region of the newly created molecule is 1, which is represented by the bit-string 01₂. The final structure of the molecule that gets added into the protocell population is therefore that shown in Table B.5.

B.3.3.3 Fissioning Factor Molecule Count Increase

In the previous section, the modifications to the molecular replication algorithm were presented. The final modification to MCS to be described is the mechanism for selectively increasing the number of protocell fissioning factor. As described in Section B.3.2, the side reaction which adds a new fissioning factor molecule depends on the final value stored in the *Comp_{state}* region of the newly created molecule. For the purposes of the work presented here, a new fissioning factor molecule is added if the *Comp_{state}* region of the newly created molecule is 0₁₀ (ie. its bit-string is 00₂) In all other cases, the new informational molecule is added to the protocell *without* causing the side reaction which adds a new fissioning factor molecule. In the case of the example molecule produced in Section B.3.3.2, the value of the *Comp_{state}* region was 1 (bit-string 01₂), so that molecule would not cause a new fissioning factor molecule to be added to the protocell. The new protocell variable, *F*, which was introduced in Section B.3.2 to represent the number of fissioning factor molecules in a given protocell, will therefore be incremented each time a new informational molecule is added to that protocell, given that the new molecule meets the appropriate *Comp_{state}* conditions laid out above.

B.3.3.4 Dynamics of the *Comp_{state}*/ *Comp_{func}* interaction

The previous two sections have outlined the mechanism by which the *Comp_{state}* and *Comp_{func}* regions are processed. In this section, the dynamics produced by the

interactions between the various possible states of $Comp_{state}$ (in the *substrate*) and $Comp_{func}$ (in the *catalyst*) are further explored.

The $Comp_{state}$ region may take any of the following values: $\{00_2, 01_2, 10_2, 11_2\}$, representing the decimal numbers from 0_{10} to 3_{10} inclusive. In the absence of active modification by the $Comp_{func}$ dynamics then, molecules would be expected to occur with uniform distribution across each of the possible $Comp_{state}$ region values due to mutations during replication, and the side-reaction which produces new fissioning factor molecules would happen, on average, once in every four replications.

On each replication event, the value of the $Comp_{state}$ of the new molecule is determined by incrementing the $Comp_{state}$ of the substrate, with a wrap-around from 3_{10} back to 0_{10} and taking the remainder when divided by the $Comp_{func}$ value of the catalyst. Taking these $Comp_{func}$ dynamics into account, however, gives rise to restrictions in the production rates of certain $Comp_{state}$ regions. z_i will be used to refer to the value of the $Comp_{func}$ region of species i from now on, and is derived from Table B.3 as before. For example, a $Comp_{func}$ of $LL(z_i = 4_{10})$ will cause the production of all possible $Comp_{state}$ regions, hence the side-reaction which produces new fissioning factor molecules would occur approximately, neglecting mutations, once every four replications:

$$\begin{aligned} 0 + 1 &= 1 \\ 1 \bmod 4 &= 1_{10} \\ Res_{state} &= 01_2 \end{aligned}$$

$$\begin{aligned} 1 + 1 &= 2 \\ 2 \bmod 4 &= 2_{10} \\ Res_{state} &= 10_2 \end{aligned}$$

$$\begin{aligned} 2 + 1 &= 3 \\ 3 \bmod 4 &= 3_{10} \\ Res_{state} &= 11_2 \end{aligned}$$

$$\begin{aligned}
3 + 1 &= 4 \\
4 \bmod 4 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

So a $Comp_{func}$ of LL can produce molecules with $Comp_{state}$ regions in the range $\{00_2, 01_2, 10_2, 11_2\}$. If the $Comp_{func}$ is $LH(z_i = 1_{10})$, however, the range of $Comp_{state}$ values that will be produced is reduced to the set $\{00_2\}$, and a new fissioning factor molecule would be added after *every* replication:

$$\begin{aligned}
0 + 1 &= 1 \\
1 \bmod 1 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

$$\begin{aligned}
1 + 1 &= 2 \\
2 \bmod 1 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

$$\begin{aligned}
2 + 1 &= 3 \\
3 \bmod 1 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

$$\begin{aligned}
3 + 1 &= 4 \\
4 \bmod 1 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

If the $Comp_{func}$ is $HL(z_i = 2_{10})$, the range of $Comp_{state}$ values that will be produced is $\{00_2, 01_2\}$ and the side-reaction which produces new fissioning factor molecules would occur approximately once every two replications:

$$\begin{aligned}
0 + 1 &= 1 \\
1 \bmod 2 &= 1_{10} \\
Res_{state} &= 01_2
\end{aligned}$$

$$\begin{aligned}
1 + 1 &= 2 \\
2 \bmod 2 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

$$\begin{aligned}
2 + 1 &= 3 \\
3 \bmod 2 &= 1_{10} \\
Res_{state} &= 01_2
\end{aligned}$$

$$\begin{aligned}
3 + 1 &= 4 \\
4 \bmod 2 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

If the $Comp_{func}$ value is $HH(z_i = 3_{10})$, the range of $Comp_{state}$ values that will be produced is $\{00, 01, 10\}$. The side-reaction which increases F will therefore occur approximately once every three replications due the distribution of available $Comp_{state}$ values in the molecular population:

$$\begin{aligned}
0 + 1 &= 1 \\
1 \bmod 3 &= 1_{10} \\
Res_{state} &= 01_2
\end{aligned}$$

$$\begin{aligned}
1 + 1 &= 2 \\
2 \bmod 3 &= 2_{10} \\
Res_{state} &= 10_2
\end{aligned}$$

$$\begin{aligned}
2 + 1 &= 3 \\
3 \bmod 3 &= 0_{10} \\
Res_{state} &= 00_2
\end{aligned}$$

$$\begin{aligned}
3 + 1 &= 4 \\
4 \bmod 3 &= 1_{10} \\
Res_{state} &= 01_2
\end{aligned}$$

We can see that a $Comp_{func}$ of $LH(z_i = 1_{10})$, produces the $Comp_{state}$ value of 00_2 upon every replication and would therefore be expected to produce qualitatively similar results to the earlier versions of MCS presented in Chapters 5 and 6, since F , the number of protocell fissioning factor molecules will remain exactly synchronised with S , the number of informational molecules. Similarly, molecules with a $Comp_{func}$ of $HL(z_i = 2_{10})$ will produce the $Comp_{state}$ value of 00_2 upon every *other* replication in the absence of mutation. F would therefore be expected to increase at half the rate of S . In Section B.3.2, it was stated that the number of protocell fissioning factor molecules in a protocell would vary between $\frac{F_{max}}{2}$ at birth and F_{max} at fission. In order to reach the fission threshold, F_{max} , an individual protocell must generate $\frac{F_{max}}{2}$ fissioning factor molecules. Neglecting the effects of mutation, we can generalise this rule to say that protocells which are homogeneously (self-replicate) species i will require $z_i(\frac{F_{max}}{2})$ molecular replications to generate $\frac{F_{max}}{2}$ fissioning factor molecules and reach the threshold of F_{max} . Table B.6 shows how S varies for homogeneous protocells where $F_{max} = 100$ and the effects of molecular mutation are neglected.

As we will see in the next section, these differences in the expected protocell size distributions, and therefore lifecycle durations, can now be used to design a protocell level fitness function that favors a certain molecular level computation function.

$Comp_{func}$		S_{min}	S_{max}
LH	$z_i = 1_{10}$	50	100
HL	$z_i = 2_{10}$	100	200
HH	$z_i = 3_{10}$	150	300
LL	$z_i = 4_{10}$	200	400

Table B.6: Range of informational molecule counts for homogeneous protocells, $F_{max} = 100$

B.3.3.5 Fitness Function Design—Artificial Selection

Given the dynamics presented in Section B.3.3.4, and assuming that these internal dynamics can be reliably observed at the protocell level, it should be possible to design a fitness function that can impose selectional pressure on protocell level traits that are phenotypic expressions of the underlying molecular dynamics. In this section, a fitness function will be described that will enable the demonstration of the directed evolution of molecular populations with respect to the $Comp_{state}$ and $Comp_{func}$ regions of the dominant molecular species of a given protocell lineage.

In Chapters 5 and 6, the size of a protocell, measured as the number of informational molecules it contains, was directly linked to the condition under which it would reproduce or fission. The modifications to the protocell chemistry described in Sections B.3.2 and B.3.3.3 break this direct link between the number of informational molecules and the fission event, and replace it with a conditional link. Protocells will then vary in size during their lifecycle depending on the nature of the internal molecular sub-populations they contain, as shown in Table B.6. Note that, by design, the variation in number of fissioning factor molecules over the protocell life-cycle is identical for all protocell strains, ranging from $\frac{F_{max}}{2}$ at birth to F_{max} just before the protocell divides by binary fission.

