**Cell states, fates and reprogramming
insights from neural networks, graphical and computational approaches**

Hannam, Ryan Michael

*Awarding institution:*
King's College London

# King's College London

# Cell states, fates and reprogramming: insights from neural networks, graphical and computational approaches

A THESIS PRESENTED
BY
RYAN HANNAM
TO
KING'S COLLEGE LONDON

FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

SUPERVISED BY
DR ALESSIA ANNIBALE
PROFESSOR REIMER KÜHN

# Abstract

Lineage specification was long thought to be an irreversible developmental process. However, with the advent of cell reprogramming and the discovery of induced pluripotent stem cells (iP-SCs), it was shown that differentiation is in fact reversible. Cell reprogramming has mainly been studied experimentally, with no universally accepted theory explaining the phenomena. The purpose of this thesis is to drive forward our understanding of cell biology, by introducing analytical models for the interaction between genes and studying the transitions between the emergent cell types. This is done by appealing to key concepts from biology and employing tools commonly used in the field of statistical physics. Inspired by models of neural networks, a model for cell reprogramming is introduced in which cell types are hierarchically related dynamical attractors corresponding to cell cycles. Stages of the cell cycle are fully characterised by the configuration of gene expression levels, and reprogramming corresponds to triggering transitions between such configurations. Two mechanisms were found for reprogramming in a two-level potency hierarchy: cycle specific perturbations and a noise-induced switching. The former corresponds to a directed perturbation that induces a transition into a cycle-state of a different cell type in the potency hierarchy (mainly a stem cell) whilst the latter is a priori undirected and could be induced, e.g. by a (stochastic) change in the cellular environment. The reprogramming model is governed by the interaction between gene expression levels, as originally hypothesised by Waddington in his Epigenetic Landscape analogy. To further develop the biological significance, a detailed mechanism for these interactions between genes, in the form of regulation through transcription factors, is studied. This consists of constructing a bipartite graph framework for gene regulatory networks. A technique that integrates the genome and transcriptome into a single regulatory network. With this perspective, we are able to deduce important features of the regulatory network that exists in every cell type, such as the typical interactions required to sustain a net gene expression profile and how regulatory interactions must change to support multicellular life.

# ACKNOWLEDGEMENTS

First and foremost, I wish to express my deepest gratitude to Dr Alessia Annibale and Professor Reimer Kühn for their constant guidance and patience throughout my PhD. Your passion for science came across in every encounter, providing both stimulating discussion and enlightening lessons. For their advice regarding certain biological aspects of this work, I would also like to recognise Dr Attila Csikász-Nagy and Dr Jens Kleinjung. I would also like to thank the disordered systems group, and wider Mathematics and Physics departments, for breaking the stereotype that academic research is an isolating experience. Similarly, my gratitude goes towards the Cross-disciplinary Approaches to Non-Equilibrium Systems (CANES) Centre for Doctoral Training (CDT), for funding this research and providing a fantastic support network, with particular recognition going to Professor Peter Sollich, Dr Chris Lorenz and Dr Joe Bhaseen.

Next, I would like to acknowledge all the other graduate students with whom I have made this journey. In many of you, I have made everlasting friendships. Carla and Riccardo, you have provided constant emotional and intellectual support. We have shared countless memories over coffee, food and music. For this, I wish you both great health and success in all your future endeavours. Davide, Edgar, Pablo and Robin, our joint enthusiasm for good science and even better beer guided our discussions to taprooms and breweries across the capital. To you, I say, "Salute!", "¡Salud!", "Saúde!" and "Santé!".

I am also indebted to my parents who have encouraged every decision I have made to reach

this point. It is truly a privilege to have a family that has supported my academic studies. Without them, I couldn't have invested the time to become the person I am today.

Above all, there is one individual whom without none of this would have been possible. You have stood by me through countless difficult times and have seen me at my highest and lowest. Your continued love and support made every step of this journey a little bit easier. For those reasons and many more, Rachael, I dedicate this thesis to you.

# Articles contributing to this thesis

1. R. Hannam, A. Annibale, R. Kühn, *Cell reprogramming modelled as transitions in a hierarchy of cell cycles*, J. Phys. A: Math & Theor. 50 425601 (2017)

2. R. Hannam, R. Kühn, A. Annibale, *Percolation in bipartite Boolean networks and its role in sustaining life* (in preparation)

Papers 1 & 2 cover the majority of the contents of chapter 2 and 4 respectively.

# CONTENTS

# FIGURES

# 1

## INTRODUCTION

The 2012 Nobel Prize for physiology or medicine was awarded to John Gurdon and Shinya Ya-
manaka, "for the discovery that mature cells can be reprogrammed to become pluripotent" [1].
That is, cells that have reached a terminal fate (cell-type) after development were converted to
cells that have the ability to differentiate into many other cell types. Gurdon pioneered the method
of somatic cell nuclear transfer in the 1960s, whilst Yamanaka and colleagues recently introduced
the method of cell reprogramming. The latter technique relies on the introduction of genes, that
encode transcription factors (TFs), which are highly expressed in the target cell type, viz. em-
bryonic stem (ES) cells, using a retrovirus. These cellular reprogramming experiments convert

somatic cells into induced pluripotent stem cells (iPSCs), which strongly resemble ES cells in both morphology and gene expression profiles.

This discovery contradicted the accepted scientific view of the time. Development was thought to be an irreversible process, with lineage specification directing the arrow of time and potency levels dropping after successive differentiation events. However, an increasing body of literature is showing evidence that a given cell can be converted to any other cell type, either directly by trans-differentiation or via guided differentiation from an iPSC. Like other stem cells, iPSCs have the ability to both self renew and differentiate into multiple different cell types (i.e. they are pluripotent). Thus, cell reprogramming opens up many possibilities with applications in personalised and regenerative medicine, disease modelling and drug development [2–6]. An increased understanding of cell reprogramming could also improve our understanding of developmental biology, in which cell fate decisions play an important role [7–10].

Although it has been more than a decade since the discovery of cell reprogramming, the underlying processes remain relatively elusive. This PhD project aims to bridge the gap between the experimental results and our understanding of molecular biology. Much of the current theoretical work treats cell types as static entities, however, cells are dynamic bodies with every cell undergoing its own cell cycle. Thus, any biologically realistic model of cell reprogramming should also produce cell cycle dynamics. This is the main goal of my PhD - to create a model that can describe how cell types emerge from the underlying molecular biology and explain transitions between distinct cell types, that is built upon key biological facts, using quantitative methods from fields such as physics, mathematics and informatics.

The remainder of this thesis is organised as follows: The rest of this chapter is dedicated to a brief review of cell reprogramming and the motivation for a statistical physics modelling approach. Chapter 2 will introduce a model for cell reprogramming as transitions between dynamical attractors that represent the gene expression profiles across cell cycles. This model is then tested on real data, with parameters for the model inferred and its assumptions tested, in chapter

(a) Waddington Landscape    (b) Underpinning gene interactions

**Figure 1.1.1:** Waddington's Epigenetic Landscape metaphor for development. A complex landscape is shaped by the interactions in a gene regulatory network. A cell's lineage is tracked by following a ball rolling down the landscape from a state of pluripotency. Successive valleys in the landscape represent different cell types with increasing levels of specificity (decreasing levels of potency). Each cell fate decision is a binary choice represented by the fork junctions of intersecting valleys. Each panel for this figure is taken separately from [11].

3. Chapter 4 delves deeper into the interaction between gene expression levels in models for cellular identity. Finally, chapter 5 summarises the main results of the thesis and outlines directions where future work may be most beneficial.

## 1.1.    A brief overview of cellular reprogramming

In the 1950s, Waddington introduced a metaphor of a ball rolling down an Epigenetic landscape to describe cell fate decisions during development [11]. Initially, the ball starts high up the landscape in a state of pluripotency. As it traverses down the landscape, it is funnelled into different valleys (see figure 1.1.1) which represent different cell types. Each of these cell fate decisions is a binary choice, with the ball rolling into either one of the valleys. The number of cell types and their specificity increases the further they are located down the landscape. Although Waddington's landscape was a metaphor, he postulated that the landscape could be shaped by the interaction between genes. Imagine the epigenetic landscape as a rubber sheet and the genes as weights. Taut strings between the weights and the landscape will govern its topology. These strings, and their tensions, then represent the interactions between the genes, such as the regulation of ex-

pression levels due to transcription factors. Thus, Waddington postulated that the interaction between sets of genes creates the valleys (cell types) and the barriers in the epigenetic landscape. With little mathematical training, and pre-dating the discovery of transcription factors, Waddington's landscape remained nothing more than an elegant metaphor.

The first formal theoretical proposal for cell types as emergent properties of a genetic network came from Kauffman in the 1960s. He attempted to model different cell types as dynamical behaviours, i.e. attractors, of random boolean networks [12]. Nodes in the network corresponded to binary gene expression levels and the edges between them correspond to their interactions. A gene was then said to be expressed, or not, depending on a randomly chosen logical function of the inputs from other genes. Distinct cell types emerge in the network as attractors of the dynamics. The number of cell types scaled with the size of the networks in a manner that was qualitatively in line with the observed number of cell types in different organisms, as a function of the number of their genes. Differentiation was modelled as a Markov process between different modes of behaviour in the network, with replication times accurately predicted as a function of the size of the gene network in a cell.

Kauffman's work was then extended by Wolpert in the 1970s. He considered cells as automata whose state is described by a set of binary gene expression levels [13]. He proposed a model for development in terms of changes in this state and asked how turning on/off individual genes affected the state of the system. Where his work truly differed from Kaufman's was in the inclusion of external cues in an attempt to model morphogenesis.

Whilst these theories were being produced, the successful retrieval of a pluripotent state was first demonstrated experimentally by John Gurdon. In his seminal work, he showed that it was possible to reprogram a cell using a technique known as somatic cell nuclear transfer (SCNT) [14]. This technique involves transferring the nucleus of a somatic cell into an enucleated embryo. The result is that the environment of the host cell alters the behaviour of the donor nucleus. The donor nucleus eventually expresses the same genes as the embryo's nucleus would have and thus

loses its specificity.

Dormant gene expression patterns were also shown to be reactivated in cell fusion experiments, in the 1980s [15, 16]. In these experiments, a somatic cell is merged with a pluripotent cell resulting in a hybrid that maintains some aspects of each of the original cell types. However, unlike somatic cell nuclear transfer, these experiments result in cells which have an increased ploidy level due to the presence of two nuclei. Thus, if the aim is to convert one cell type directly into another, cell fusion is less successful than somatic cell nuclear transfer. However, the hybrid cells can become diploid once more if one of the nuclei is removed after the fusion process.

Three decades later, cell reprogramming via retroviral transduction was reported by Yamanaka and Takahashi [17]. They hypothesised that embryonic stem cells and oocytes contain the necessary signals for promoting and sustaining pluripotency, because of the success of SCNT and cell fusion experiments. Genes that were considered to be important for pluripotency were identified as those which are highly expressed in embryonic stem cells. It was hypothesised that introducing large quantities of products of these genes to a cell would encourage them to recover the property of pluripotency. These experiments were successful, first with mice [17] and then human [18] fibroblast cells being reprogrammed to the iPSC state in 2-3 weeks. The iPSCs reprogrammed from mice cells were shown to contribute to chimera formation when inserted into developing mice embryos, further establishing their status as stem cells, along with their gene expression profiles and morphology. Remarkably, the authors were able to narrow down the initial concoction of 24 genes, thought to be important for pluripotency, to just 4 reprogramming ingredients: Oct3/4, Sox2, Klf4 and c-Myc. These factors are now commonly known as the OSKM or Yamanaka factors. Despite the success of these experiments, reprogramming was very inefficient. The reprogramming of cells to iPSCs in these original experiments took approximately 2-3 weeks, with only around 0.001% of cells successfully reaching the iPSC state.

The inefficiency of cell reprogramming raised a debate about whether only certain cells in a population had the capacity to be reprogrammed. Jacob Hanna quashed this idea by showing that

almost all cells have the potential to reprogrammed using the Yamanaka factors, given sufficient time and culture conditions [19]. This result suggests that the potential for pluripotency must be common to all cells. Because it is the same in every cell type, the genome of an organism is a strong candidate for the source of pluripotency, as first predicted by Waddington. Furthermore, SCNT, cell fusion and cell reprogramming are all consistent with the idea that the mechanisms converting cells to a pluripotent state occur in the nucleus, not the cytoplasm. It is now widely accepted that cell reprogramming arises from a change in gene expression patterns within a cell.

Since the discovery of iPSCs many protocols have been found for cell reprogramming using small molecules (such as mRNA) in place of the Yamanaka factors [20, 21]. Protocols are also consistently being developed to implement trans-differentiation. Previously, if a conversion between two terminally differentiated cell types was required one would first need to generate iPSCs and then guide their differentiation to the desired final cell type using specific culture conditions.

The "omics" revolution and big data have dramatically changed the field of biology. Informatics and data analysis are becoming increasingly valuable tools due to the large data sets that are collected from experiments. Historically, different cell types were classified in a qualitative manner based on morphology, by studying individual molecular components (such as proteins or RNA molecules) using biofluorescent markers, or a combination of the two. However, the big data methodologies have improved the classification of cell types on a molecular level [22, 23]. In recent years, technologies have improved to the stage at which individual cells can have their entire gene expression profile analysed using sequencing techniques [24–26]. Furthermore, by tracking the transcriptomics of a single cell, it is possible to follow molecular changes during lineage specification (analogous to mapping the ball's trajectory down the Waddington landscape) or identify components that could be crucial for cell fate decisions [27–32].

Enabled by these improved technologies, Sui Huang et al. were the first to demonstrate experimental evidence that cell types are high dimensional attractors in a gene expression space, as originally hypothesised by Waddington. By tracking the gene expression profile of differenti-

ating neutrophil cells under a variety of conditions, they showed that differentiation trajectories converge from many different directions to give the same expression profile in the final cell state [33]. This work paved the way for a new wave of dynamical systems approaches to modelling cell fates on complex gene interaction networks. Since this experimental result was published, the Huang group has continued to model cell fates using a complex systems approach. Typically, they consider cell fate decisions during development as a hierarchy of binary switches that represent developmental branching points [34–37].

Other than gene regulation, there is evidence to show that epigenetics plays an important roll in cell reprogramming. The changes in chromatin structure due to epigenetic modifiers allow different genes to be accessed by transcription factors [38, 39]. This structural change could facilitate the activation (or inhibition) of genes that are important for pluripotency (or specificity). Artyomov et al. developed a computational model in an attempt to capture the interplay between gene expression and epigenetic marks across the cell cycle, in order to predict pathways for reprogramming [40]. They model each cell cycle as a two-stage process. Gene expression levels and epigenetic marks are updated separately in a two-phase cell cycle. The genetic and epigenetic networks are able to interact in the sense that the epigenetic state provides feedback for the gene expression levels and vice versa. In their model, the authors assume that each cell type is defined by the expression of very few genes, referred to as a module, whilst the rest of the genome is inactive. However, this approach does not allow for any overlap in the gene expression profiles of different cell types. Furthermore, they use these distinct modules to explain the cell potency hierarchy. This implies that each module also has a potency level associated with it, i.e. the potency hierarchy of different cell types is defined in the DNA sequence of the genome and is not an emergent property of the interactions between gene expression levels. Despite some shortcomings in their model, this work is the only theoretical model for cell reprogramming that includes cell cycle dynamics that I am currently aware of.

The cell cycle may play a more important role than that suggested by Artyomov. There are

specific stages of the cell cycle in which many molecular components, including the genome, are duplicated. These stages use broadly the same molecular mechanisms across the different cell types of an organism, and thus represent some level of molecular uniformity. Furthermore, during asymmetric division two daughter cells are produced with different fates, so the checkpoints of the cell cycle ought to be capable of distinguishing between the components required for different cell types. Whether or not the cell cycle plays an important role in cell reprogramming remains an open question. However, technologies are now reaching a stage where it is possible to sequence cells during specific cycle phases [41, 42] opening the door to study correlations between cell cycle and reprogramming events. Unsurprisingly protein concentrations vary during the cell cycle, and hence it is reasonable to expect gene expression levels do so too. Therefore, regardless of any relation to reprogramming, any accurate model that defines a cell type as the combined effect of gene expression levels should be able to reproduce the variability that arises during the cell cycle.

There is no current accepted theory for the mechanisms behind cell reprogramming, that accurately predicts transitions between cell types in a quantitatively meaningful manner. Such a theory would clearly be hugely significant to fundamental molecular biology as well as development and biomedicine. This is partly due to the complexity of gene regulation and the role that noise plays in gene expression. Ordinary differential equation models of gene regulation can be cumbersome from an analytic point of view even when modelling a simple feedback loop [43]. Recent studies have also shown that the logic in gene regulatory networks is highly susceptible to even low levels of noise [44]. Thus, amplification of the transcription level noise may be a mechanism through which a cell can de-differentiate from a robust terminal state [45]. Noise is also thought to play key roles on a cell population level. Variations in gene expression levels in populations of cells (of a single cell type) may allow cells to rapidly respond to environmental changes [46]. This heterogeneity across many cells of a single type may also result from individual cells approaching differentiation at different rates [34], with variation in expression levels due

to preparation for differentiation or the events of the cell cycle.

Undeterred by the uncertainty surrounding cell reprogramming, iPSCs have already entered the medical world with the first clinical trials involving iPSCs occurring in the recent past [47]. Thus, it is vital that the theoretical models catch up with the fast-paced experimental developments.

There are various mathematical and computational models for gene regulation, cell fate decisions, and cell reprogramming, other than those touched on above. For an overview of these models the reader is directed to references [7, 34, 48–52]. Many of these works oversimplify the problem or take Waddington's metaphor too seriously, by trying to define a developmental landscape in terms of incomplete network entropies or pseudo-potentials [53–55]. Not only is the latter case too literal an interpretation of a metaphor, but both of these types of models typically result in no meaningful measurable quantities for experimental biologists to investigate. The most significant theoretical work on cell reprogramming comes in the form of Mogrify, a software tool to predict the most efficient transcription factors for direct reprogramming between human cell types [56]. Mogrify combines gene expression data with regulatory network information to predict the most effective transcription factors to induce a cell conversion. It has successfully identified known, and predicted new, reprogramming factors for trans-differentiation in human cells. Despite being a powerful predictive tool for reprogramming, Mogrify does not explain the mechanisms behind cell reprogramming transitions.

Parallel to the improved understanding of molecular biology that came with experimental technologies, there was the birth of complex systems theory which arguably originates as far back as the Ising model in statistical physics. Initially, many of the models in complex systems were developed independently in different fields, but have culminated in the recent development of predictive models that have had success in capturing the emergent properties of cells from the underlying interacting genes - as originally hypothesised by Waddington. In the next section, the origins of the most significant of these models are discussed to set the scene for many of the tech-

niques that are used in the remainder of this thesis. For a direct comparison between the models touched on in the next section, and their relation to models in the field of machine learning, the reader is pointed to the review by Fierst and Phillips [57].

## 1.2. FROM THE ISING MODEL TO GENE REGULATION VIA NEURAL NETWORKS

The Ising model is a statistical mechanics model for (ferro)magnetism, invented in the 1920s by Wilhelm Lenz [58]. It is named after his student, Ernst Ising, who solved it in one dimension [59] and can be found in many modern undergraduate physics courses. The model consists of binary variables, $\sigma_i \in \{\pm 1\}$, that represent the magnetic dipole moments of atomic spins in a material. The spins are arranged on a lattice or network, and each spin is able to interact with its neighbours. The state of the system at any time is fully characterised by the configuration of spins $\boldsymbol{\sigma}(t) = (\sigma_1, \ldots \sigma_N)$, where $N$ is the number spins in the system. Although, a very simplified model, it captures the behaviour of ferromagnetic materials that are able to maintain their magnetisation long after they are exposed to an external magnetic field (unlike paramagnetic materials that lose theirs when the external field is removed). The dynamics of the Ising model is governed by the interactions between spins $J_{ij}$, where $i$ and $j$ are used to denote the site labels of the spins. If $J_{ij} > 0$, the interaction is said to be ferromagnetic and spins prefer to have their magnetic moments aligned. If $J_{ij} < 0$, the interaction is antiferromagnetic and spins will prefer to have their magnetic moments oriented in opposing directions. If $J_{ij} = 0$ the spins $i$ and $j$ do not interact. Similarly, the sign local field $h_j$ of each site governs how each spin reacts to an external field.

In the 1940s, McCulloch and Pitts developed a model for logic processing in the brain [60]. They developed a model for artificial biological neurons. Each neuron has a set of inputs with weights determining the output of the neuron. The inputs can either be excitatory or inhibitory. If the inhibitory inputs are activated then the neuron is repressed, otherwise, its output is calculated by summing across all excitatory inputs and comparing this value with some threshold. The model can be viewed as a network of interacting neurons that are connected via their inputs and

outputs. Each neuron in the network can be in one of two states, either repressed "$0$" or activated "$1$". The network of neurons is able to reproduce gated logic systems that can perform Boolean operations. The threshold used in the interactions is justified as many biochemical reactions appear to behave similarly to a Hill function.

Inspired by works of bistability in biochemical feedback switches, Kauffman created a network model for gene regulation in which cell types are attractors of the dynamics [12]. Whilst Kauffman did not directly cite the Ising model in his original Boolean network model, there are clearly many similarities between the two. Kauffman's Boolean network model consists of a network of interacting nodes that represent genes. Each node can be described by a binary variable, "on" or "off", which describes whether a gene is activated or inhibited. The state of the system at any time is then described by the configuration of the $N$ binary genes. The key difference between the Ising model and Kauffman's work lies in the dynamics of the two models. The Ising model updates each magnetic spin based on the sum of the interactions between neighbours, whereas random Boolean networks assign each node, or gene, with a random Boolean function of its inputs. Hence, Boolean networks can reproduce logical rules (e.g. AND, OR, NAND, NOR, etc.). Although the interactions of the Ising model can be constructed from Boolean functions, random Boolean networks are much more general and can have high levels of non-linearity in their dynamics. The interactions in Boolean networks can also be asymmetric because the network is constructed by randomly assigning a gene's inputs from the rest of the genes in the network. Kauffman mainly studied the stability and behaviour of attractors in the dynamics, investigating the effects of noise on the attractor lengths. Differentiation in Boolean network models corresponds to transitions between two attractors either by a signal or noise. Kauffman found that the number of attractors, or cell types, scales with the amount of DNA content (number of nodes) with a power of $0.63$. Since Kauffman's original paper in 1969, there has been a large body of literature investigating Boolean networks and their relationship to biological networks, such as gene regulatory networks [61].

Spin glasses are disordered magnetic systems, named after the positional disorder of amorphous structures such as glasses [62]. In a spin glass, a distribution of interactions $J_{ij}$ results in irregular patterns in the configuration of atomic spins. Any opposing interactions compete with one another to align a spin in different orientations. This "frustration" results in multiple (meta)stable configurations compared to ferromagnetic systems. Spin glasses have been well studied in the statistical physics community with interactions over different length scales such as between neighbours, in the Edward-Anderson (EA) model [63], or over any length scales in the Sherrington Kirkpatrick (SK) model [64]. Spin glasses have a rich behaviour in which the system can capture a number of magnetic properties, such as paramagnetism, hysteresis and remnant magnetism, depending on the strength of the interactions between spins relative to the thermal noise in the system.

The robust nature of the metastable states in spin glasses led to the theory being adapted to understand neural networks (NNs) [65]. Binary spin states were reinterpreted as firing/quiescent neurons with synaptic interactions $J_{ij}$. A popular recurrent neural network model is the Hopfield model [66] of associative memory, that uses the Hebbian learning rule to store memories or configurations of neuronal activity in the synaptic interactions. These configurations are a form of quenched disorder because they do not change with respect to the level of noise in the system. The Hopfield model is an associative neural network because its memory retrieval capabilities are robust to input errors and hardware failure. The Hopfield model has a relatively low storage capacity: the number of memories/patterns $P$ that can be stored in a network of $N$ neurons scales as $P \sim \alpha N$, where $\alpha \approx 0.14$ [67, 68]. The model has multiple regimes of behaviour, acting as a spin glass, paramagnetic system or capable of retrieving memories, depending on the level of noise in the system and the storage level $\alpha$ [67, 68].

Another notable model for gene regulatory networks, that operates in the same manner as some neural network models, was introduced by Wagner in the 1990s [69, 70]. Similar to Kauffman's Boolean network the state of the system is described by the configuration of binary "on/off"

genes, $(G_1, \ldots, G_N)$, that are nodes in the network. Genes interact with one another in the model through the proteins that they express $(P_1, \ldots, P_N)$. Each gene can express a single protein or not, i.e. $P_i \in \{0, 1\}$. Those proteins then determine the expression level of the genes through activation/inhibition $S_i = 2P_i - 1$, where $S_i = -1$ and $S_i = 1$ correspond to repression and activation respectively. The state of each gene at a given time is updated by a step function of a weighted sum of the expression states. This is analogous to the update rule used in the dynamics of artificial neural networks, such as those developed by McCulloch and Pitts [60]. Wagner's gene network model can converge to a fixed point or cycle between different expression states, from an initial state of protein levels. Wagner focused his study on networks that converge to a stable equilibrium. The realisation of a stable equilibrium in a network was later termed "developmental stability" [71, 72], because Wagner gene networks are typically used to study the evolution of a population due to genetic mutations, genetic drift or environmental selection [73–79].

There are clearly many parallels that can be drawn between spin glasses, neural networks and gene regulatory networks (GRNs). Just as neural networks are robust to disruption of the synaptic efficacies and neuron failure, the morphological and phenotypic properties of the different cell types that arise from distinguishable gene expression profiles are largely robust to genetic variations and DNA replication errors. The individual components interact with one another through regulatory interactions of activation and inhibition, that encourage and suppress gene expression, analogous to magnetic interactions aligning spin orientations. Furthermore, the number of stable cell types is just a small fraction of all possible gene expression profiles, which is similar to the low storage capacity of neural networks. Next, external signals, such as morphogen gradients, can alter the expression of genes similar to spins aligning with local and external fields. Lastly, spin systems and neural networks can be studied on complex networks that have non-trivial topologies, similar to gene regulatory networks. The fields of the statistical physics of spin glasses and neural networks are well established. Thus, given the previously listed similarities between magnets, memories and cell types in relation to the configuration of interacting units that comprises

them, it is only logical to borrow from the library of results, tools, and techniques of statistical physics. In fact, the similarity between cell types and spin glasses has been made before. Derrida and Flyvbjerg first noted the similarity between Kauffman's models and spin-glasses [80].