B.3.3.6 Molecular Activity Factor

Protocell size then, as measured by the number of informational molecules within the protocell (S), is a suitable protocell level target for selection to encourage the directed evolution of the internal molecular chemistry of protocells. Through the use of an arbitrary mapping function, we can calculate a “molecular activity factor”, $g \in [0, 1]$, which is now made a function of S and can be taken as a probability of molecular reaction occurring for that protocell, thereby influencing the reaction rate of protocells based on their physical size. An analogy could be made here between this system and, perhaps, a cell which gathers and exploits light energy at a rate which is proportional to the size of the cell—bigger cells might be better at gathering energy due to the higher surface area presented to the energy source.

This molecular activity factor is implemented as an extension to the molecular interaction algorithm presented in Sections 4.4.5 and 6.2.2 and applies at the point immediately following the selection of two molecules as substrate and catalyst respectively.

- the instantaneous count of the number of informational molecules in a protocell, S , serves as a parameter to a function, g , which is essentially a biased coin toss
- If $g(S)$ returns the value 0, then the interaction is considered elastic, otherwise, the interaction proceeds as previously described in Section 6.2.2 and the product is subject to the $Comp_{state} / Comp_{func}$ interaction described above in Section B.3.3.4.

Figure B.1 shows the range of the activity factor which is used for the experiments presented later in this chapter. It has been set up to favour protocells whose dominant molecular species has a $Comp_{func}$ of $HH(z_i = 3_{10})$. Such protocells will typically exist with instantaneous S values between 150 and 300, as shown in Table B.6, though this will also be affected by molecular mutation rates and the resultant mutational load. Furthermore, the S values for $(z_i = 4_{10})$ and $(z_i = 2_{10})$ protocells overlap with the $(z_i = 3_{10})$ range in parts, and will therefore benefit from increased molecular activity for *some* of the time.

B.3.3.7 Effect of Imposed Molecular Activity Factor

The expected gestation period for protocells under the newly applied molecular activity function, neglecting mutations during molecular replication, can be calculated as follows:

$$\begin{aligned}\frac{dS}{dt} &= Sg(S) \\ dt &= \frac{dS}{Sg(S)} \\ T &= \int_{S_{min}}^{S_{max}} \frac{dS}{Sg(S)}\end{aligned}$$

From this we can calculate the average protocell birth rate as $\frac{1}{T}$. The average cell size over its life-cycle will be given by:

$$\bar{S} = \frac{1}{T} \int_0^T S dt$$

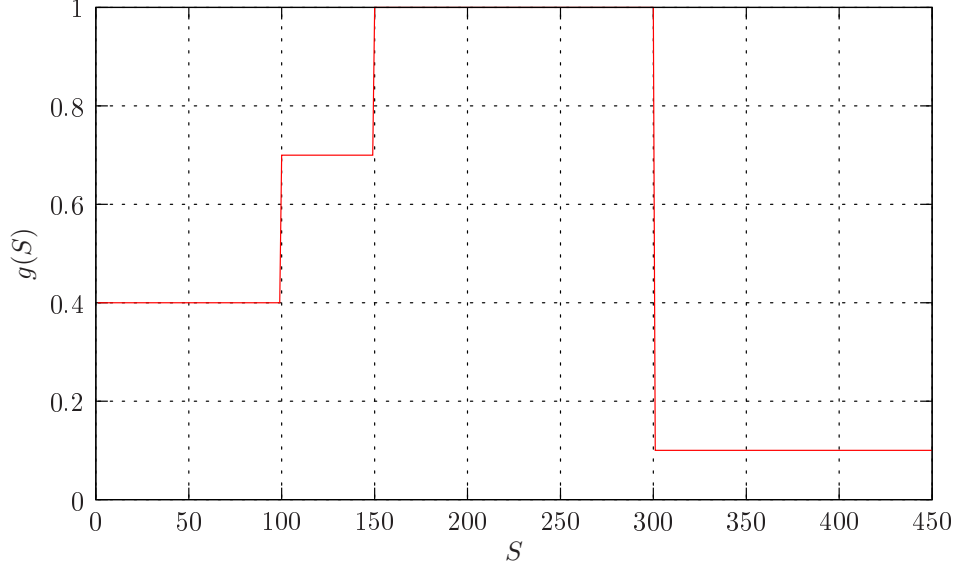


Figure B.1: Molecular activity factor, $g(S)$, configured to favour $z_i = 3$

Given that the total number of informational molecules is subject to a fixed maximum, M_{max} (Section 4.4.2), the protocell population size will vary depending on the average size of protocells present. If populated by protocells of a single $Comp_{func}$ type, the expected population size will be given by $\frac{M_{max}}{S}$. Applying these derivations to the molecular activity factor from Figure B.3.3.6, and taking $M_{max} = 1.5 \times 10^4$, we can calculate (using numerical methods as required) a variety of nominal protocell parameters, as shown in Table B.7. Actual values will vary, due to mutational load at the molecular level, depending on:

1. the exact length of the particular dominant molecular species,
2. the characteristics of the spectrum of accessible mutants for the particular dominant molecular species

both of which contribute to the overall reproduction and mutation rates for protocells (Section 6.3.2.1).

B.4 MCS-2: The Experiments

In this final set of experiments, we will demonstrate the directed evolution of a rudimentary counting behaviour, from among the (small) repertoire of available behaviours, using the artificial protocell system that has been described thus far.

z_i	Gestation Time	Birth Rate	Relative Fitness	Mean Size	Mean Population
1	1.733	0.577	0.400	72	208
2	0.867	1.154	0.800	139	107
3	0.693	1.442	1.000	216	69
4	3.282	0.305	0.211	334	45

Table B.7: Nominal parameters for protocells of the four possible distinct counting base values

This will take the form of a modular counting function with a specific modulus (corresponding to a specific z_i / $Comp_{func}$ region of the principal molecular species of protocell lineages).

B.4.1 Evolution of Protocell-embedded molecular computation

In the following experiment, the molecular activity factor, $g(S)$, has been set up as shown in Figure B.1. Each of the experiments will be initialised with protocell strains which have a principal molecular species with $z_i = 4$, which yields the lowest possible fitness for a protocell lineage according to the molecular activity factor. We can therefore expect that evolution will drive this initial population towards higher fitness protocell strains, specifically those whose principal molecular species has $z_i = 3$. As the data in Table B.7 show, the relative fitnesses between protocells subject to this molecular activity factor can be expressed as follows: $(z_i = 4) < (z_i = 1) < (z_i = 2) < (z_i = 3)$.

- $v = 0.01$; $\beta = 0.4$
- $M_{max} = 15000$; $S_{max} = 100$
- An initialiser molecular species of length $n = 20$ is randomly chosen, with the added conditions that:
 1. this initialiser species is capable of self-replication (produced by an algorithm similar to that described in Section 6.2.2.1), and
 2. this initialiser species has an $Comp_{func}$ of $LL(z_i = 4_{10})$, which produces the least fit protocell strains per the molecular activity factor shown in Figure B.1
- Cells are repeatedly initialised with normalised molecular populations of the initialiser species until the global $maxMolecules$ threshold is reached. The initialisation algorithm respects the $Comp_{func}$ region of the initialiser molecule and guarantees that cells initially have an appropriate level of protocell fissioning factor, F .
- The process for initialising individual protocells ensures that at $t = 0$, the protocell population is in an approximately normalised state. This means that protocell sizes, measured as total number of molecules, will be uniformly random, and that within protocells, a suitable mutational load is present, and furthermore, that the amount of protocell fissioning factor, F , in protocells is appropriate for the $Comp_{func}$ region of their principal molecular species and size in terms of S , the number of molecules in the protocell.
- Repeated molecular interactions, as described in Section 4.4.5, are then carried out, including the appropriate $Comp_{state} / Comp_{func}$ dynamics as described in Section B.3.
- For clarity in presenting and analysing the results, the class-6 parasite threshold for determining whether a protocell should be labeled “mixed” or not (Section 5.1.4.4) has been increased to 10 here. Experiments using a lower threshold led to the majority of protocells being labeled “mixed”. This is conjectured to be due to the fact that MCS-2 catalyst molecules fold into shorter binding patterns than MCS-1 catalysts of the same overall length, due to the reserved regions in MCS-2. This shorter binding pattern makes it more likely that catalysts in MCS-2 can bind to substrates, and thus, MCS-2 protocells may have a generally higher number of class-6 parasites of their dominant molecules.

B.4.1.1 2-Step Evolution of Optimal z_i

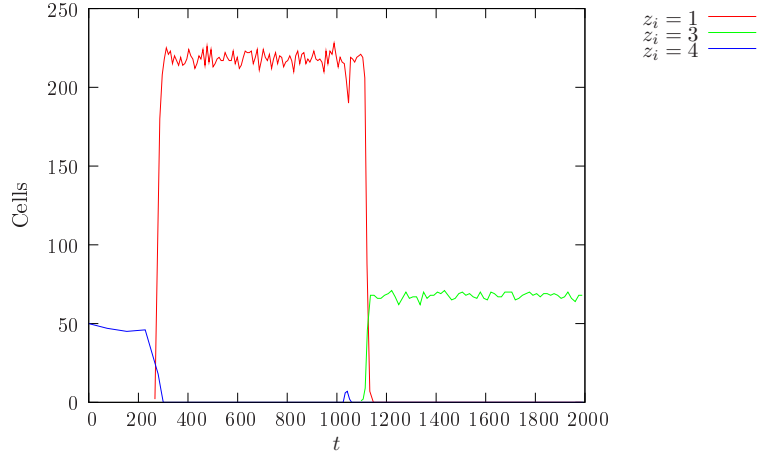


Figure B.2: Successful evolution of protocell-embedded 2-bit molecular counting via 2 transitions. Numbers of protocells, classified by counting base (z_i) of principal molecular species in each protocell.