Neural network like models have also been applied to the problem of cell reprogramming before. Lang et al. adapted the Hopfield model to construct a complex epigenetic landscape from a protein-protein interaction network [81]. In their model different cell types exist as attractors in the high dimensional epigenetic landscape. They were successfully able to reproduce known reprogramming protocols and showed that partially reprogrammed cells, commonly observed in experiments, emerge directly from the dynamics of their model converging to a state that is a mixture of multiple attractors. Despite the success of their model, it lacks some key aspects of cell biology, such as variation in gene expression levels due to cell cycles. Their model is also hard to interpret from an experimental point of view. For example, they model the effect of culture conditions or different transcription factors with field contributions in the Hamiltonian of their model. Then, to model the effect of the OSKM factors during reprogramming experiments, they send the corresponding components of the transcription factor field to infinity whilst setting all other components to zero. Whilst their model is able to accurately capture experimental observations, it is difficult to interpret from a biological point of view.

Similarly, Szedlak et al. have used Hopfield-like networks to model the cell cycle [82] and explore signalling patterns in cancerous cells [83]. Once more, whilst these models are capable of capturing aspects of the behaviour that they are attempting to model, the parameters of the models can be difficult to interpret in terms of measurable quantities. Therefore, it is also worth emphasising that genes are not neurons or magnetic moments. Careful considerations should be made when constructing models for the interaction of their expression levels by leaning on known results from theoretical and experimental biology. Parameters and interactions should have purposeful meanings that are rooted in biological facts. It is not sufficient simply to change the meaning of variables in these models to suit one's own purpose. Thus, success relies on build-

ing models that are motivated from molecular biology, whilst also learning lessons from known results in statistical physics where we can.

# 2

# A HIERARCHICAL MODEL FOR CELL

# REPROGRAMMING

## 2.1. INTRODUCTION

The retrieval of pluripotent cells was first pioneered by John Gurdon in the 1960s, using nuclear transfer to clone a frog from the nuclei of somatic cells extracted from a xenopus tadpole [14]. More recently, cell reprogramming has shown that it is possible to obtain induced pluripotent stem cells (iPSCs), which strongly resemble embryonic stem cells (ES), from somatic cells via the introduction of just 4 transcription factors (Oct3/4, Sox2, Klf4 and c-Myc - now known as

the Yamanaka or OSKM factors) [17, 18]. It has also been demonstrated that nearly all somatic cells can be reprogrammed in this manner [19], suggesting that the "code" for pluripotency lies in the genome common to all cells of an organism. Once reprogrammed it is possible to guide iPSCs to differentiate into a desired cell type using specific culture conditions [84]. Due to their ability to self renew and differentiate into many different cell types, stem cells (including iPSCs) hold great potential for both personalised and regenerative medicine [2]. Furthermore, iPSCs can act as a model environment for studying disease and testing drug delivery mechanisms [3, 5]. Since the original reprogramming experiments, multiple protocols have been uncovered by replacing certain Yamanaka factors with other proteins or small molecules. For an extensive biomedical review of iPSCs, and how they differ from other stem cells, see Takahashi 2015 [9].

However, despite the great potential of iPSCs, and the evolution of cell reprogramming experiments, much is still unknown about the decisions governing the fate of a cell. Cell fate decisions were first modelled by Waddington using his idea of an epigenetic landscape. This model describes development using the analogy of a ball rolling down a hill from states of high potency to fully differentiated ones. Different cell types are represented as valleys in the landscape, and a cell's fate is determined by the valley which the ball falls into [85]. The number of valleys increases the further the ball moves down the landscape, representing the increasing diversity of cell types during development. Whilst this model provides an interesting metaphor for differentiation, it lacks some key aspects of cell biology, such as cell cycles. Recent experimental work has suggested that cell types can be considered as high dimensional attractors of a gene regulatory network [33], paving the way for a dynamical systems approach to cell reprogramming.

Many of the current models describing cell fate decisions focus on specific small gene regulatory networks (GRN) that are believed to govern pluripotency or differentiation, or, they approach the problem from a cell population perspective. For further details of the current mathematical and computational models of cell reprogramming the reader is directed towards the references Morris et al. 2014 [7] and Herberg and Roeder 2015 [49] respectively.

In this chapter, a mathematical model is presented which models cell reprogramming in terms of transitions between attractors of a high dimensional dynamical system. The attractors of the dynamics represent the gene expression levels throughout the cell cycles of different cell types, and they are related to one another in a hierarchical manner. The rest of this chapter is organised as follows: first, the theory behind the model is formulated by appealing to a small set of key observations concerning cell chemistry, before being applied to a specific type of hierarchy. Next, evidence for reprogramming in the model is presented and discussed, with the main findings of the work summarised at the end of the chapter. The majority of the mathematical details have been relegated to the appendices with the aim of making this chapter accessible to interested readers from various scientific backgrounds.

## 2.2. Theory

Cells are the fundamental units of structure and reproduction in most organisms [86]. They are complex and dense building blocks which contain a rich tapestry of biochemical reactions involving a multitude of chemical species (e.g. proteins, sugars, lipids, etc.). Metabolic pathways, such as glycolysis, involve many intermediate steps converting the product of one reaction into the substrate for another. Enzyme reactions, like those involved in glycolysis, can be described generally by a set of differential equations, known as the Michaelis-Menten equations [87]. Thus, to fully describe the dynamics of cell chemistry one would need to incorporate the Michaelis-Menten equations for all possible reactions into a theory which would describe reaction and diffusion mechanisms, self-organisation, biochemical signalling, etc. One component of such a theory would be the transcription of the $N$ genes of the organism's genome, which alone represents a vast state space. For example, even if one assumes binary gene expression levels (i.e. genes are either expressed or not expressed), there are $2^N$ possible configurations of gene expression levels. A natural question then arises from the complex chemistry of cellular life: How do so many reactions, of so many species, give rise to a comparably low number of different cell types? For

example, the human genome is comprised of approximately $25,000$ genes yet only gives rise to around $300$ different cell types. A plausible realisation of this fact is to suppose that stable cell types emerge as attractors of the full reaction dynamics of the cell.

For the purpose of modelling cell reprogramming, we propose to construct a reduced model by appeal to the following line of reasoning. Suppose one was able to integrate out all components of the complete theory other than gene expression levels, the result would be a reduced model which will have the following two features: ($i$) it will involve interactions between genes; ($ii$) the interactions will exhibit memory effects. The interaction of genes would result in a feedback mechanism that could explain the existence of stable attractors. In the reduced model, memory would be a result of the interplay between genes and proteins. Transcription factors (TFs) are proteins that regulate the expression of genes (through activation/inhibition). These proteins are translated from RNA, which is transcribed from the genes in the cell's nucleus. Thus, the expression level of a gene will depend on the previous expression levels through gene regulation. Furthermore, proteins can regulate the genes which they were synthesised from, other genes, and/or combine with other proteins to form complexes which are transcription factors. Hence the expression level of a given gene will depend on the previous expression levels of many (or all) other genes. Memory is, in fact, *required* to create dynamic cell cycle attractors with different durations for each of the phases of the cell cycle, in a model based on gene expression levels only.

Based on these observations, we build a *minimal* model that describes cell types in terms of gene expression levels across their cell cycles. We can make simplifications to the reduced model, that do not change the intuition behind, or nature of, the model but make the mathematics easier to work with. One such simplification is the discretisation of time, which allows one to neglect the effects of memory. To do this we measure time in terms of stages passed through the cell cycle (e.g. $G_1$, $S$, $G_2$, ...). This allows one to ignore the different durations of each cycle phase by concentrating on which phase of the cycle a cell is in. Another assumption is that the gene expression levels are binary variables, $n_i$ (with $i = 1 \ldots N$), i.e. genes can exist in one of only

two states: they are either expressed or are not. These states may be represented by the binary values $n_i = 1$ and $n_i = 0$ respectively, hence the common terminology Boolean, or "on/off", genes. Again, it is important to stress that these assumptions make the mathematics of the model much simpler, but can be relaxed if a more comprehensive description of cell cycle regulation is required.

A general model for the dynamics of interacting binary genes would have the following form,

$$n_i(t+1) = \Theta\left[h_i(t) - \theta_i - T\xi_i(t)\right] , \tag{2.1}$$

where $n_i$ is the gene expression level of the $i^{\text{th}}$ gene, with the effect of the gene interactions encoded in a local field, $h_i(t)$ of the form,

$$h_i(t) = \sum_j J_{ij} n_j(t) + \sum_{j,k} J_{ijk} n_j(t) n_k(t) + \dots . \tag{2.2}$$

Here $J_{ij}$ is the effect of the interaction between genes $i$ and $j$, and $J_{ijk}$ is likewise the effect of the triplet interactions between the 3 genes $i$, $j$ and $k$, (there could also be higher order interactions which are represented by the ... ). Any constant contributions to the local field, such as self regulation, can be absorbed into the definition of $\theta_i$. The $\xi_i$ are random variables with zero mean and a suitably normalised variance, which mimic noise to represent the fundamental stochasticity of reaction events. Popular noise models are Gaussian and thermal noise. We use $T$ to vary the strength of the noise. Anticipating our later choice of the thermal noise model, I may refer to $T$ as temperature or noise strength interchangeably. The $\Theta[x]$ is the Heaviside step function: $\Theta[x] = 1$ for $x > 0$ and $\Theta[x] = 0$ otherwise. Thus, (2.1) states that a gene will be expressed in the next phase of the cell cycle (i.e. $n_i(t+1) = 1$) if the combined effect of all interactions and stochastic noise at time $t$ exceeds a gene-specific threshold, $\theta_i$. At each time step every gene expression level is updated according to this rule, and the state of the system is fully described at any time $t$ by the *instantaneous* configuration $\mathbf{n}(t) = (n_1(t), \dots, n_N(t))$ of gene expression levels.

The network of effective interacting gene expression levels one should consider is a subset of the entire genome. Only the regulatory genes (i.e. the genes that encode for proteins that are transcription factors or form complexes that are transcription factors) need to be considered. Whilst the expression of other genes may contribute to identifying a given cell fate, their expression is driven by that of the regulatory genes. That is the expression of regulatory genes is independent of that of non-regulatory genes. Thus, the total number of genes $N$ in our model should be thought of as the total number of *regulatory* genes.

### 2.2.1. MINIMAL MODEL

We restrict ourselves to consider a system involving pair interactions only, and simplify matters further by assuming uniform thresholds, i.e. $\theta_i = \theta \; \forall \; i$. Thus the dynamics of the minimal model is given by the following simple expression,

$$n_i(t+1) = \Theta \left[ \sum_j J_{ij} n_j(t) - \theta - T\xi_i(t) \right] . \tag{2.3}$$

This expression is reminiscent of the models used in the field of neural networks (NNs) for associative memory, with a post-synaptic potential (PSP), $h_i = \sum_j J_{ij} n_j$ and a neuron fires ($n_i(t+1) = 1$) if the PSP exceeds a given threshold $\theta$. In associative memory, configurations of neuronal activity representing some memories are stored in the synaptic efficacies $J_{ij}$, such that they are attractors of the dynamics. The NN is then said to recall a pattern when the system converges to the corresponding configuration from some initial condition. Such NNs are said to be content addressable because the attractor to which they converge is given by the (content of the) initial state. Associative NNs of this type are robust to input errors and hardware failures (such as disruption of the synaptic efficacies and thresholds). Analogously, in our model, specific configurations of gene expression levels, which represent the different cell types of an organism, are stored in the gene interactions, $J_{ij}$, which therefore govern the dynamics. Using a temporal ordering of

the cycle state specific configurations, the attractors become dynamic attractors that represent the cell cycles of each cell type (as will be shown in detail in section 2.3). It is also desirable that the gene interaction network is robust to variation or errors in gene expression levels, because, despite variation in gene expression levels across human individuals and mutations in the human genome, individuals of a population have the same set of cell types.

The analogy between associative memory and cell reprogramming has been made recently in the context of protein interaction networks [81]. However, that work differs from that presented here due to the absence of cell cycles and the potency hierarchy. In their work, Lang et al. extend Waddington's developmental landscape metaphor into an epigenetic landscape describing the interactions between proteins. Instead, by working with gene interactions it is possible to neglect the specific activation/inhibition nature of transcription factors, which will be encoded in the effective interactions between genes. Thus, the interaction between genes, as opposed to proteins, is a more natural approach to model cell fates. However, in chapter 4, the combined dynamics of gene expression levels and TFs is studied to provide a deeper understanding of the gene regulatory dynamics encoded in the effective interactions between gene expression levels.

### 2.2.2. CELL CYCLE SIMILARITIES, LINEAGES AND REPROGRAMMING

Cell reprogramming requires transitions between cell states, which can be either a trans- or de-differentiation, i.e. either across or up the potency cascade [9]. With this idea in mind, the cell cycles stored in the model are related to one another in a hierarchical manner. Specifically - apart from the stem state which sits at the top of the hierarchy and thus has no ancestor (see figure 2.2.1) - the gene expression levels of all cell types are conditionally dependent on their parents. This set up is inspired by the storage of memories, in a Markovian hierarchy, for associative memory NNs [88]. Parallels between the hierarchical relation of cell cycles and Waddington's epigenetic landscape can be made. However, the present approach does not directly model a landscape of the cell states.

**Figure 2.2.1:** A schematic diagram of the hierarchy of cell types, in terms of cell potency. Stem cells, e.g. embryonic stem cells (ES) or iPSCs, sit at the top of the hierarchy due to their ability to differentiate into many different cell types. The further down the hierarchy a cell type is, the lower its level of potency (or higher its level of specialisation). Differentiation corresponds to moving down one level of the hierarchy (green arrow); de-differentiation is equivalent to moving up the hierarchy (red arrow); trans-differentiation (blue arrow) corresponds to transitions between cell types of the same level.