Figures B.2 and B.3 show the outcome of an experimental run where the optimal $z_i = 3$ was reached after two population level fitness transitions. The run was initialised with all protocells having $z_i = 4$, hence the lowest possible fitness with respect to the imposed molecular activity factor (Section B.3.3.6). For clarity, Figure B.3 was plotted using a minimum threshold for cell populations and only shows cell lineages which reached at least 10 cells in population.

Figure B.2 shows the protocell population throughout this run with protocells grouped by the z_i value of the principal molecule. Note how the actual number of protocells depends on the predominant z_i values of the protocells present (Table B.7). The global protocell reactor (Section 4.4.2) has a fixed capacity in terms of the number of informational molecules it can contain. Given the range of expected protocell sizes shown in Table B.6, it would therefore be expected that the total number of protocells would vary depending on the predominant z_i (Table B.7). The significant events in Figure B.2 are:

- At $t \simeq 260$, the predominant z_i of the protocell population shifted to $z_i = 1$
- At $t \simeq 1100$, the predominant z_i of the protocell population shifted to $z_i = 3$, the optimal z_i value with respect to the imposed molecular activity factor
- Thereafter, the population remained dominated by protocells with the optimal z_i

Later, in Section B.4.1.2, we will see a run which required three fitness transitions to reach the optimal z_i , though it should be noted that this is not typical of MCS-2

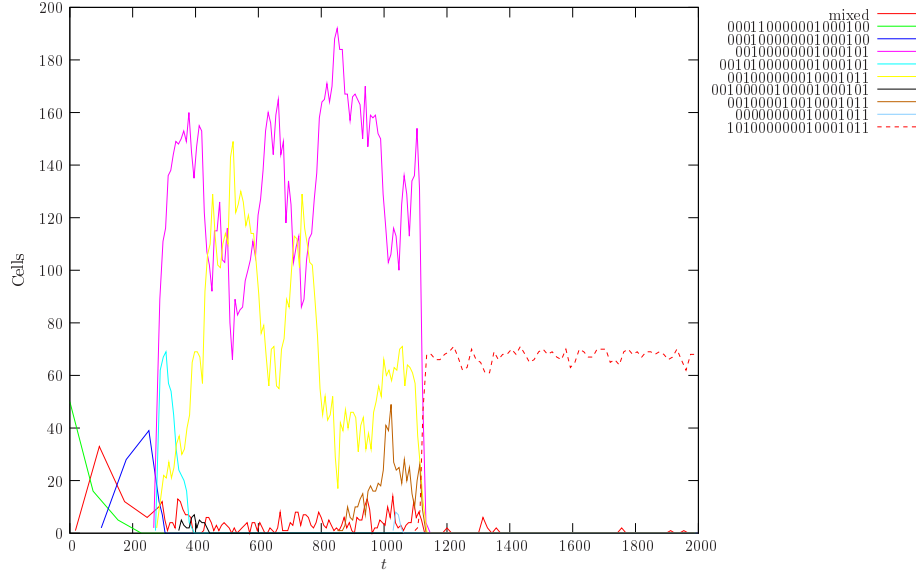


Figure B.3: Successful evolution of protocell-embedded 2-bit molecular counting via 2 transitions. Numbers of protocells, classified by principal molecular species in each protocell.

runs and is particularly sensitive to initial conditions. For the most part, MCS-2 experiments reach optimal fitness after one or two fitness transitions as is the case with this set.

Figure B.3 shows the protocell population throughout this run with protocells grouped by the molecular sequence of the principal molecule. It also shows the population of “mixed” cells—those which have more than 10 instances of molecules which are class-6 facultative parasites of the principal molecular species. The significant events in Figure B.3 are:

- Between $t = 0$ and $t \simeq 200$, the initialiser molecular species diversifies into a sub-population of “mixed” cells due to parasitic invasion at the molecular level.
- At $t \simeq 100$, a lineage emerges from the mixed population whose principal molecular species has the same z_i as the initialiser molecular species. This new lineage, as expected, has the “class 6” facultative parasite relationship with the previous lineage.
- At $t \simeq 250$, three protocell lineages emerge which have the “fitter” $z_i = 1$. Each of these three lineages has the expected class-6 relationship with the previous dominant lineage, suggesting that this is an example of three independent “discoveries” of the fitter $z_i = 1$. The lengths of the principal molecules in these three new lineages is unique in each case—relatively, $n - 1$, n , and $n + 1$, where

n is the length of the principal molecular species of the previous dominant lineage. From Sections 5.1.3 and 6.3.2.1, we know that the relative lengths of the principal molecular species of protocell lineages can have a marked effect on their relative fitnesses, though a quantification of this difference is significantly more difficult in MCS-2. It is worth noting though that the lineage which is displaced at $t \simeq 400$ is the $n + 1$ lineage, and that the lineage which peaks to > 180 cells at $t \simeq 850$ is the $n - 1$ lineage, highlighting again the qualitative difference in relative fitnesses due to dominant molecular length.

- At $t \simeq 850$, a protocell lineage emerges and peaks at $t \simeq 1000$. This lineage again had the sub-optimal $z_i = 1$, and emerged as a class-6 facultative parasite from the $n - 1$ lineage described above.
- At $t \simeq 1150$, a protocell lineage with the optimal $z_i = 3$ emerged. This lineage emerged as a class-6 facultative parasite from the $n - 1$ lineage, and proceeded to dominate for the remainder of the run. This period of dominance is accompanied by a significant lack of “mixed” cells, suggesting that this lineage is relatively safe from invasion by other lineages with $z_i = 3$.

B.4.1.2 3-Step Evolution of Optimal z_i

Figure B.4 shows the results of a run in which the optimal $z_i = 3$ was reached after three fitness transitions. This result is very sensitive to initial conditions and further attempts to reproduce the result with different conditions all resulted in runs which reached optimal z_i after 1 or 2 fitness transitions. The overall length of this run is also significantly longer than the run presented in Section B.4.1.1, by a factor of about 10. This run was again initialised with all protocells having $z_i = 4$, the lowest possible fitness with respect to the imposed molecular activity factor (Section B.3.3.6). Figure B.4 shows the protocell population throughout this run with protocells grouped by the z_i value of the dominant molecule.

The significant events in Figure B.4 are:

- At $t \simeq 500$, the predominant z_i of the protocell population shifted to $z_i = 1$
- At $t \simeq 1000$, the predominant z_i of the protocell population shifted to $z_i = 2$
- At $t \simeq 8000$ and again between $t \simeq 12000$ and $t \simeq 13500$, some protocells appeared which had optimum $z_i = 3$, though none of these were able to successfully invade the population
- At $t \simeq 16000$, the predominant z_i of the protocell population shifted to $z_i = 3$, the optimal z_i value with respect to the imposed molecular activity factor

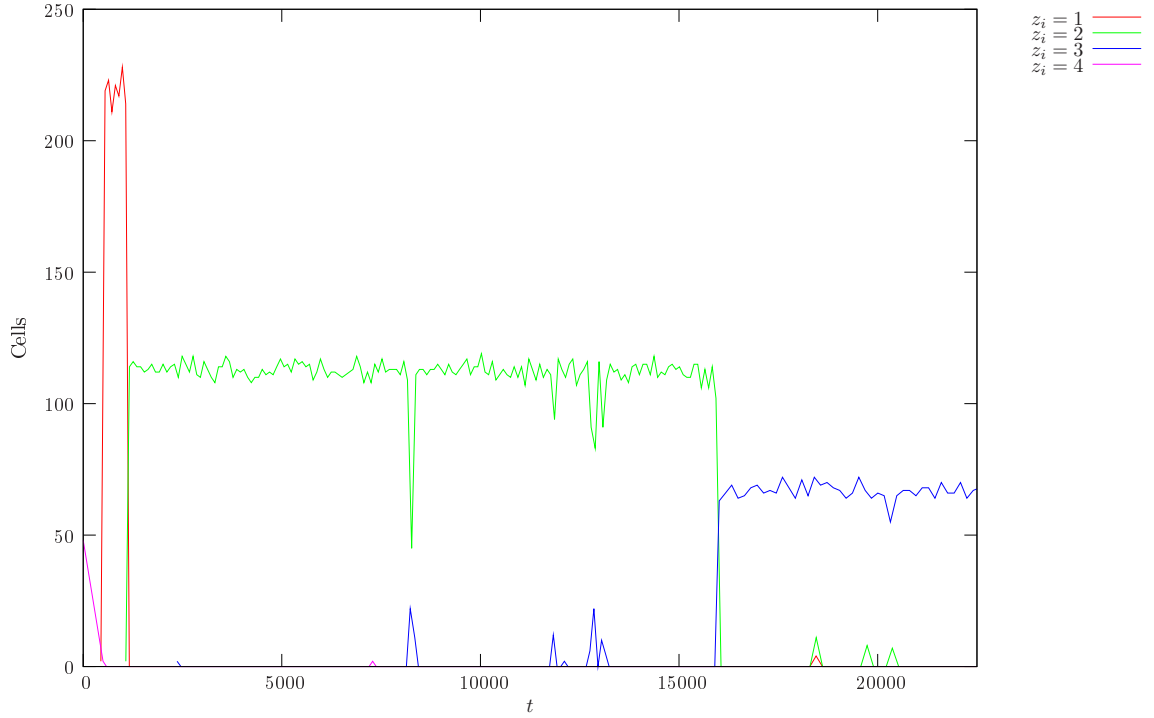


Figure B.4: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Numbers of protocells, classified by counting base (z_i) of principal molecular species in each protocell.