It has been demonstrated that protein and mRNA levels vary across the cell cycle [41], thus it is reasonable to infer that the gene expression levels of a cell also changes throughout its cycle. It is thus plausible to conceive of a situation in which the global expression levels of different cell types are more similar in certain stages of the cell cycle than in others. For example, during the S-phase the gene expression levels could likely be vastly reduced in all cells types (as suggested by Cho et al. [89]), as the DNA is otherwise occupied through replication. It is also the case that most of an organism's cells undergo the mitotic phase via broadly the same mechanisms. Hence, there could exist (at least) one phase of the cell cycle in which different cell types are more similar than others (see figure 2.2.2). These stages of the cell cycle would represent a *natural* target in which it is easier to induce switches between cell types, i.e. to reprogram a cell. This is one of the main hypotheses that we will be testing in the present study.

At the time of writing, I am aware of only one other model which includes a hierarchy of cell

**Figure 2.2.2:** A schematic diagram of the possible similarities between the cell cycles of two different cell types (labelled A and B). The horizontal distances between the analogous phases of the two cell cycles represent the level of similarity - the closer the phases the more similar they are. Two different cell types could be more similar during the S- and/or M-phases, in which the biological processes are broadly similar across different cell types of an organism.

states and the cell cycle. In their model, Artyomov et al. defined a cell type through the expression levels of a small ensemble of master regulatory genes referred to as a *module* [40]. They included the cell cycle as an interplay between the gene expression levels of a cell and the epigenetic state of the cell. On the other hand, here, we treat each cell type as a dynamic entity, which transitions through different configurations of gene expression levels that correspond to stages of its cell cycle. Each configuration describes the *entire* transcriptome of a cell in a given cycle phase. In this work we do *not* directly model the epigenetics of a cell. However, the similarity between different cell types, during specific phases of the cell cycle discussed above, could be a result of epigenetic changes, such as changes in chromatin structure or the presence of histone markers.

### 2.2.3. Two-level hierarchy

To validate the principles of our approach we apply the neural-network-like model above to a simplified version of the biology. This makes the mathematics easier to implement and keeps the

Stem cell



$$\eta^\rho$$

$$\eta^{\rho,1} \quad \eta^{\rho,2} \quad \eta^{\rho\mu}$$

corresponding ρ-phase progeny

**Figure 2.2.3:** A cell potency hierarchy for a two-level system. The stem cell sits at the top of the hierarchy and is given by the configuration of gene expression levels $\eta^\rho$, where the superscript $\rho$ labels the stage of the cell cycle (e.g. S-phase). The second level of the hierarchy consists of $M$ daughter cell configurations. In general, the configuration of the daughter cell is given by $\eta^{\rho\mu}$, where the superscript $\rho$ and $\mu$ label the stage of the cell cycle (e.g. S-phase) and the type of daughter cell respectively (e.g. neuron, B-cell, etc.). A hierarchy of this form exists for every stage of the cell cycle, so that, every cell type has the same number of cycle phases.

notation simple and transparent. Such simplified scenarios still capture the main principles of the biology, and in practice, the mathematics can easily be extended to more realistic systems. We therefore consider a two level hierarchy in which fully differentiated cells are direct descendants, or daughters, of the stem cell (see figure 2.2.3). Each of the cell cycles is coarse grained into 3 stages, with a single cycle phase made more similar across the different cell types. The coarse graining of the cell cycles was carried out purely for computational efficiency and generalisation to 4 or 5 stage cell cycles is straightforward.

## 2.3. Probabilistic framework

We will now introduce a specific model which implements our generic reasoning within a probabilistic framework. Consider a system of $N$ genes, each labelled by $i = 1 \ldots N$, and $M$ daughter cell types, labelled by $\mu = 1 \ldots M$. Each cell type, daughter or stem, undergoes a cell cycle of

length $C$ and we denote each phase of the cycle by $\rho = 1 \ldots C$ (with $C + 1 \equiv 1$). The expression of the $i^{\text{th}}$ gene, in the $\rho^{\text{th}}$ phase of the stem cell cycle, is denoted by $\eta_i^\rho \in \{0, 1\}$, where $\eta_i^\rho = 1$ corresponds to that gene being expressed. We denote by $a^\rho$ the fraction of genes expressed during the $\rho^{\text{th}}$ cell cycle phase, also referred to as the activity of that cycle phase. Thus, the probability of the gene expression levels in the $\rho^{\text{th}}$ phase of the stem cell cycle are given by

$$p(\eta^\rho) = \begin{cases} a^\rho & \text{for } \eta^\rho = 1 \,, \\ 1 - a^\rho & \text{for } \eta^\rho = 0 \,. \end{cases} \tag{2.4}$$

Then the configuration of the stem cell state, in the $\rho^{\text{th}}$ cycle phase, is given by $\boldsymbol{\eta}^\rho = (\eta_1^\rho, \ldots, \eta_N^\rho)$. Furthermore, for every state, $\boldsymbol{\eta}^\rho$, of the stem cell cycle there is a corresponding set of descendants, in the same stage of the cell cycle, $\boldsymbol{\eta}^{\rho\mu} = (\eta_1^{\rho\mu}, \ldots, \eta_N^{\rho\mu})$. Similar to the stem cell cycle, each daughter cell $\mu$ has an activity given by $a^{\rho\mu}$, which governs the probability of expressing a gene in each stage of the cell cycle

$$p(\eta^{\rho\mu}) = \begin{cases} a^{\rho\mu} & \text{for } \eta^{\rho\mu} = 1 \,, \\ 1 - a^{\rho\mu} & \text{for } \eta^{\rho\mu} = 0 \,. \end{cases} \tag{2.5}$$

We assume the gene expression levels in the stem cell, $\eta_i^\rho$, are independent, identically distributed (i.i.d) random variables, which implies that the configurations $\boldsymbol{\eta}^\rho$ are independent along the cell cycle. This assumption was made to simplify the mathematics, but may be relaxed if a more comprehensive model is desired. The configurations of the daughter cells, on the other hand, are derived from the corresponding phases of stem cell. We define the probability of turning a gene off during the differentiation-transition from a stem cell to a daughter cell, in the same phases of the cell cycle, with equal activities $a^\rho = a^{\rho\mu}$, as $\gamma^{\rho\mu}$. Thus, the transition probability of a gene being expressed in both the stem and daughter cell (of the same cell cycle phase) is given by $(1 - \gamma^{\rho\mu}) \frac{a^{\rho\mu}}{a^\rho}$, i.e. the ratio of probabilities the gene is on in both states multiplied by the proba-

bility it was not turned on in differentiation. The full transition matrix used to construct the gene expression levels of the daughter cells from the same phase of the stem cell cycle can be found in appendix 2.A. Due to this construction of the daughter cell cycles in terms of the stem cell cycle, the different daughter cell configurations are *conditionally dependent* on the stem cell state.

The interactions between the gene expression levels of the system should be chosen such that the cell cycles of the daughter and stem cells are attractors of the dynamics. To construct the interactions of the model we combine a set of known results from the field of NNs. Hopfield originally showed that multiple configurations can be stored in the synaptic couplings of a neural network using the Hebb rule [66]. It is known that dynamic attractors can be stored in the couplings by adapting the Hebb rule to include a temporal order to the stored configurations, i.e the interactions have a contribution from the current pattern and its successor [90, 91]. Thus, in our notation of gene expression levels, a sequence of stem cell cycle phases may be stored in the interactions in the following manner.

$$J_{ij}^{(cycle)} = \frac{1}{N} \sum_{\rho=1}^{C} \frac{\left(\eta_i^{\rho+1} - a^{\rho+1}\right)\left(\eta_j^\rho - a^\rho\right)}{a^\rho\left(1 - a^\rho\right)} \, . \tag{2.6}$$

This choice of interaction ensures that if the gene expression levels evolve according to (2.3) and are initialised, at time $t$, in the configuration $\mathbf{n}(t) = \boldsymbol{\eta}^\rho$, their configuration in the next time step will, with high probability, be $\mathbf{n}(t+1) = \boldsymbol{\eta}^{\rho+1}$. To ensure the sequence of configurations retrieved by the system is a closed cycle, the successor of the final configuration must be equivalent to the initial configuration. For low activity configurations, it is required that one removes the bias from each of the cycle phases, in order to achieve stable limit cycle attractors in the dynamics. Here this is done by subtracting the average gene expression of each of the cell cycle phases, resulting in the contributions from each cycle phase having zero mean.

Information can also be stored in the interactions in a hierarchical manner [88, 92, 93] (equivalent to the structure shown in figure 2.2.3). This is done by including contributions from each

state in the hierarchy in the interactions. Each pattern must then be weighted by a factor determined by its position in the hierarchy [94]. Combining these two ingredients, the interactions that stabilise a hierarchy of cell cycles can be written as follows,

$$J_{ij} = \frac{1}{N} \sum_{\rho=1}^{C} \left\{ \frac{(\eta_i^{\rho+1} - a^{\rho+1})(\eta_j^{\rho} - a^{\rho})}{a^{\rho}(1 - a^{\rho})} + \sum_{\mu=1}^{M} \frac{(\eta_i^{\rho+1,\mu} - a_{\mu}(\eta_i^{\rho+1}))(\eta_j^{\rho\mu} - a_{\mu}(\eta_j^{\rho}))}{a^{\rho\mu}(1 - a^{\rho\mu})} \right\} . \quad (2.7)$$

Here the summations are over cycle phases, $\rho$, and daughter cell types, $\mu$. We chose to remove the bias from the daughter cell type by subtracting the conditional average of the gene expression levels, $a_{\mu}(\eta^{\rho}) = \mathbb{E}\left[\eta^{\rho\mu}|\eta^{\rho}\right]$, i.e. the average gene expression level of the daughter cell given the expression levels in the same cell cycle phase of the stem cell. However, the bias could also be removed using the activity of the daughter cells, $a^{\rho\mu}$, in place of the conditional averages in (2.7). The weights in the denominators are the variances of the gene expression levels in the corresponding cell cycle stage for the stem or daughter cells. Note that, if the $\rho+1$ was replaced with $\rho$ in (2.7) then this would be the standard prescription for storing a hierarchy of configurations. Since they are included we have in fact stored a *hierarchy of cell cycles*.

### 2.3.1. INTRODUCING THE DYNAMICS

Given the form (2.7) of the interactions, one can express the local fields $h_i(t)$, appearing in the dynamics (2.3), concisely in terms of a set of macroscopic dynamical order parameters, namely

$$\widetilde{m}_{\rho}(t) = \widetilde{m}_{\rho}(\mathbf{n}(t)) = \frac{1}{N} \sum_{i=1}^{N} \frac{\eta_i^{\rho} - a^{\rho}}{a^{\rho}(1 - a^{\rho})} n_i(t) , \quad (2.8)$$

$$\widetilde{m}_{\rho\mu}(t) = \widetilde{m}_{\rho\mu}(\mathbf{n}(t)) = \frac{1}{N} \sum_{i=1}^{N} \frac{\eta_i^{\rho\mu} - a_{\mu}(\eta_i^{\rho})}{a^{\rho\mu}(1 - a^{\rho\mu})} n_i(t) . \quad (2.9)$$

If we absorb the threshold $\theta$, appearing in (2.3), into the definition of $h_i(t)$, this gives

$$h_i(t) = \sum_\rho \left\{ (\eta_i^{\rho+1} - a^{\rho+1})\widetilde{m}_\rho(t) + \sum_\mu \left[\eta_i^{\rho+1,\mu} - a_\mu(\eta_i^{\rho+1})\right] \widetilde{m}_{\rho\mu}(t) \right\} - \theta \,. \qquad (2.10)$$

The value of $\theta$ is fixed, such that the stable cell cycle attractors exist at sufficiently low noise levels. Note that $h_i(t) = h_i(\widetilde{\mathbf{m}}(t))$, in which $\widetilde{\mathbf{m}}$ is the vector containing the set of all dynamical order parameters $\{\widetilde{m}_\rho(\mathbf{n}(t))\}$ and $\{\widetilde{m}_{\rho\mu}(\mathbf{n}(t))\}$. As a consequence, and by appeal to the law of large numbers in the limit of a large number $N$ of genes, one can formulate the dynamics of our model in closed form *entirely* in terms of the dynamic order parameters, giving

$$\widetilde{m}_\rho(t+1) = \left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1-a^\rho)} \mathbb{P}\left[\xi \leq \beta h(t)\right] \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} , \qquad (2.11)$$

$$\widetilde{m}_{\rho\mu}(t+1) = \left\langle \frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1-a^{\rho\mu})} \mathbb{P}\left[\xi \leq \beta h(t)\right] \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} , \qquad (2.12)$$

provided that $N \gg M$ in this limit (we will see later why this condition must be satisfied). In (2.11) and (2.12) $\beta = 1/T$ is the inverse of the noise strength. The $\mathbb{P}(\xi \leq z)$ in (2.11) and (2.12) is the cumulative distribution function (CDF) for the noise probability, $p(\xi)$ (i.e. the probability that $\xi$ will take a value less than or equal to $z$). Popular choices for the $p(\xi)$ are the Gaussian distribution, and the qualitatively and quantitatively similar, logistic distribution $\mathbb{P}(\xi \leq z) = \frac{1}{2}(1 + \tanh \frac{z}{2})$. We will use the latter, for which (2.11) and (2.12) can be written in the following form (for details of this calculation see appendix 2.B),

$$\widetilde{m}_\rho(t+1) = \frac{1}{2}\left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1-a^\rho)} \tanh\left(\frac{\beta h(t)}{2}\right) \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} , \qquad (2.13)$$

$$\widetilde{m}_{\rho\mu}(t+1) = \frac{1}{2}\left\langle \frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1-a^{\rho\mu})} \tanh\left(\frac{\beta h(t)}{2}\right) \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} , \qquad (2.14)$$

where the angle brackets, $\langle \ldots \rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}}$, represent the average and conditional averages over all stem and daughter cycle states. These equations of motion are easily solved numerically by forward iteration, starting from suitable initial conditions.