- Thereafter, the population remained dominated by protocells with the optimal z_i

During this run, a total of 214 independent protocell lineages were founded. Figure B.5 shows a plot over the entire run of the protocells grouped by the molecular sequence of their principal molecular species. This plot was produced using a minimum threshold for protocell populations. Due to the length of the run and number of species involved, a single such threshold is insufficient to present all of the data in a meaningful way.

B.5 Detailed Analysis of Figure B.4

During the experimental run presented in Section B.4.1.2, a total of 214 independent protocell lineages were founded. Due to the length of the run and number of species involved, a single figure using a particular plotting threshold for the protocell population is insufficient to present all of the data in a meaningful way. Here, multiple sub-plots are presented to highlight the events that occurred in more detail.

Below we highlight the events in figures B.6, B.7, B.8, B.9 and B.10.

The events in Figure B.6 are:

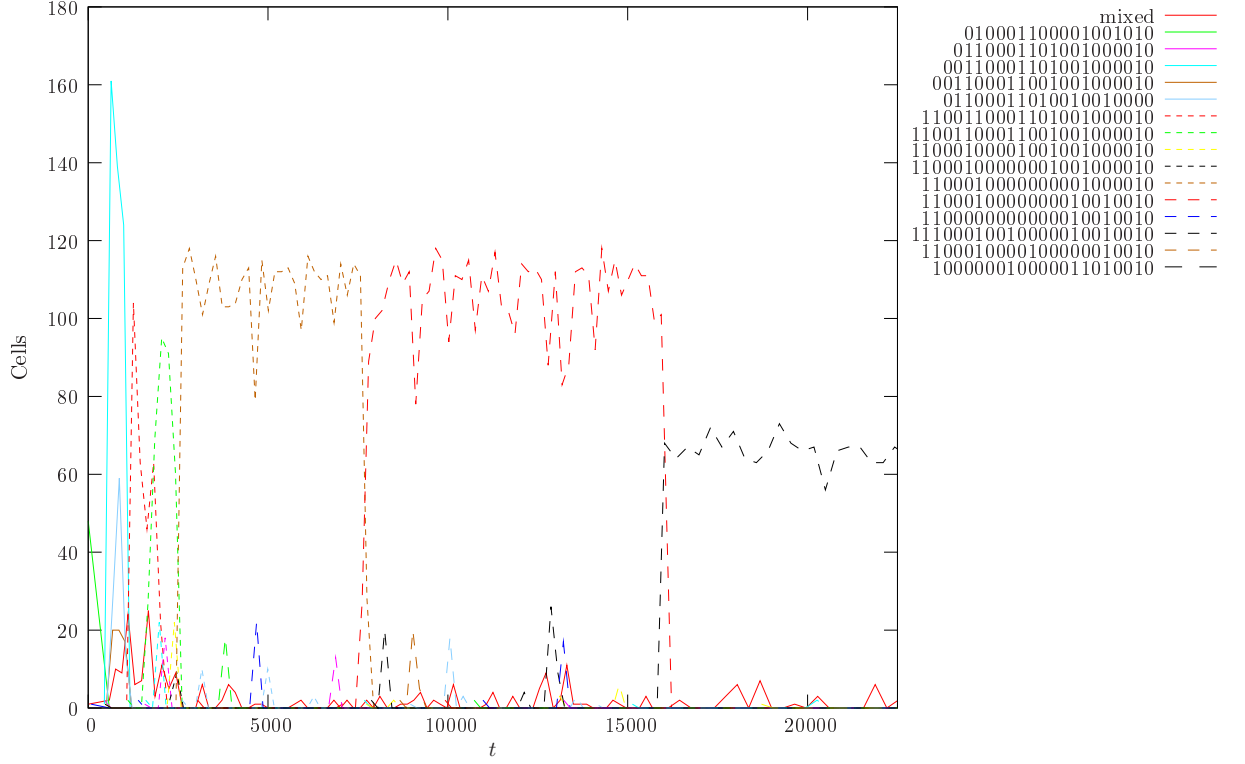


Figure B.5: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

- At $t \simeq 200$, the initialiser lineage is displaced by another lineage with the same sub-optimal $z_i = 0$.
- At $t \simeq 420$, four protocell lineages emerge with $z_i = 1$. As expected, the principal molecular species of each new lineage shares the pairwise class-6 facultative parasite relationship with the principal species of the lineage that was dominant before. This fitness transition was followed by a period of diversity between protocell lineages of approximately the same fitness. Molecular sequence length and the neighbourhood of accessible mutants will be a significant determining factor for the dynamics between these lineages (Section 5.1.3 and 6.3.2.1).
- At $t \simeq 1100$, a protocell lineage emerged with $z_i = 2$. This lineage was a class-6 facultative parasite of the protocell lineage which reached a peak population of ≈ 185 at $t \simeq 1050$.
- The decline in the protocell population from $t \simeq 1700$ onwards is due to overlap in the plots.

The events in Figure B.7 are:

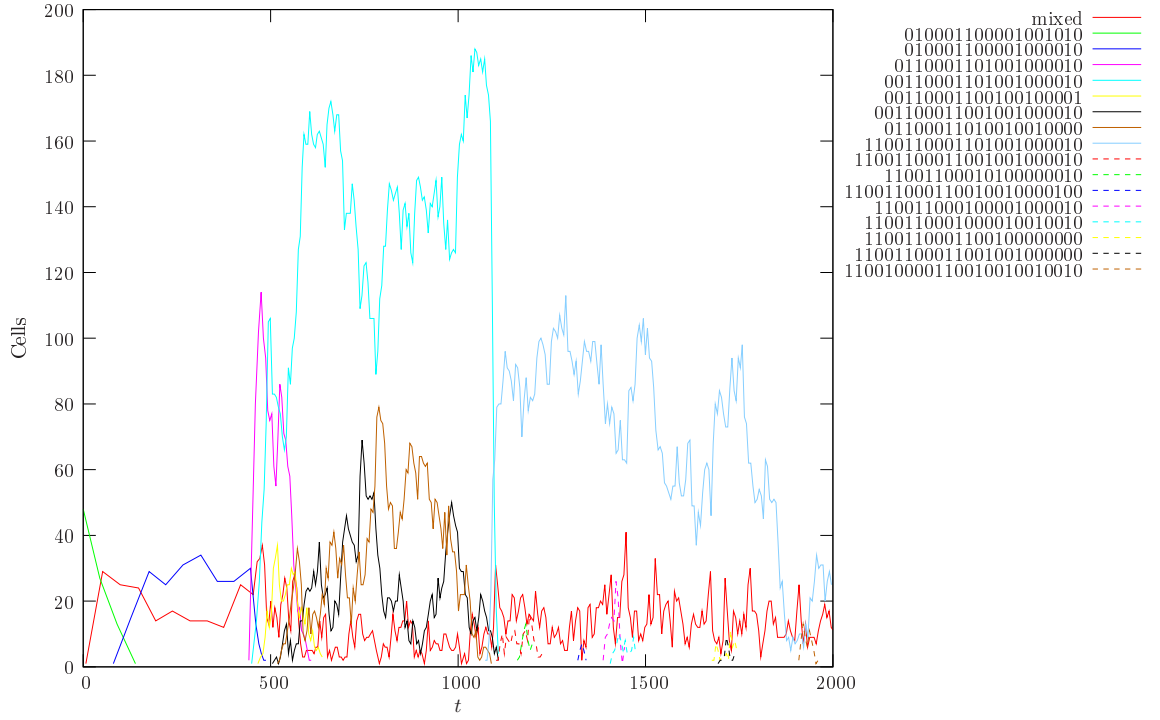


Figure B.6: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

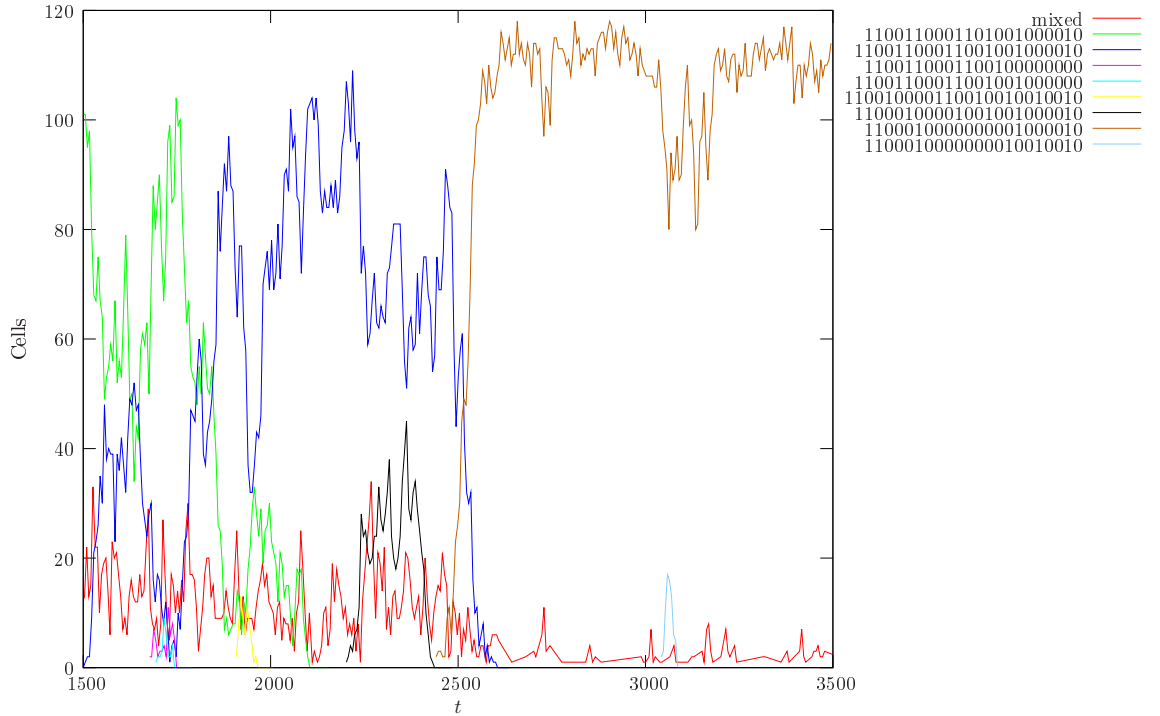


Figure B.7: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

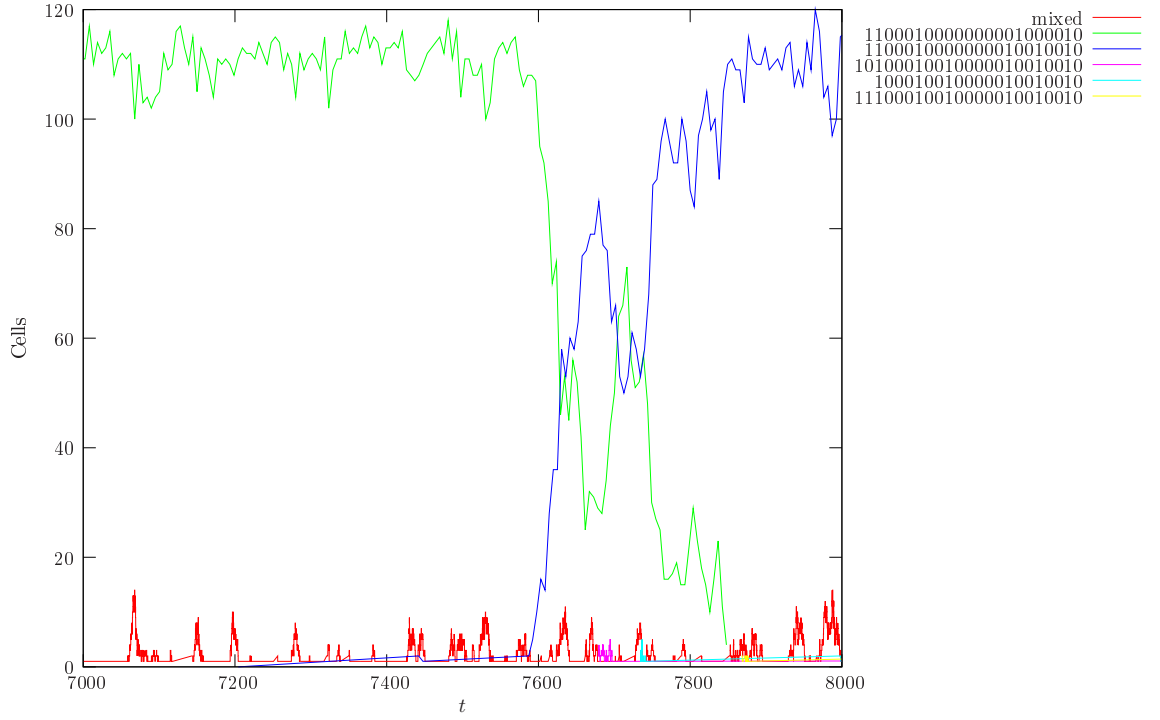


Figure B.8: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

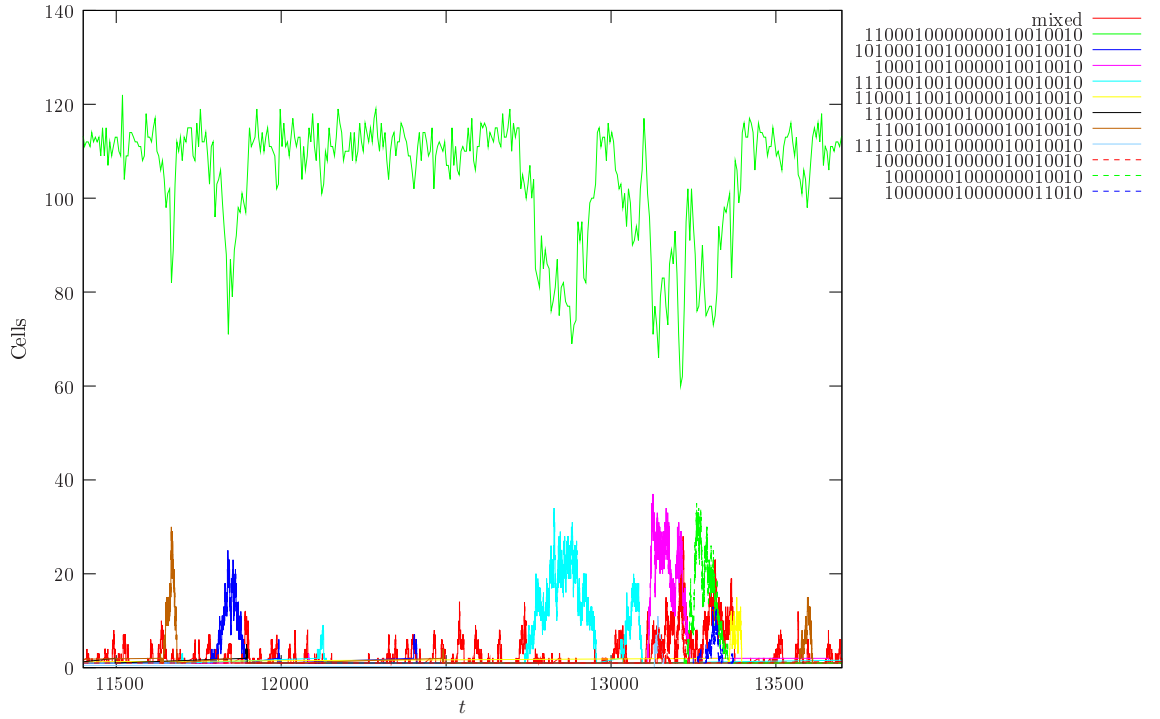


Figure B.9: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

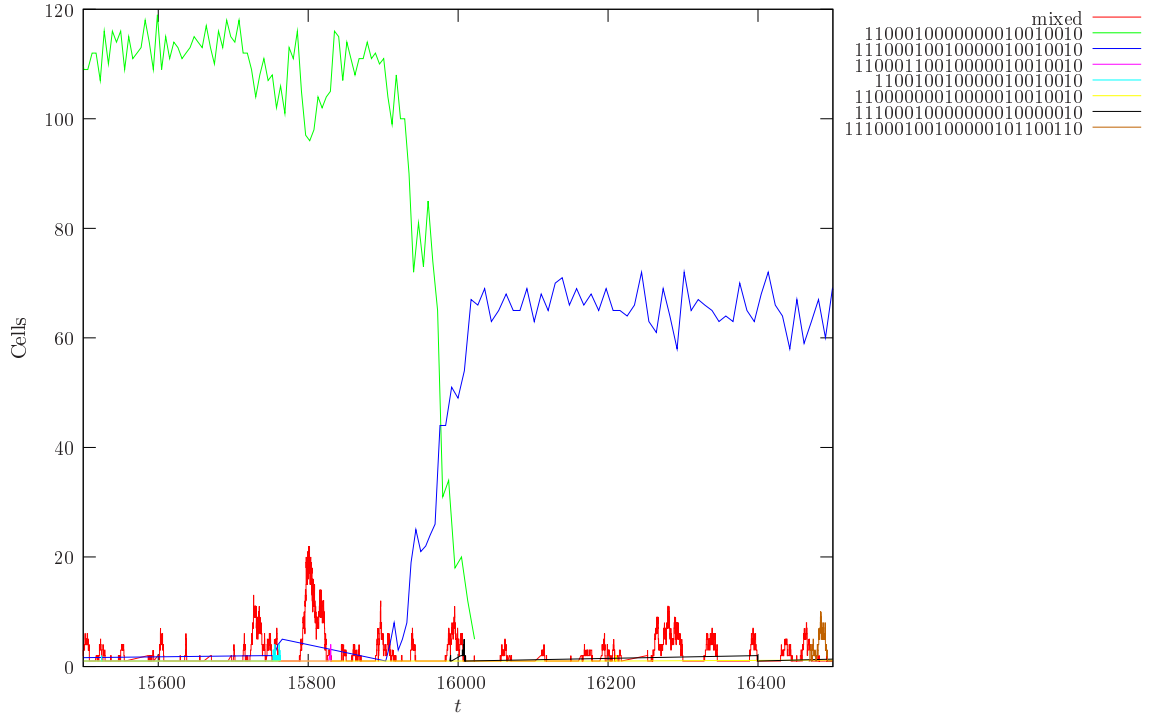


Figure B.10: Successful evolution of protocell-embedded 2-bit molecular counting via 3 transitions. Each plot represents a protocell lineage with a different principal molecular species.

- At $t \simeq 1750$, a protocell lineage with $z_i = 2$ successfully invades the population. Though there is no change in fitness with respect to the imposed molecular activity factor, the invading protocell species has a principal molecule which is a class-6 facultative parasite of the previous lineage, and has a molecular bit-string which is one bit *longer* than the previous lineage. This event again highlights that the modifications made in Chapter 6 make it possible to have a population displacement based on something other than molecular length, namely the spectrum of available mutants.
- At $t \simeq 2200$, a protocell lineage emerges (black plot in Figure B.7) which is a class-6 facultative parasite of the predominant cell lineage. This lineage gave rise to a new lineage whose principal molecular species was a single mutation away, though this new protocell species does not appear in the plot as it had already exceeded the number of parasitic molecules to cross the “mixed” cell threshold. It was this transient protocell species which gave rise to the emergent protocell lineage at $t \simeq 2500$ which proceeded to invade the protocell population. The dominance of this new lineage is accompanied by a low number of “mixed” protocells and consequently, low protocell diversity.

The events in Figure B.8 are:

- At $t \simeq 7600$, a protocell lineage emerges with the same $z_i = 2$ as the previously dominant lineage. The principal molecular species of the new lineage shares the class-6 facultative parasite relationship with the principal molecular species of the previously dominant lineage.

The events in Figure B.9 are:

- At $t \simeq 11600$, a protocell lineage emerges with the same $z_i = 2$ as the currently dominant lineage. The principal molecular species of the new lineage is one bit shorter than the currently dominant lineage, and this new lineage enjoys a brief period of protocell population growth.
- At $t \simeq 11800$, a protocell lineage emerges with the fitter $z_i = 3$ according to the imposed molecular activity factor. The principal molecular species of the new lineage is one bit longer than the currently dominant lineage and this new lineage does not successfully invade the protocell population.
- At $t \simeq 12750$, a protocell lineage emerges with the fitter $z_i = 3$ according to the imposed molecular activity factor. As before, the principal molecular species of the new lineage is one bit longer than the currently dominant lineage which once again apparently prevents this new species from successfully invading the protocell population.
- At $t \simeq 13200$ and $t \simeq 13350$, two protocell lineages appear which have the same $z_i = 2$ as the currently dominant lineage, and are one and two bits shorter than the currently dominant lineage respectively.

The events in Figure B.10 are:

- At $t \simeq 16000$, a protocell lineage emerges with the fittest $z_i = 3$ according to the imposed molecular activity factor. The principal molecular species of the new lineage shares the class-6 facultative parasite relationship with the principal species of the previously dominant lineage. This fitter lineage continues to dominate the protocell population for the remainder of the run.

B.5.1 Discussion of Results

The two experiments described above demonstrate the directed evolution of a protocell level trait which is emergent from the composition of the molecular network that is contained within that protocell. The essential evolutionary features of these two runs are typical of MCS-2 experiments carried out across a range of system parameters. The experiments presented here demonstrate evolution to optimal fitness via two or three protocell level fitness transitions. During system testing, it