### 2.3.2. SIGNAL-TO-NOISE ANALYSIS

Models of NNs, like those that inspired the choice of interactions leading to (2.7), typically have a finite storage capacity. That is, there is an upper limit on the number of patterns that can be stored in the interactions. If one exceeds this storage capacity, the model will enter a spin glass regime where it is no longer capable of accurately recalling any patterns. Here we use a signal-to-noise analysis for (2.7), in order to assess the criteria under which a hierarchy of cell cycles can be stored with this type of modelling procedure.

The local fields $h_i = \sum_j J_{ij} n_j - \theta_i$ determine the evolution of the system. In a signal-to-noise analysis the local fields are dissected into two parts that represent the contribution to the couplings that arises from the condensed pattern $\mathbf{n}$ (i.e. the pattern the system currently finds itself in) known as the signal, and the contribution to the couplings from all other patterns termed the noise. Commonly, the condensed pattern $\mathbf{n}$ is assumed to be the only configuration that the system has a finite correlation with in the thermodynamic limit. By studying the noise and signal contributions to the local field, we can determine the conditions under which the model is successfully able to recall the condensed pattern [95], or, in the case of the gene expression level dynamics, progress through the cell cycle.

If we chose the condensed pattern to be $\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}}$, then the local fields (in the absence of any

gene-specific thresholds) can be decomposed as,

$$
\begin{aligned}
h_i(\boldsymbol{\eta}^{\bar{\rho}}) = & \frac{\eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1}}{N} \sum_{j(\neq i)} \frac{\eta_j^{\bar{\rho}} - a^{\bar{\rho}}}{a^{\bar{\rho}}(1 - a^{\bar{\rho}})} \eta_j^{\bar{\rho}} \\
& + \sum_{\mu} \frac{\eta_i^{\bar{\rho}+1,\mu} - a_\mu(\eta_i^{\bar{\rho}+1})}{N} \sum_{j(\neq i)} \frac{\eta_j^{\bar{\rho}\mu} - a_\mu(\eta_j^{\bar{\rho}})}{a^{\bar{\rho}\mu}(1 - a^{\bar{\rho}\mu})} \eta_j^{\bar{\rho}} \\
& + \sum_{\rho(\neq\bar{\rho})} \left\{ \frac{\eta_i^{\rho+1} - a^{\rho+1}}{N} \sum_{j(\neq i)} \frac{\eta_j^{\rho} - a^{\rho}}{a^{\rho}(1 - a^{\rho})} \eta_j^{\bar{\rho}} \right. \\
& \left. + \sum_{\mu} \frac{\eta_i^{\rho+1,\mu} - a_\mu(\eta_i^{\rho+1})}{N} \sum_{j(\neq i)} \frac{\eta_j^{\rho\mu} - a_\mu(\eta_i^{\rho})}{a^{\rho\mu}(1 - a^{\rho\mu})} \eta_j^{\bar{\rho}} \right\}, \quad (2.15)
\end{aligned}
$$

where the top lines and bottom lines are the contribution to the local fields from the signal $S$ and the noise $R$ respectively. The noise includes contributions to the local field from all cell cycle phases different to the condensed pattern, i.e. all $\rho \neq \bar{\rho}$. In the limit of a large number of regulatory genes, $N$, the signal contribution to the local fields is simply $(\eta^{\bar{\rho}+1} - a^{\bar{\rho}+1})$, and acts to progress the stem cell cycle from one phase to the next according to the dynamics (2.3) - further details can be found in appendix 2.D. In the same limit, the contribution to the local field from the noise is a sum of random variables, and is itself a zero-mean random variable. Thus, using the central limit theorem, the noise distribution is a zero mean Gaussian that is completely characterised by its variance. The variance of the noise is found by squaring the bottom lines of the local field above, and averaging over all gene expression level patterns (see appendix (2.D)). It can be shown that the variance of the noise contribution to the local fields behaves as

$$
\langle R^2 \rangle \propto \frac{(C-1)M}{N} . \tag{2.16}
$$

Therefore, for the signal to dominate in the local fields, allowing successful progression of the cell cycle, $\langle R^2 \rangle$ should be small. This occurs provided that $N \gg M(C-1)$ (see appendix 2.D for full details). Remarkably, this condition is in line with what is typically observed in multicellular organisms, e.g. for Humans $N \sim 25,000$, $M \sim 300$ and $C = 5$. Furthermore, the same

**Figure 2.4.1:** Numerical solutions of (2.13) and (2.14), at a low effective temperature $T = 0.01$, when the system is initialised with a high overlap with the daughter cell cycle. The peaks in $m_\rho$ and $m_{\rho\mu}$ correspond to the system passing through states correlated with the different cell cycle stages of the daughter and stem cell - i.e. each peak is a successive $\rho$, and the dashed lines are the overlap with a single daughter type, $\mu$, and the solid lines are the overlap with the stem cell. The following parameters were used: $C = 3$, $a^1 = a^3 = 0.7$, $a^2 = 0.3$, $a^{1\mu} = a^{3\mu} = 0.6$, $a^{2\mu} = 0.2$, $\gamma^{\rho\mu} = 0.2$ for all $\rho$ and $\mu$ and $\theta = 0$.

condition is also found when a daughter cell cycle is chosen as the condensed pattern.

## 2.4. Results

In the results that follow a single cell cycle phase was made more similar between the daughter and stem cell cycles. This was done by having a lower value of $a^\rho$ and $a^{\rho\mu}$ in that particular phase. The similarity between the cycle phases of the stem and daughter cells can be seen from the covariance between their expression levels,

$$\mathrm{cov}\left[\eta^\rho, \eta^{\rho\mu}\right] = a^{\rho\mu}(1 - \gamma^{\rho\mu} - a^\rho).$$ (2.17)

Here, the probabilities of expressing a gene in both cell cycle stages ($a^{\rho\mu}$ and $a^\rho$ respectively) and the probability that a gene is not turned on during differentiation (i.e. $1 - \gamma^{\rho\mu}$) govern the

similarity between the gene expression levels of the same cycle phase across the two generations of cell types. Thus, it is possible to make certain stages of the cell cycle more similar across the two levels of the hierarchy by tuning the parameter values used in (2.17). At this point, it should also be noted that there are restrictions on the values that $\gamma^{\rho\mu}$ can take, in order for the transition probability from $\boldsymbol{\eta}^{\rho}$ to $\boldsymbol{\eta}^{\rho\mu}$ to be correctly defined as a probability (for details see the appendix 2.A).

Our choice of the parameters is based on the analysis carried out in Ramskold 2009 that suggests that 60-70% of all genes are expressed in human cells [22]. In addition, results in the literature suggest that the activity in stem cells is higher than its progeny [96], so we will choose $a^{\rho} > a^{\rho\mu}$ for all $\rho$, $\mu$. This leads us to a choose $a^{\rho} = 0.7$ and $a^{\rho\mu} = 0.6$ in all cells and all phases except for the phase $\rho = 2$, that we aim to make more similar between the stem and daughter cells. In this phase, we chose the parameters $a^{\rho} = 0.3$ and $a^{\rho\mu} = 0.2$, which lead, via (2.17), to a higher similarity between the stem and daughter cells. This is also consistent with the expectation that gene expression is lower in the S-phase due to the genome being occupied with other processes, such as DNA synthesis. The value $\gamma^{\rho\mu} = 0.2$ was used for all cell cycle phases, $\rho$, and daughter cell types, $\mu$, and the threshold values were set to zero ($\theta = 0$). Unless stated otherwise, these parameter values are used in all of the results and analysis that follow in the remainder of this chapter.

For the presentation of the results we use the so-called overlaps, which measure the correlation between the state of the system $\mathbf{n}(t)$ and the gene expression patterns that are characteristic of the cell cycle states of the stem and daughter cells, respectively. They are closely related to the dynamic order parameters, and in fact identical for the stem cell cycles, and are defined as,

$$m_{\rho}(\mathbf{n}(t)) = \widetilde{m}_{\rho}(\mathbf{n}(t)) \,, \tag{2.18}$$

$$m_{\rho\mu}(\mathbf{n}(t)) = \frac{1}{N} \sum_{i=1}^{N} \frac{\eta_i^{\rho\mu} - a^{\rho\mu}}{a^{\rho\mu}(1 - a^{\rho\mu})} n_i(t) \,, \tag{2.19}$$

The overlaps are normalized to have values in the interval $[-1, 1]$. An overlap of $m_{\rho\mu} = 1$ (or $-1$) means the system is fully correlated (or anti-correlated) with the cell type $\mu$, in the cell cycle phase $\rho$, whereas $m_{\rho\mu} = 0$ implies that they are completely uncorrelated. The same is true for the stem cell cycle and corresponding values of $m_\rho$.

In figure 2.4.1 we plot the numerical solutions of (2.13) and (2.14) when the system is initialised in a daughter cell cycle, at a low noise level, $T = 10^{-2}$. Peaks of the dashed line correspond to the system transitioning through states with a high overlap with the daughter cell cycle phases $m_{\rho\mu}$. Similarly, peaks in the solid line correspond to the overlaps with different phases of the stem cell cycle, $m_\rho$.

Because of correlations between the gene expression patterns of the same cycle phases of stem and daughter cells, one observes non-zero mutual overlaps between them (for details see appendix 2.E). Specifically, if the system is in a (perfect) daughter cell state, $n_i = \eta_i^{\rho\mu} \, \forall i$, the overlap with the corresponding stem cell state is

$$m_\rho(\mathbf{n}(t) = \boldsymbol{\eta}^{\rho\mu}) = \left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1 - a^\rho)} \eta^{\rho\mu} \right\rangle = \frac{\text{cov}\left[\eta^\rho, \eta^{\rho\mu}\right]}{\text{var}\left[\eta^\rho\right]} \,. \tag{2.20}$$

Conversely if the system is in a stem cell state, $n_i = \eta_i^\rho$, the overlap with the daughter cell state $\eta_i^{\rho\mu}$ is

$$m_{\rho\mu}(\mathbf{n}(t) = \boldsymbol{\eta}^\rho) = \left\langle \frac{\eta^{\rho\mu} - a^{\rho\mu}}{a^{\rho\mu}(1 - a^{\rho\mu})} \eta^\rho \right\rangle = \frac{\text{cov}\left[\eta^\rho, \eta^{\rho\mu}\right]}{\text{var}\left[\eta^{\rho\mu}\right]} \,. \tag{2.21}$$

The peaks corresponding to the phase $\rho = 2$ in figure 2.4.1 are higher than those corresponding to the other two phases, because the gene expression activity in this phase was chosen to have a higher covariance, hence a higher value of the overlap (2.20), between the stem and daughter cells. The other two phases have identical gene expression activities, hence they have identical values for their overlaps. The initial value of $m_\rho(\mathbf{n}(0))$ in figure 2.4.1 was determined using (2.20).

**Figure 2.4.2:** Left: Numerical solutions of (2.13) and (2.14), at a noise level $T = 0.14$. The initial condition for the overlap with the stem cell was determined using (2.20). Right: Monte-Carlo simulation dynamics at the same temperature for $N = 25,000$ genes. The system was initialised in a configuration with a high overlap with the $\rho = 1$ phase of the daughter cell $\mu$, but as the dynamics progress this decays and the system converges to a high value for the overlap with the stem cell cycle. This transition takes multiple generations of the cell cycle and the system passes through an intermediate state with equal overlap with both stem and daughter cell cycles where the two lines intersect. Only the envelope of the trajectories is shown, i.e. the cycle phase, $\rho(t) = 1 + (t \mod C)$, which the system is expected to be in. For both panels the same parameter values were used as in figure 2.4.1.

### 2.4.1. NOISE INDUCED SWITCHING

At a low noise level, if the system is initialised in a daughter cell it will transition along that cell cycle indefinitely. However, as $T$ is increased above some critical value the noise will take the system away from the daughter cell and it will fall into the attractor corresponding to the stem cell cycle. If the noise is then reduced to a sufficiently low level, the system will become fully correlated with the stem cell cycle. The noise-induced transition from the daughter cell cycle to the stem cell cycle is shown in the left panel of figure 2.4.2, for a value of the temperature $T$ at which the daughter cell cycle is no longer stable. Monte Carlo simulations of the dynamics for $N = 25,000$ confirm the validity of our analytic solution formulated in terms of the macroscopic dynamic order parameters in (2.13) and (2.14) - right panel of figure 2.4.2. In this figure we do not plot all time-dependent overlaps as in figure 2.4.1 but only the "envelope" of the overlaps defined as the overlaps $m_{\rho(t)}$ and $m_{\rho(t)\mu}$ with the expected cycle state, given by $\rho(t) = 1 + (t$

mod $C$).

The de-differentiation transition takes multiple time steps before a steady state is reached, where the system is in the stem cell cycle attractor. This kind of dynamics is in line with that seen in re-programming experiments, which take multiple generations of the cell cycle before the iPSCs strongly resemble embryonic stem cells [18, 19].

If, however, the noise level is too high the system quickly loses any correlation with all cell cycles - i.e. all the overlaps become zero. To find the range of noise levels over which it is pos-sible to retrieve the stem cell from the daughter cell cycle one can investigate the stability of the solutions of (2.13) and (2.14) as in appendix 2.F. One can also carry out the following numer-ical experiment: the noise level, $T$, was incremented from zero and at each $T$ the equations of motion were solved numerically, the steady state values of the overlaps ($m_{\rho(t)}$ and $m_{\rho(t)\mu}$) which the dynamics converged to were then recorded. These steady state values are plotted against the corresponding noise level in figure 2.4.3. It is clear that above some critical $T$ reprogramming via de-differentiation to the stem cell occurs due to the noise in the system. The value of this critical $T$ will depend on $a^{\rho}$, $a^{\rho\mu}$ and $\gamma^{\rho\mu}$.

The critical value of $T$ for noise-induced reprogramming may also be found numerically, by performing a stability analysis on the equations of motion, the details of which can be found in appendix 2.F. The stability analysis accurately identifies the two transitions in the system - the value of $T$ at which the stem cell pattern is retrieved and the high value of $T$ at which the sys-tem loses all overlap with all cell cycle phases. The numerical values of $T$ at which both of these transitions occur agree perfectly with those observed in figure 2.4.3.

## 2.4.2. DIRECT PERTURBATIONS

The noise induced de-differentiation is different from reprogramming experiments, in which the de-differentiation is due to a direct perturbation using factors that are common to embryonic stem cells (i.e. the Yamanaka factors). Such a directed perturbation can be modelled in our system by

**Figure 2.4.3:** Steady state solutions of (2.13) and (2.14), showing the overlaps with stem and daughter cell cycle stages ($m_\rho$ and $m_{\rho\mu}$ respectively) as a function of noise strength, $T$. The dashed and solid lines correspond to the daughter and stem cell cycle overlaps respectively. At low $T$ the phases of the stem cell cycle that have the same activity result in identical overlaps $m_{\rho\mu}$. However, as $T$ increases the overlaps for each phase of the daughter and stem cell cycles become distinguishable, before reconverging at a critical $T$. Above this critical $T$, the stem cell cycle is retrieved and all $m_\rho(\mathbf{n})$ collapse into a single curve, whilst the $\rho = 1$ and $3$ curves of $m_{\rho\mu}(\mathbf{n})$ recombine. At sufficiently high values of the overlap, the system becomes uncorrelated with all cell cycles. The same parameter values were used as in figure 2.4.1.

introducing an extra contribution $\lambda_i(t)$ to the local field, which pushes the system in the direction of the stem cell cycle, and has the form

$$\lambda_i(t) = k(\eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1})\widetilde{m}_{\bar{\rho}\mu}(t)c_i \,. \tag{2.22}$$

Here $k$ is the strength of the perturbation, $\bar{\rho}$ is the stage of the cycle to which the perturbation is applied, and $c_i$ is a logical variable representing whether or not the perturbation is applied to gene $i$ ($c_i = 1$ with probability $q$, and 0 otherwise).

Since one of the phases of the cell cycle is more similar across different cell types, it is an obvious target for perturbations when attempting to reprogram a cell. The perturbations should be applied just prior to the most similar phase so as to only minimally disrupt the progression of the cell cycle. So choosing $\bar{\rho}$ as the cycle phase prior to the maximally similar one, is expected to be the optimal reprogramming protocol at a given temperature.

In figure 2.4.4 we are carrying out the same numerical experiment as in figure 2.4.3, except the probability, $q$, that a perturbation is applied, is incremented rather than the noise level, $T$. This experiment shows that de-differentiation is possible with a directed perturbation even at low noise levels where the daughter cell cycles are stable. The retrieval of the stem cell cycle is only possible above some critical value of the fraction of perturbed genes, that we call the reprogramming threshold, $q_r$. Because the $\rho = 2$ cell cycle stage is more similar across different cell types, perturbations applied to $\bar{\rho} = 1$ should have a lower $q_r$ value compared with perturbations applied to other phases, i.e. $\bar{\rho} = 2$ or $\bar{\rho} = 3$. This is indeed borne out by the theory.

Increasing the noise level towards the critical $T$ required for noise induced de-differentiation can dramatically change $q_r$, see figure 2.4.5. The critical value $q_r$ has a non-monotonic dependence on the noise level, $T$. This is a direct result of the non-linear nature of the system and the dependence of the perturbation (19) on the dynamical order parameters $\widetilde{m}_{\bar{\rho}\mu}$. As expected the $\bar{\rho} = 1$ perturbations have the lowest $q_r$ values at any given $T$. This is because the $\rho = 2$ phase was made

**Figure 2.4.4:** Steady state stem and daughter cell cycle overlaps, $m_\rho$ and $m_{\rho\mu}$, versus the fraction of genes to which a perturbation of the form (2.22) is applied to, $q$. Here, the perturbations are applied prior to the most similar phase ($\bar{\rho} = 1$). The system was kept at a low noise level of $T = 0.01$, whilst all other parameter values are the same as in figure 2.4.1. The dashed and solid lines correspond to the daughter and stem cell cycle overlaps respectively. At low $q$ the phases of the stem cell cycle that have the same activity result in identical overlaps $m_{\rho\mu}$. However, as $q$ increases, the overlaps for each phase of the daughter and stem cell cycles become distinguishable, before reconverging at critical $q$. Above the critical value of $q$ the stem cell cycle is retrieved and the overlaps for all stem cell phases become identical, whereas overlap with the $\rho = 2$ phase of the daughter cell remains separate from the overlaps for the other phases of the daughter cell cycle.

**Figure 2.4.5:** The fraction, $q_r$, of genes that a perturbation of the form 2.22 is applied to in order to retrieve the stem cell cycle versus noise levels, $T$, increasingly close to that required for the noise induced switching (see figure 2.4.3). The different curves represent different $\bar{\rho}$ protocols for the perturbations. For each protocol a perturbation strength $k = 1$ was used, whilst all other parameter values used are the same as in figure 2.4.1. The relative $q_r$ values at a given $T$ can be explained in terms of the Hamming distances between the states involved in the perturbation ($\boldsymbol{\eta}^{\bar{\rho}\mu}$ and $\boldsymbol{\eta}^{\bar{\rho}+1}$). As this distance increases so does $q_r$. As expected $\bar{\rho} = 1$ is the most efficient perturbation for a single target phase - left panel. For the perturbations applied to two phases ($\bar{\rho} = 1$ and $2$) - right panel - the perturbations were applied to each phase with an equal probability $q_r$.

to exhibit the largest degree of mutual similarity among cell types, due to a decreased activity in this phase. The fact that the $\bar{\rho} = 3$ perturbations have a lower $q_r$ than the $\bar{\rho} = 2$ perturbations follows from the Hamming distance between the state in which the perturbation is applied and the stem cell state targeted by that perturbation, which is smaller for a perturbation applied in the $3\mu$ state than for a perturbation applied in the $2\mu$ state. That is, $d[\boldsymbol{\eta}^{3\mu}, \boldsymbol{\eta}^1] < d[\boldsymbol{\eta}^{2\mu}, \boldsymbol{\eta}^3]$, where the normalized Hamming distance between states is defined as

$$\mathrm{d}\left[\boldsymbol{\eta}^{\bar{\rho}\mu}, \boldsymbol{\eta}^{\bar{\rho}+1}\right] = \frac{1}{N}\sum_{i=1}^{N}\left|\eta_i^{\bar{\rho}+1} - \eta_i^{\bar{\rho}\mu}\right| . \tag{2.23}$$

Hence a higher fraction of genes need to be perturbed to achieve de-differentiation using $\bar{\rho} = 2$ compared with $\bar{\rho} = 3$. The Hamming distances between different cell cycle states are calculated in terms of the activities in Appendix 2.G.

We have also looked at a case where perturbations of the form (2.22) are acting during multiple

stages of the cell cycle in the reprogramming experiments. For example, during the most similar phase and the one prior to it. In this case, the effects of each perturbation are combined and $q_r$ may decrease compared to applying perturbations to a single phase - right panel of figure 2.4.5.

The critical fraction $\mathcal{O}(0.1)$ of genes that need to be perturbed to reprogram a cell, at low $T$, may initially seem much greater than the four Yamanaka factors introduced in the reprogramming experiments (see figure 2.4.4). However, this order of magnitude is actually in line with the experimental results, as can be seen by considering the following argument. From the $20 - 25,000$ genes in the human genome only around $10\%$ are thought to be responsible for synthesising transcription factors [97–99]. This imbalance in numbers requires each transcription factor to interact with many more genes than is needed to synthesise it on average. If we consider the interactions between genes and transcription factors as a bipartite graph (see figure 2.4.6), as we will in later chapter 4, then the TFs have on average an out-degree of $\mathcal{O}(100)$ and the genes have an average in-degree of $\mathcal{O}(10)$ to ensure that there is a conservation in the number of interactions. These numbers agree with the median in- and out-degrees for genes and transcription factors found from a computational analysis of the human gene regulatory network [100]. Assuming each regulatory gene contributes to the synthesis of a single transcription factor, then perturbing 10% of the regulatory genes equates to perturbing roughly $250$ genes, which could be achieved by perturbing just 2-3 transcription factors. Thus the fraction $q_r = \mathcal{O}(0.1)$ of perturbed gene expression levels is *perfectly in line* with the number of transcription factors used to achieve pluripotency in reprogramming experiments.

### 2.4.3. The $G_0$ phase and cell cycle arrest

The local fields defined by (2.10) can be adapted to include other plausible biological features. The main feature discussed here is the resting phase of the cell cycle, known as the $G_0$ phase. Up until this point, all of the cells in the model are set-up such that they transition periodically through a cell cycle. However, most terminally differentiated, or "adult", cells exist in a single

**Figure 2.4.6:** A bipartite graph representing the interactions between genes and transcription factors (TFs). The number of TFs scales as $P = \alpha N$ where $\alpha < 1$. For conservation, the sum of in-degrees of all genes must equal the sum of out-degrees for all TFs. The number of genes that a TF regulates could then be an order of magnitude larger than the number of TFs that interact with a given gene. Thus introducing a small number of TFs could have a significant effect on the gene expression state of the network. Not all of the nodes and connections of the network are shown.

phase known as $G_0$, which is either an extended $G_1$ phase with the cell unable to pass the checkpoint into the $S$ phase or a completely distinct phase from the cell cycle (see figure 2.2.2). This phenomenon can be included in our modelling approach by adapting the form of $J_{ij}$ constructed in section 2.3. A few possibilities are detailed below.

For a two level potency hierarchy in which the stem cell undergoes a cell cycle whereas its progeny are modelled in terms of fixed point attractors representing terminally differentiated cells in the $G_0$ phase, the interactions would take the form

$$J_{ij} = \frac{1}{N} \left\{ \sum_{\rho=1}^{C} \frac{(\eta_i^{\rho+1} - a^{\rho+1})(\eta_j^{\rho} - a^{\rho})}{a^{\rho}(1 - a^{\rho})} + \sum_{\mu=1}^{M} \frac{(\eta_i^{G_0\mu} - a_\mu(\eta_i^{G_0}))(\eta_j^{G_0\mu} - a_\mu(\eta_j^{G_0}))}{a^{G_0\mu}(1 - a^{G_0\mu})} \right\}.$$

$$(2.24)$$

Here, unlike (2.7), the summation over the cell cycle phases applies to only the first term corresponding to the stem cell patterns. In the second term, the contribution to the interactions from each the daughter cell pattern is given in a single phase $\rho = G_0$. For interactions of this type

**Figure 2.4.7:** The overlap for the daughter cell cycle showing activation of a cycle from a fixed $G_0$ phase (left) and, conversely, the arrest of a cell cycle into a fixed $G_0$ phase (right), at a low noise level $T = 0.01$. All other parameter values are the same as in figure 2.4.1.

directed reprogramming would be most efficient when targeting the stem cell cycle phases with the smallest hamming distance, or largest overlap, with the of the initial adult cell state.

An alternative to directly reprogramming a cell fixed in the $G_0$ phase would be to first "restart" that cell's cell cycle, and then going on to use one of the reprogramming mechanisms described in the previous sections. To model this kind of behaviour would require interactions of the form (2.7), but including a dominant contribution for the daughter cell cycles, such as,

$$\epsilon \sum_{\mu=1}^{M} \frac{(\eta_i^{G_0\mu} - a_\mu(\eta_i^{G_0}))(\eta_j^{G_0\mu} - a_\mu(\eta_j^{G_0}))}{a^{G_0\mu}(1 - a^{G_0\mu})} \tag{2.25}$$

where epsilon is a parameter that controls the depth of the fixed point $G_0$ attractors. Then by tuning the value of $\epsilon$ one would be able to restart or arrest the cell cycle, as shown in figure 2.4.7. For daughter cells that are terminally differentiated into a cell cycle phase $\rho = 1$, i.e. $\rho = 1$ corresponds to the $G_0$ phase of the cell cycle, and starting with an initially high value of $\epsilon$, the daughter cell cycle can be activated by decreasing the value of $\epsilon$ linearly in time towards $\epsilon = 0$. Similarly, a cell cycle can be arrested, if the value of $\epsilon$ is increased from an initially negligible contribution. Interactions of this type could be useful for modelling diseases related to cell cycle arrest as well as reprogramming.