was also observed that optimal fitness could be reached via a single population level fitness transition. Detailed analysis of these runs at the level of individual protocell strains revealed that the complicated internal molecular dynamics of protocells, as described in Sections 5.1.3 and 6.3.2.1, strongly influence the evolutionary trajectories that are available. Careful attention has been given to the presentation of the results of the above experiments, especially the detail in Section B.5, to show that the results of these experiments are explained in the context of the various concepts of protocell fitness that have been identified throughout this work. It does not automatically follow then that the outcome of an evolutionary experiment in such a system will be wholly, or even primarily, determined by the externally applied selection pressure, and is more often than not due to a combination of the kinds of complicated molecular dynamics presented in Chapters 5 and 6.

B.6 Conclusion

This chapter presented some significant modifications to the MCS to enable the directed evolution of protocell embedded molecular computation. This (minimal) computation was identified as the binary counting behaviour of molecular species, specifically the directed modification of a portion of the structure of newly replicated molecules. Aggregates of molecules were able to modulate the production rate of particular molecular sequences which in turn had an effect on the growth rates and average size of the protocells which contained them, through a new linkage between molecular populations and protocell fission. A detailed analysis was presented for two runs of the system, both of which evolved the specified optimum computational behaviour via qualitatively different pathways, with respect to an imposed molecular activity factor which essentially acted as a fitness function. The nature of this molecular activity factor was completely arbitrary, therefore the possibility to exchange it for a new fitness function existed, but was not explored.

The work presented earlier in Chapters 4, 5 and 6 provided a solid understanding of the molecular and protocellular level dynamics in the absence of any external perturbations. This understanding provided the foundation for the application of an externally provided arbitrary fitness function, and provided a toolkit for analysing the outcome of experimental runs of the system. The experimental verification of the convergent evolution of isolated sub-populations of molecules toward an arbitrary externally determined “goal state”, through the use of an external fitness function applied only at the protocell level, can be considered a supplementary contribution of this thesis. However, the presentation of these results would be incomplete without a discussion about the modelling methodology adopted. Section 4.1 justified

the general approach in the context of the literature, but as with all such modelling of real world phenomena, ideally, one would like to close the loop and validate the results achieved against a real-world system. It is an unfortunate fact however that a significant technical difficulty remains, in that there are no real-world examples of these kinds of artificial protocells against which we might validate the results. Throughout the thesis, various complex behaviours of the MCS have been presented, and the results obtained are not clear cut. It is not obvious then whether the observations made are artifacts of the system design, or whether they represent a realistic abstraction of the corresponding real-world phenomena. In conclusion, significant further work would be required both in the MCS itself, and more importantly in the realisation of wet-lab artificial protocells in order to make classical validation of the obtained results feasible.

Bibliography

- Adami, C. and Brown, C. T. (1994). Evolutionary learning in the 2d artificial life system avida. *Artificial Life*, IV:377–381.
- Adler, J. and Tso, W. (1974). Decision-making in bacteria: chemotactic response of escherichia coli to conflicting stimuli. *Science*, 184(4143):1292–2294.
- Altmeyer, S., Füchslin, R., and McCaskill, J. (2004). Folding stabilizes the evolution of catalysts. *Artificial Life*, 10(1):23–38.
- Barricelli, N. (1957). Symbiogenetic evolution processes realized by artificial methods. *Methodos*, 9(35-36):143–182.
- Barricelli, N. (1963). Numerical testing of evolution theories. *Acta Biotheoretica*, 16(3):99–126.
- Bartel, D. and Unrau, P. (1999). Constructing an RNA world. *Trends in Biochemical Sciences*, 24:M9–M13.
- Bedau, M. A. (1996). The nature of life. *The Philosophy of Artificial Life*, pages 332–357.
- Bedau, M. A. (1998). Four puzzles about life. *Artificial Life*, 4(2):125–140.
- Bedau, M. A. (1999). Can unrealistic computer models illuminate theoretical biology. *Proceedings of the 1999 Genetic and Evolutionary Computation Conference Workshop Program*, pages 20–23.
- Bedau, M. A., McCaskill, J. S., Packard, N. H., Rasmussen, S., Adami, C., Green, D. G., Ikegami, T., Kaneko, K., and Ray, T. S. (2001). Open problems in artificial life. *Artificial Life*, 6(4):363–376.
- Bentley, P. (2001). *Digital biology*. Headline, London.
- Bentley, P. J. (2004). Fractal proteins. *Genetic Programming and Evolvable Machines*, 5(1):71–101.