## 2.5. SUMMARY

In this chapter, we presented a general (minimal) model for cell reprogramming as transitions between attractors of a dynamical system. The principles of the model are derived from a set of key facts concerning cell chemistry, which suggest that cell types, and their associated cell cycles, can be considered as attractors of the dynamics of interacting gene expression levels. The specific form of gene interactions used to achieve this goal is inspired by combining two strands of neural network modelling: the storage of (limit) cycles and the storage of hierarchically organised attractors.

This chapter is intended to provide a proof of concept of this type of modelling approach. We thus decided to investigate the simplest possible hierarchy of cell types that allows us to test our approach, viz. a two-level hierarchy consisting of the stem cell and a single layer of differentiated cells derived from it. Furthermore, we chose to consider only interactions between pairs of binary gene expression levels.

We have shown that cell reprogramming is possible using either an undirected approach, which consists of increasing the noise level in the dynamics, or an approach that relies on direct perturbations between specific phases of the cell cycle. Two key non-trivial results appear from our model. Firstly, it takes multiple generations of the cell cycle for a progenitor to be reprogrammed to a stem cell, as it transitions through intermediate states which show similarity with both the initial and final state. Also, a finite fraction of gene expression levels need to be perturbed in order to reprogram a cell.

We assume that there are states in the cell cycle where the mutual similarity in gene expression levels between different cell types is large. These stages of the cell cycle are then natural targets for perturbations to induce changes in cell type. The fraction of genes that need to be perturbed in order to reprogram a cell depends on the stage of the cell cycle to which the perturbations are applied, as well as the noise level of the system. At low noise levels, this number was in line with

that required in the Yamanaka reprogramming experiments and was found to decrease (substantially for $\bar{\rho} = 2$ & $3$) with increasing noise levels. The "true" noise level of a cell is difficult to quantify, but our model allows for reprogramming in both low and high noise regimes.

As far as the authors are aware, gene expression levels in different levels of the cell potency hierarchy or in different phases of the cell cycle are still not well characterised. Throughout this chapter we have used a scenario where gene expression levels in differentiated cells are slightly lower than in a stem cell during the same phase of the cell cycle, and we have taken one of the cycle phases to have lower levels of gene expression than the others (thereby increasing mutual similarities of different cells during this cycle phase). We have checked that de-differentiation along the two different routes, noise-induced and via directed perturbations, does not depend on these specific choices, although details, such as critical thresholds for reprogramming, do change as scenarios are modified.

At the time of writing, I am aware of only a single other study that models the cell cycle as configurations of gene expression patterns using a cyclic Hopfield-like model [82]. However, said work does not investigate transitions between different cell types and is not concerned with reprogramming dynamics.

There are some limitations to the modelling approach presented in this chapter. Firstly, we consider only pairwise interactions between regulatory gene expression levels. Including higher-order interactions in NN models typically stabilise the dynamics and allow for storage of a greater number of attractors. Incorporating higher order interactions between genes would be biologically reasonable, since proteins expressed from multiple genes can form transcription factors complexes. Also, genes can often require proteins binding to promoter sites and enhancer regions before they are activated. Using discrete time dynamics excludes the possibility of variability in gene expression levels in a given cell cycle phase. Therefore any in-cycle dynamics is missed, such as any cell signalling cascades. Finally, we only use rough estimates for the average gene expression levels in numerical experiments and simulations.

One possibility for extending the model presented in this chapter would be to relax the choice of independent cell cycle states. Correlated cell cycle phases can be incorporated into the two-level hierarchy by changing the way in which the cell cycles are constructed in the model. One could achieve this using a three-level hierarchy to store all cell cycles, whilst maintaining the feature that all descendants are a single differentiation from the stem cell cycle. In this situation the root of the hierarchy would be a template of the stem cell expression levels, the second level would then be constructed from this and represent each stage of the stem cell cycle. The newly included third level would consist of each daughter cell type branching off from the corresponding stem cell cycle phase. Such a set up would then be analogous to a two-level hierarchy for each cell cycle phase with correlations in the gene expression levels along the cell cycles of each cell type.

In the next chapter, the parameter choices used in this chapter are compared against real data, in order to justify their choice. However, since there is still much to be learnt about cell reprogramming and the decisions of cell fates in developmental biology, the model has been presented here under the belief that it captures aspects of cellular reprogramming both qualitatively and quantitatively, and it therefore represents a possible stepping-stone for developing future theoretical models and as a motivation for experimental work to be directed towards creating a catalogue of gene expression profiles around the cell cycle for different cell types.

# APPENDICES

APPENDIX 2.A    DIFFERENTIATION TRANSITION MATRICES

To derive the daughter cell expression levels from the stem cell's a transition matrix, $\mathbf{W}$, was used. When $\mathbf{W}$ is applied to the probability distribution of the a stem cell phase the result corresponds to the distribution of the daughter state, i.e. $\mathbf{W} \begin{bmatrix} a^\rho \\ 1 - a^\rho \end{bmatrix} = \begin{bmatrix} a^{\rho\mu} \\ 1 - a^{\rho\mu} \end{bmatrix}$. If we define the probability that a gene is switched off during differentiation from the stem to daughter cell as $\gamma^{\rho\mu}$, but the activities remain equal in the same cell cycle phases after differentiation $\left( a^{\rho\mu} = a^\rho \right)$, then we can construct $\mathbf{W}$ as,

$$W(\eta^{\rho\mu}|\eta^\rho) = \begin{bmatrix} W(1|1) & W(1|0) \\ W(0|1) & W(0|0) \end{bmatrix} = \begin{bmatrix} (1 - \gamma^{\rho\mu}) & \dfrac{\gamma^{\rho\mu} a^\rho}{1 - a^\rho} \\ \gamma^{\rho\mu} & 1 - \left( \dfrac{\gamma^{\rho\mu} a^\rho}{1 - a^\rho} \right) \end{bmatrix} . \tag{2.26}$$

Applying this transition matrix to the distribution of stem cell gene expression levels $\begin{bmatrix} a^\rho \\ 1 - a^\rho \end{bmatrix}$ returns the same activities for the daughter cell. Thus, for different activities in the same cell cycle

phases of the daughter and stem cells, one can easily adapt $\mathbf{W}$ to be of the form,

$$W(\eta^{\rho\mu}|\eta^{\rho}) = \begin{bmatrix} (1 - \gamma^{\rho\mu})\dfrac{a^{\rho\mu}}{a^{\rho}} & \dfrac{\gamma^{\rho\mu}a^{\rho\mu}}{1 - a^{\rho}} \\ 1 - (1 - \gamma^{\rho\mu})\dfrac{a^{\rho\mu}}{a^{\rho}} & 1 - \left(\dfrac{\gamma^{\rho\mu}a^{\rho\mu}}{1 - a^{\rho}}\right) \end{bmatrix}. \tag{2.27}$$

To have (2.27) defined as a transition matrix, its columns must sum to one - as they clearly do - with each element $W(\eta^{\rho\mu}|\eta^{\rho}) \in [0, 1]$. Then, since $a^{\rho}$ and $a^{\rho\mu} \in [0, 1]$ by definition, we must obey the following constraint on $\gamma^{\rho\mu}$ for itself and (2.27) to be correctly defined as a probability,

$$\max\left[0, \frac{a^{\rho\mu} - a^{\rho}}{a^{\rho\mu}}\right] \leq \gamma^{\rho\mu} \leq \min\left[1, \frac{1 - a^{\rho}}{a^{\rho\mu}}\right]. \tag{2.28}$$

## APPENDIX 2.B  DERIVATION OF THE EQUATIONS OF MOTION

In each time step the state of the system is updated based on the local fields at each site. That is, the expression level of gene $i$, at time $t$, depends on the value of the field at time $t - 1$, i.e.

$$n_i(t + 1) = \Theta\left[h_i(\mathbf{m}(t)) - T\xi_i(t)\right], \tag{2.29}$$

where $\Theta(x)$ is the Heaviside step function ($\Theta(x) = 0$ for $x \leq 0$ and $\Theta(x) = 1$ for $x > 0$), and $\xi_i(t)$ is thermal noise at the site $i$ (with $\mathbb{P}\left[\xi_i(t) < z\right] = \frac{1}{2}\left[1 + \tanh(\beta z/2)\right]$).

The expected value of site $i$ can be obtained by averaging (2.29),

$$\langle n_i(t + 1)\rangle = \mathbb{P}\left[\Theta\left[h_i(\mathbf{m}(t)) - \xi_i(t)\right] > 0\right] \tag{2.30}$$

$$= \mathbb{P}\left[\xi_i(t) < h_i(\mathbf{m}(t))\right] \tag{2.31}$$

$$= \frac{1}{2}\left[1 + \tanh\left(\frac{\beta h_i(\mathbf{m}(t))}{2}\right)\right]. \tag{2.32}$$

Then using the definitions of the dynamic order parameters (2.8) and (2.9), the following expressions can be obtained, when the $N \to \infty$ limit is taken, by making use of the law of large

numbers:

$$\widetilde{m}_\rho(t+1) = \frac{1}{2}\left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1-a^\rho)}\tanh\left(\frac{\beta h(t)}{2}\right)\right\rangle_{\eta^\rho,\eta^{\rho\mu}},\tag{2.33}$$

$$\widetilde{m}_{\rho\mu}(t+1) = \frac{1}{2}\left\langle \frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1-a^{\rho\mu})}\tanh\left(\frac{\beta h(t)}{2}\right)\right\rangle_{\eta^\rho,\eta^{\rho\mu}},\tag{2.34}$$

where $\langle\ldots\rangle_{\eta^\rho,\eta^{\rho\mu}}$ is shorthand for an average over the statistics of (correlated) stem and daughter cell expression levels throughout their cycles, i.e.

$$\prod_{\rho=1}^{C}\left[\sum_{\eta^\rho\in\{0,1\}} p(\eta^\rho)\prod_{\mu=1}^{M}\left[\sum_{\eta^{\rho\mu}\in\{0,1\}} W(\eta^{\rho\mu}|\eta^\rho)(\ldots)\right]\right].$$

A more rigorous calculation to determine these equations of motion can be done by following the reasoning in Coolen et al [101], and is shown in Appendix 2.C.

## APPENDIX 2.C   MACROSCOPIC DYNAMICS FOR THE ORDER PARAMETERS

Here we derive the equations of motion for the dynamic order parameters from the macroscopic dynamics. This is a more rigorous calculation that can be done to obtain the same results as in Appendix 2.B and is thus only shown for the interested reader. Starting from the dynamics (2.3), we know that a gene will be expressed in the next time step if the effects of all the interactions is greater than the biological noise in the system, i.e. $\mathbb{P}[n_i(t+1)] = \mathbb{P}[-\xi_i < \frac{h_i(t)}{T}]$. Similarly, a gene is not expressed if its interactions with all other genes is less than the biological noise in the system. For a symmetric noise distribution $\mathbb{P}(\xi) = \mathbb{P}(-\xi)\ \forall i$ we can write the probabilities that a gene is expressed or not as

$$\mathbb{P}[n_i(t+1)=1] = \mathbb{P}\left[\xi_i < \frac{h_i(t)}{T}\right] = \int_{-\infty}^{\frac{h_i(t)}{T}} d\xi P(\xi),\tag{2.35}$$

$$\mathbb{P}[n_i(t+1) = 0] = \mathbb{P}\left[\xi_i \leq \frac{-h_i(t)}{T}\right] = \int_{-\infty}^{\frac{-h_i(t)}{T}} d\xi P(\xi)\,, \tag{2.36}$$

which combine to give the probability

$$\mathbb{P}[n_i(t+1)] = \int_{-\infty}^{\frac{(2n_i(t+1)-1)h_i(t)}{T}} d\xi P(\xi) = \frac{1}{2} + \int_0^{\frac{(2n_i(t+1)-1)h_i(t)}{T}} d\xi P(\xi)$$
$$= g\left(\frac{(2n_i(t+1)-1)\,h_i(t)}{T}\right)\,. \tag{2.37}$$

There are many choices one could make for the function $g$, however we choose the function $g(x) = \frac{1}{2}\left[1 + \tanh(\frac{x}{2})\right]$, which has similar behaviour to a Gaussian distribution, but is much easier to work with. Assuming the system undergoes parallel updates, with each gene treated independently, the probability of the future state of the system is given by

$$\mathbb{P}[\mathbf{n}(t+1)] = \prod_{i=1}^N \frac{1}{2}\left[1 + (2n_i(t+1)-1)\tanh\left(\frac{\beta h_i(t)}{2}\right)\right]\,, \tag{2.38}$$

where we have made use of $\tanh[\pm x] = \pm \tanh[x]$ because $n_i \in \{0, 1\}$, and the shorthand of $\beta = T^{-1}$ and $h_i(\mathbf{n}(t)) = h_i(t)$.

We would like to write the stochastic dynamics in terms of the evolving microscopic probabilities of state. First we introduce the notation $p_t(n) = \mathbb{P}[\mathbf{n}(t) = \mathbf{n}]$ and rewrite the state probability as,

$$p_{t+1}(\mathbf{n}) = \prod_i^N \frac{\exp\left[\frac{(2n_i(t)-1)\beta h_i(t)}{2}\right]}{2\cosh\left(\frac{\beta h_i(t)}{2}\right)}\,, \tag{2.39}$$

where the hyperbolic tangent has been written in its exponential form, $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$, and we have made use of the identity $\frac{1}{2}\left[1 \pm \tanh(x)\right] = \frac{e^{\pm x}}{2\cosh(x)}$.