- Bentley, P. J. (2007). Systemic computation: A model of interacting systems with natural characteristics. *Parallel Algorithms and Applications*, 22(2):103–121.
- Bernal, J. (1949). The physical basis of life. *Proceedings of the Physical Society. Section B*, 62:597.
- Bernal, J. (1959). The problem of stages in biopoesis. *Proceedings of the first international symposium on the origin of life on the earth: held at Moscow 19-24 August 1957*, page 38.
- Boerlijst, M. and Hogeweg, P. (1991). Spiral wave structure in pre-biotic evolution: Hypercycles stable against parasites. *Physica D: Nonlinear Phenomena*, 48(1):17–28.
- Brenner, S., Jacob, F., and Meselson, M. (1961). An unstable intermediate carrying information from genes to ribosomes for protein synthesis. *Nature*, 190(4776):576–581.
- Brimacombe, R., Trupin, J., Nirenberg, M., Leder, P., Bernfield, M., and Jaouni, T. (1965). RNA codewords and protein synthesis, viii. nucleotide sequences of synonym codons for arginine, valine, cysteine, and alanine. *Proceedings of the National Academy of Sciences of the United States of America*, 54(3):954–960.
- Buss, L. (1987). *The Evolution Of Individuality*. Princeton University Press.
- Cairns-Smith, A. (1982). *Genetic takeover and the mineral origins of life*. Cambridge University Press.
- Church, A. (1936). An unsolvable problem of elementary number theory. *American Journal of Mathematics*, 58(2):345–363.
- Darwin, C. (1845). *Journal of researches into the natural history and geology of the countries visited during the voyage of H.M.S. Beagle round the world*. John Murray.
- Darwin, C. (1859). *Origin of Species by Means of Natural Selection, Or the Preservation of Favored Races in the Struggle for Life*. John Murray.
- Dawkins, R. (1991). *The blind watchmaker*. Penguin, Harmondsworth.
- Dawkins, R. (2009). *The Greatest Show on Earth: The Evidence for Evolution*. Bantam Press.

- Decraene, J., Mitchell, G., and McMullin, B. (2006). Evolving artificial cell signaling networks using molecular classifier systems. In *1st IEEE/ACM International Conference on Bio-Inspired Models of Network, Information and Computing Systems (IEEE/ACM BIONETICS 2006)*, Cavalese, Italy.
- Dewdney, A. K. (1984). Core wars. *Scientific American*, 250:14–22.
- Dittrich, P. and Speroni, P. (2007). Chemical organisation theory. *Bulletin of Mathematical Biology*, 69(4):1199–1231.
- Dittrich, P., Ziegler, J., and Banzhaf, W. (2001). Artificial chemistries - a review. *Artificial Life*, 7(3):225–275.
- Eigen, M. (1971). Selforganization of matter and the evolution of biological macromolecules. *Die Naturwissenschaften*, 58(10):465–523.
- Eigen, M., McCaskill, J., and Schuster, P. (1989). The molecular quasi-species. *Advances in Chemical Physics*, 75:149–263.
- Eigen, M. and Schuster, P. (1977). The hypercycle, a principle of natural self-organization. *Die Naturwissenschaften*, 64(11):541–565.
- Farmer, J. and Belin, A. (1990). Artificial life: The coming evolution. In *Artificial Life II: Proceedings of The Workshop On Artificial Life Held February, 1990 In Santa Fe, New Mexico*.
- Fontana, W. (1992). Algorithmic chemistry. In Langton, C. G., Taylor, C., Doyne-Farmer, J., and Rasmussen, S., editors, *Artificial Life II*, pages 159–209. Addison-Wesley, Redwood City, CA.
- Fontana, W. and Buss, L. (1994a). The arrival of the fittest: Toward a theory of biological organization. *Bulletin of Mathematical Biology*, 56(1):1–64.
- Fontana, W. and Buss, L. (1994b). What would be conserved if "the tape were played twice"? *Proceedings of the National Academy of Sciences*, 91(2):757–761.
- Forterre, P. and Philippe, H. (1999). The last universal common ancestor (LUCA), simple or complex? *Biological Bulletin*, 196(3):373–377.
- Fraser, C., Gocayne, J., White, O., Adams, M., Clayton, R., Fleischmann, R., Bult, C., Kerlavage, A., Sutton, G., Kelley, J., et al. (1995). The minimal gene complement of *Mycoplasma genitalium*. *Science*, 270(5235):397.

- Fuqua, W., Winans, S., Greenberg, E., et al. (1994). Quorum sensing in bacteria: the LuxR-LuxL family of cell density-responsive transcriptional regulators. *Journal of bacteriology*, 176(2):269–275.
- Gánti, T. (2003). *The Principles of Life*. New York : Oxford University Press.
- Gardner, M. (1970). Mathematical games: The fantastic combinations of john conway’s new solitaire game ‘life’. *Scientific American*, 223(4):120–123.
- Gibson, D., Benders, G., Andrews-Pfannkoch, C., Denisova, E., Baden-Tillson, H., Zaveri, J., Stockwell, T., Brownley, A., Thomas, D., Algire, M., et al. (2008). Complete chemical synthesis, assembly, and cloning of a *Mycoplasma genitalium* genome. *Science’s STKE*, 319(5867):1215.
- Gilbert, W. (1986). Origin of life: The RNA world. *Nature*, 319(6055):618.
- Glass, J., Assad-Garcia, N., Alperovich, N., Yooseph, S., Lewis, M., Maruf, M., Hutchison, C., Smith, H., and Venter, J. (2006). Essential genes of a minimal bacterium. *Proceedings of the National Academy of Sciences*, 103(2):425–430.
- Hanczyc, M. and Szostak, J. (2004). Replicating vesicles as models of primitive cell growth and division. *Current Opinion in Chemical Biology*, 8:660–664.
- Hogeweg, P. (1994). Multilevel evolution: replicators and the evolution of diversity. *Physica D: Nonlinear Phenomena*, 75(1-3):275–291.
- Hogeweg, P. and Takeuchi, N. (2003). Multilevel selection in models of prebiotic evolution: Compartments and spatial self-organization. *Origins of Life and Evolution of the Biosphere*, 33:375–403.
- Holland, J. (1998). *Emergence: From chaos to order*. Oxford University Press.
- Holland, J. H. (1975). *Adaptation in natural and artificial systems*. Ann Arbor MI: University of Michigan Press.
- Holland, J. H. (1976). Studies of the spontaneous emergence of self-replicating systems using cellular automata and formal grammars. In Lindenmayer, A. and Rozenberg, G., editors, *Automata, Languages, Development*, pages 385–404. North-Holland.
- Holland, J. H. (2006). Studying complex adaptive systems. *Journal of Systems Science and Complexity*, 19(1):1–8.
- Holland, J. H. and Reitman, J. (1977). Cognitive systems based on adaptive algorithms. *ACM SIGART Bulletin*, pages 49–49.

- Johnston, W. K., Unrau, P. J., Lawrence, M. S., Glasner, M. E., and Bartel, D. P. (2001). RNA-catalyzed RNA polymerization: Accurate and general RNA-templated primer extension. *Science*, 292(5520):1319–1325.
- Joyce, G. F. (1991). The rise and fall of the RNA world. *The New Biologist*, 3(4):399–407.
- Jukes, T. and Osawa, S. (1993). Evolutionary changes in the genetic code. *Comparative biochemistry and physiology, B*, 106(3):489.
- Kelly, C., Decraene, J., Mitchell, G. G., McMullin, B., and O’Brien, D. (2006a). Cellular computation using classifier systems. In *2006 International Workshop on Systems Biology*, Maynooth, Ireland.
- Kelly, C., McMullin, B., and O’Brien, D. (2006b). On protocell ‘computation’. In *European Conference on Complex Systems (ECCS-06)*.
- Kelly, C., McMullin, B., and O’Brien, D. (2008). Enrichment of interaction rules in a string-based artificial chemistry. In Bullock, S., Noble, J., Watson, R. A., and Bedau, M. A., editors, *Proceedings of the Eleventh International Conference on Artificial Life*.
- Kelly, C., McMullin, B., O’Brien, D., Mitchell, G. G., and Speroni di Fenizio, P. (2007). The Evolution of Complexity in a Multi-Level Artificial Chemistry. *European Conference on Complex Systems (ECCS-07)*, Dresden, Germany.
- Kruger, K., Grabowski, P. J., Zaug, A. J., Sands, J., Gottschling, D. E., and Cech, T. R. (1982). Self-Splicing RNA: Autoexcision and autocyclization of the ribosomal RNA intervening sequence of tetrahymena. *Cell*, 31(1):147–157.
- Langton, C. G., editor (1989a). *Artificial Life*. Addison-Wesley Publishing Company, Inc., Redwood City, California. Proceedings of an interdisciplinary workshop, Los Alamos, New Mexico, September, 1987.
- Langton, C. G. (1989b). Artificial life. In Langton (1989a), pages 1–47. Proceedings of an interdisciplinary workshop, Los Alamos, New Mexico, September, 1987.
- Lepot, K., Benzerara, K., Brown, G., and Philippot, P. (2008). Microbially influenced formation of 2,724-million-yr-old stromatolites. *Nature Geoscience*, 1(2):118.
- Lincoln, T. and Joyce, G. (2009). Self-sustained replication of an RNA enzyme. *Science*, 323(5918):1229.

- Luisi, P. (2002). Toward the engineering of minimal living cells. *The Anatomical Record*, 268(3):208–214.
- Martin, W. (2003). On the origins of cells: a hypothesis for the evolutionary transitions from abiotic geochemistry to chemoautotrophic prokaryotes, and from prokaryotes to nucleated cells. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 358(1429):59–85.
- Maynard Smith, J. (1969). The status of neo-darwinism. *Towards a theoretical biology*, 2:82–89.
- Maynard Smith, J. (1989). *Evolutionary Genetics*. Oxford University Press.
- Maynard Smith, J. and Szathmáry, E. (1997). *The Major Transitions in Evolution*. Oxford Press.
- Maynard Smith, J. and Szathmáry, E. (2000). *The origins of life: From the birth of life to the origin of language*. Oxford University Press, USA.
- McCaskill, J. et al., editors (2008). *Final Report on the Progress of the EU FP6 funded project: PACE (FP6-IST-FET-002035)*. PACE Consortium.
- McIlroy, D., Morris, R., and Vyssotsky, V. (1972). Darwin: A game of survival and (hopefully) evolution. *Software—Practice and Experience*, 2:91–96.
- McMullin, B. (1992). The holland α -universes revisited. *Proceedings of the First European Conference on Artificial Life*, MIT Press, Cambridge, pages 317–326.
- McMullin, B. (1997). Computational autopoiesis: The original algorithm. Technical report, Working Paper 97-01-001, Santa Fe Institute, Santa Fe, NM 87501, USA, January 1997.
- McMullin, B. (2000). John von neumann and the evolutionary growth of complexity: Looking backward, looking forward. *Artificial Life*, 6(4):347–361.
- McMullin, B., Kelly, C., and O’Brien, D. (2007a). Multi-level selectional stalemate in a simple artificial chemistry. *Advances in Artificial Life*, pages 22–31.
- McMullin, B., Kelly, C., and O’Brien, D. (2007b). Preliminary steps toward artificial protocell computation. *International Conference on Morphological Computation, Venice, Italy*.
- McMullin, B., Kelly, C., and O’Brien, D. (2008). Evolution of protocell-embedded molecular computation. In McCaskill et al. (2008).

- McMullin, B. and Varela, F. (1997). Rediscovering computational autopoiesis. *Fourth European Conference on Artificial Life*, pages 38–47.
- Mendel, G. (1866). *Versuche über Pflanzenhybriden [Experiments on plant hybrids]*. Akademische Verl.-Gesellschaft, Berlin.
- Miller, S. (1955). Production of some organic compounds under possible primitive earth conditions. *Journal of the American Chemical Society*, 77(9):2351–2361.
- Nachman, M. and Crowell, S. (2000). Estimate of the mutation rate per nucleotide in humans. *Genetics*, 156(1):297.
- Nielsen, P. (1993). Peptide nucleic acid (PNA): a model structure for the primordial genetic material? *Origins of Life and Evolution of Biospheres*, 23(5):323–327.
- Nirenberg, M., Matthae, J., and Jones, O. (1962). An intermediate in the biosynthesis of polyphenylalanine directed by synthetic template RNA. *Proceedings of the National Academy of Sciences of the United States of America*, 48(1):104.
- Olson, J. (2006). Photosynthesis in the archaean era. *Photosynthesis Research*, 88(2):109–117.
- Orgel, L. (2004). Prebiotic chemistry and the origin of the RNA world. *Critical Reviews in Biochemistry and Molecular Biology*, 39(2):99–123.
- Osawa, S., Jukes, T., Watanabe, K., and Muto, A. (1992). Recent evidence for evolution of the genetic code. *Microbiology and Molecular Biology Reviews*, 56(1):229.
- Overton, C. (1895). "Über die osmotischen eigenschaften der lebenden pflanzen-und tierzelle—about the osmotic properties of living plant and animal cells. *Vierteljahresschr. Naturforsch. Ges. Zürich*, 40:159–201.
- Packard, N., Crutchfield, J., Farmer, J., and Shaw, R. (1980). Geometry from a time series. *Physical Review Letters*, 45(9):712–716.
- Palsson, B. (2006). *Systems biology: properties of reconstructed networks*. Cambridge University Press New York, NY, USA.
- Pargellis, A. (2001). Digital life behavior in the amoeba world. *Artificial Life*, 7(1):63–75.

- Rasmussen, S., Bedau, M., Chen, L., Deamer, D., Krakauer, D., Packard, N., and Stadler, P. (2008). *Protocells: bridging nonliving and living matter*. MIT Press, USA.
- Rasmussen, S., Chen, L., Deamer, D., Krakauer, D. C., Packard, N. H., Stadler, P. F., and Bedau, M. A. (2004a). Transitions from nonliving to living matter. *Science*, 303(5660):963–965.
- Rasmussen, S., Chen, L., Stadler, B., and Stadler, P. (2004b). Proto-organism kinetics: Evolutionary dynamics of lipid aggregates with genes and metabolism. *Origins of Life and Evolution of Biospheres*, 34(1):171–180.
- Rasmussen, S., Knudsen, C., Feldberg, R., and Hindsholm, M. (1990). The core-world: Emergence and evolution of cooperative structures in a computational chemistry. *Physica D: Nonlinear Phenomena*, 42(1-3):111–134.
- Ray, T. S. (1991). An approach to the synthesis of life. In Langton, C. G., Tayler, C., Farmer, J. D., and Rasmussen, S., editors, *Artificial Life II*, pages 371–408, Reading, MA. Addison-Wesley.
- Reed, J., Toombs, R., and Barricelli, N. (1967). Simulation of biological evolution and machine learning. i. selection of self-reproducing numeric patterns by data processing machines, effects of hereditary control, mutation type and crossing. *Journal of Theoretical Biology*, 17(3):319–42.
- Reich, K. (2000). The search for essential genes. *Research in Microbiology*, 151(5):319–324.
- Rendell, P. (2002). Turing universality of the game of life. *Collision-based computing*, pages 513–539.
- Segré, D., Ben-Eli, D., Deamer, D. W., and Lancet, D. (2001). The lipid world. *Origins of Life and Evolution of Biospheres*, 31(1):119–145.
- Shannon, C. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423.
- Smith, J. and Martin, L. (1973). Do cells cycle? *Proceedings of the National Academy of Sciences*, 70(4):1263.
- Stadler, P., Fontana, W., and Miller, J. (1993). Random catalytic reaction networks. *Physica D*, 63(3-4):378–392.

- Stegmann, U. E. (2004). The arbitrariness of the genetic code. *Biology and Philosophy*, 19(2):205–222.
- Szathmáry, E. (1986). The eukaryotic cell as an information integrator. *Endocytological Cell Research*, 3:113–132.
- Szathmáry, E. and Maynard Smith, J. (1997). From replicators to reproducers: the first major transitions leading to life. *Journal of Theoretical Biology*, 187:555–571.
- Szostak, J., Bartel, D., and Luisi, P. (2001). Synthesizing life. *Nature*, 409:387–390.
- Szybalski, W. and Skalka, A. (1978). Nobel prizes and restriction enzymes. *Gene*, 4(3):181.
- Taylor, T. (1999). *From Artificial Evolution to Artificial Life*. PhD thesis, School of Informatics, University of Edinburgh.
- Theis, M., Gazzola, G., Forlin, M., Poli, I., Hanczyc, M., and Bedau, M. (2006). Optimal formulation of complex chemical systems with a genetic algorithm. *ECCS06 online Proceedings (p. 193)*, Oxford, UK.
- Turing, A. (1937). On computable numbers. *Proceedings of the London Mathematical Society*, 2(42):230–65.
- Varela, F., Maturana, H., and Uribe, R. (1974). Autopoiesis: The organization of living systems, its characterization and a model. *BioSystems*, 5(4):187–196.
- Von Neumann, J. and Burks, A. W. (1966). *Theory of Self-Reproducing Automata*. University of Illinois Press Champaign, IL, USA.
- Voytek, S. and Joyce, G. (2007). Emergence of a fast-reacting ribozyme that is capable of undergoing continuous evolution. *Proceedings of the National Academy of Sciences*, 104(39):15288.
- Wächtershäuser, G. (1990). Evolution of the first metabolic cycles. *Proceedings of the National Academy of Sciences*, 87(1):200–204.
- Wächtershäuser, G. (2007). On the chemistry and evolution of the pioneer organism. *Chemistry & Biodiversity*, 4(4).
- Wang, J. (1991). DNA topoisomerases: why so many. *J. Biol. Chem*, 266(11):6659–6662.
- Watson, J. and Berry, A. (2003). *DNA: The secret of life*. Alfred a Knopf Inc.

Watson, J. and Crick, F. (1953). A structure for deoxyribonucleic acid. *Nature*, 171:737–738.

Wei, S. and Young, R. (1989). Development of symbiotic bacterial bioluminescence in a nearshore cephalopod, *euprymna scolopes*. *Marine Biology*, 103(4):541–546.