If we do not know the exact state that the system is in but we know the probability distribution

over all possible configurations, then we can rewrite the dynamics as a Markov process.

$$p_{t+1}(\mathbf{n}) = \sum_{\mathbf{n}'} W(\mathbf{n}, \mathbf{n}') \, p_t(\mathbf{n}') \tag{2.40}$$

We can do this because the future state of the system only depends on the current configuration $\mathbf{n}'$. The Markov transition probability from state $\mathbf{n} \to \mathbf{n}'$ is then given by

$$W(\mathbf{n}, \mathbf{n}') = \prod_i^N \frac{\exp\left[\frac{(2n_i(t)-1)\beta h_i(\mathbf{n}')(t)}{2}\right]}{2\cosh\left(\frac{\beta h_i(\mathbf{n}')(t)}{2}\right)}. \tag{2.41}$$

We know from the definition of the local fields (2.10) that the $h_i(t)$ depends on the state of the system only through the dynamical order parameters. The order parameters are a natural candidate to describe the systems behaviour on a macroscopic level. The probability of the system having values $\widetilde{m}_\rho(\mathbf{n})$ and $\widetilde{m}_{\rho\mu}(\mathbf{n})$ at a time $t$ is then given by

$$P_t(\widetilde{m}_\rho(\mathbf{n}), \widetilde{m}_{\rho\mu}(\mathbf{n})) = \sum_{\mathbf{n}} p_t(n)\delta\left(\widetilde{m}_\rho - \widetilde{m}_\rho(\mathbf{n})\right)\delta\left(\widetilde{m}_{\rho\mu} - \widetilde{m}_{\rho\mu}(\mathbf{n})\right) \tag{2.42}$$

Thus, making use of our Markov property we can rewrite the dynamics as

$$P_{t+1}(\widetilde{m}_\rho(\mathbf{n}), \widetilde{m}_{\rho\mu}(\mathbf{n})) = \int d\widetilde{\mathbf{m}}' \widetilde{W}_t(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}')P_t(\widetilde{\mathbf{m}}') \tag{2.43}$$

with the kernel $\widetilde{W}_t(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}')$ defined as,

$$\frac{\sum_{\mathbf{n},\mathbf{n}'} \delta\left(\widetilde{m}_\rho - \widetilde{m}_\rho(\mathbf{n})\right)\delta\left(\widetilde{m}'_\rho - \widetilde{m}'_\rho(\mathbf{n}')\right)\delta\left(\widetilde{m}_{\rho\mu} - \widetilde{m}_{\rho\mu}(\mathbf{n})\right)\delta\left(\widetilde{m}'_{\rho\mu} - \widetilde{m}'_{\rho\mu}(\mathbf{n}')\right) W(\mathbf{n}, \mathbf{n}')p_t(\mathbf{n}')}{\sum_{\mathbf{n}'} \delta\left(\widetilde{m}'_\rho - \widetilde{m}'_\rho(\mathbf{n}')\right)\delta\left(\widetilde{m}'_{\rho\mu} - \widetilde{m}'_{\rho\mu}(\mathbf{n}')\right) p_t(\mathbf{n}')}.$$

$$\tag{2.44}$$

Using the definition of the Markov transition probability (given above) and the local fields we

can rewrite this kernel in the form,

$$\widetilde{W}(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}') = \sum_{\mathbf{n}} \delta\left(\widetilde{\mathbf{m}} - \widetilde{\mathbf{m}}(\mathbf{n})\right) \left[ \prod_{i=1}^{N} \exp\left\{ -\ln 2\cosh\left(\frac{\beta h_i(\mathbf{n}')}{2}\right) e^{\frac{\beta(2n_i - 1)h_i(\mathbf{n}')}{2}} \right\} \right]$$

(2.45)

where $\widetilde{\mathbf{m}}$ contains all dynamic order parameters, and the average over $p_t(\mathbf{n}')$ vanishes. This is because the local fields depend on $\mathbf{n}'$ only through $\widetilde{\mathbf{m}}'_\rho$ and $\widetilde{\mathbf{m}}'_{\rho\mu}$, which have been fixed by the $\delta$-functions. Hence, the time dependence of the kernel can also be dropped and it can be written as the average,

$$\widetilde{W}(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}') = \left\langle \delta(\widetilde{\mathbf{m}} - \widetilde{\mathbf{m}}(\mathbf{n})) e^{-\sum_i \ln\cosh\frac{\beta h_i(\mathbf{n}')}{2}} e^{\frac{\beta}{2}\sum_i (2n_i - 1)h_i(\mathbf{n}')} \right\rangle_{\mathbf{n}},$$

(2.46)

where $\langle \ldots \rangle_{\mathbf{n}} = \frac{1}{2^N} \sum_{\mathbf{n}}(\ldots)$. The components of $\widetilde{W}(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}')$ can be rewritten by using,

$$\delta\left(\widetilde{\mathbf{m}}_\rho - \widetilde{\mathbf{m}}_\rho(\mathbf{n})\right) = \delta\left(\beta N \widetilde{m}_\rho - \beta \sum_i \frac{\eta_i^\rho - a^\rho}{a^\rho(1 - a^\rho)}\right),$$

$$\delta\left(\widetilde{\mathbf{m}}_{\rho\mu} - \widetilde{\mathbf{m}}_{\rho\mu}\right) = \delta\left(\beta N \widetilde{m}_{\rho\mu} - \beta \sum_i \frac{\eta_i^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})}\right),$$

and the integral representation of the delta function, in the form

$$\widetilde{W}(\widetilde{\mathbf{m}}_\rho, \widetilde{\mathbf{m}}'_\rho) = \frac{\beta N}{2\pi} \int dk \exp\left[ N\left( ik\beta m_\rho - \left\langle \frac{\beta h(\widetilde{\mathbf{m}}')}{2} \right\rangle - \left\langle \ln\cosh\frac{\beta h(\widetilde{\mathbf{m}}')}{2} \right\rangle \right.\right.$$

$$\left.\left. + \left\langle \ln\left( 1 + e^{-ik\beta \frac{\eta^\rho - a^\rho}{a^\rho(1 - a^\rho)} + \beta h(\widetilde{\mathbf{m}}')} \right) \right\rangle \right) \right],$$

(2.47)

and

$$
\widetilde{W}(\widetilde{\mathbf{m}}_{\rho\mu}, \widetilde{\mathbf{m}}'_{\rho\mu}) = \frac{\beta N}{2\pi} \int dk \exp\left[ N\left( ik\beta m_{\rho\mu} - \left\langle \frac{\beta h(\widetilde{\mathbf{m}}')}{2} \right\rangle - \left\langle \ln\cosh \frac{\beta h(\widetilde{\mathbf{m}}')}{2} \right\rangle \right.\right.
$$
$$
\left.\left. + \left\langle \ln\left( 1 + e^{-ik\beta \frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})} + \beta h(\widetilde{\mathbf{m}}')} \right) \right\rangle \right) \right].
$$

$$(2.48)$$

Because we are in the limit of large $N$, these integrals are in a form that can be evaluated using the saddle point method. Evaluating these integrals, one finds that at the saddle point the dynamic order parameters are

$$
\widetilde{m}_\rho(\mathbf{n}(t+1)) = \frac{1}{2}\left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1 - a^\rho)} \tanh\left( \frac{\beta h(\widetilde{\mathbf{m}}'(t))}{2} \right) \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} \tag{2.49}
$$

$$
\widetilde{m}_{\rho\mu}(\mathbf{n}(t+1)) = \frac{1}{2}\left\langle \frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})} \tanh\left( \frac{\beta h(\widetilde{\mathbf{m}}'(t))}{2} \right) \right\rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}} \tag{2.50}
$$

where $\langle \ldots \rangle_{\boldsymbol{\eta}^\rho, \boldsymbol{\eta}^{\rho\mu}}$ is shorthand for an average over the statistics of stem and daughter cell expression levels throughout their cycles, i.e.

$$
\prod_{\rho=1}^{C}\left[ \sum_{\eta^\rho \in \{0,1\}} p(\eta^\rho) \prod_{\mu=1}^{M}\left[ \sum_{\eta^{\rho\mu} \in \{0,1\}} W(\eta^{\rho\mu}|\eta^\rho)(\ldots) \right] \right].
$$

Thus, in the large $N$ limit, the kernel $\widetilde{W}(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}')$ becomes $\delta\left( \widetilde{\mathbf{m}} - \widetilde{\mathbf{m}}^* \right)$ where $\widetilde{\mathbf{m}}^*$ is the saddle point solution with the components given above. The evolution of the distribution

$$
P_{t+1}(\widetilde{m}_\rho, \widetilde{m}_{\rho\mu}) = \int d\widetilde{\mathbf{m}}' \widetilde{W}(\widetilde{\mathbf{m}}, \widetilde{\mathbf{m}}') P_t(\widetilde{\mathbf{m}}')
$$

is therefore deterministic, and the dynamic order parameters will evolve according to the sad-

dle point solutions (2.49) & (2.50). Hence, if we know the initial values of the dynamic order parameters we can calculate their values at any future time.

## APPENDIX 2.D    SIGNAL-TO-NOISE ANALYSIS

Starting from the decomposition of the local fields (2.15) with the condensed pattern chosen to be $\eta^{\bar{\rho}}$, in the large system limit the signal contribution to the local fields is

$$
S = (\eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1}) \left\langle \frac{\eta^{\bar{\rho}} - a^{\bar{\rho}}}{a^{\bar{\rho}}(1 - a^{\bar{\rho}})} \eta^{\bar{\rho}} \right\rangle + \sum_\mu (\eta_i^{\bar{\rho}+1,\mu} - a_\mu(\eta_i^{\bar{\rho}+1})) \left\langle \frac{\eta^{\bar{\rho}\mu} - a_\mu(\eta^{\bar{\rho}})}{a^{\bar{\rho}\mu}(1 - a^{\bar{\rho}\mu})} \eta^{\bar{\rho}} \right\rangle ,
$$
(2.51)

where the law of large numbers has been used to replace the sum over genes. The first average in $S$ can be evaluated to 1 simply by using $\langle (\eta^{\bar{\rho}})^2 \rangle = \langle \eta^{\bar{\rho}} \rangle$ for $\eta^{\bar{\rho}} \in \{0, 1\}$. Next, the second average term evaluates to zero for all $\mu$, this can be seen by first averaging $\eta^{\bar{\rho}\mu}$ over the conditional probability $W(\eta^{\bar{\rho}\mu}|\eta^{\bar{\rho}})$ to get $a_\mu(\eta^{\bar{\rho}})$ by definition. Thus, for $N \to \infty$ the signal contribution to the local field $h_i$ is simply,

$$
S = \eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1} .
$$
(2.52)

If $\eta^{\bar{\rho}+1} = 1$ the signal is positive and will result in $n_i(t+1) = 1$ according to the gene expression level dynamics (2.3). Contrarily, when $\eta^{\bar{\rho}+1} = 0$ the signal is negative and acts to silence the gene expression of gene $i$ at the next time step according to (2.3). Thus, the signal acts to progress the gene expression levels of the cell cycle and will do successfully provided that the noise from all other contributions to the local field do not overcome it.

The noise contribution to the local fields $R$ is given by,

$$
R = \frac{1}{N} \sum_{\rho(\neq\bar{\rho})} \left\{ (\eta_i^{\rho+1} - a^{\rho+1}) \sum_{j(\neq i)} \frac{\eta_j^{\rho} - a^{\rho}}{a^{\rho}(1 - a^{\rho})} \eta_j^{\bar{\rho}} \right.
$$
$$
\left. + \sum_{\mu} (\eta_i^{\rho+1,\mu} - a_{\mu}(\eta_i^{\rho+1})) \sum_{j(\neq i)} \frac{\eta_j^{\rho\mu} - a_{\mu}(\eta_j^{\rho})}{a^{\rho\mu}(1 - a^{\rho\mu})} \eta_j^{\bar{\rho}} \right\},
$$

$$(2.53)$$

which is a sum of independent random variables each of which are zero mean due to the independence of gene expression levels in different stages of the cell cycle. Therefore, $R$ itself is a zero-mean random variable ($\langle R \rangle = 0$) that will have a Gaussian distribution according to the central limit theorem. Its variance is then roughly given by $\langle R^2 \rangle$. Hence, in order for the noise contributions not to corrupt the progression of the cell cycle, due to $S$, the variance of $R$ should be small. Writing $R$ out in full and then squaring it one obtains

$$
R^2 = N^{-2} \sum_{\substack{\rho(\neq\bar{\rho}),\\\rho'(\neq\bar{\rho})}} (\eta_i^{\rho+1} - a^{\rho+1})(\eta_i^{\rho'+1} - a^{\rho'+1}) \sum_{\substack{j(\neq i),\\k(\neq i)}} \frac{(\eta_j^{\rho} - a^{\rho})(\eta_k^{\rho'} - a^{\rho'})}{a^{\rho}a^{\rho'}(1 - a^{\rho})(1 - a^{\rho'})} \eta_j^{\bar{\rho}}\eta_k^{\bar{\rho}} +
$$
$$
N^{-2} \sum_{\substack{\rho(\neq\bar{\rho}),\\\rho'(\neq\bar{\rho}),\\\mu,\nu,\\j(\neq i),\\k(\neq i)}} (\eta_i^{\rho+1,\mu} - a_{\mu}(\eta_i^{\rho+1}))(\eta_i^{\rho'+1,\nu} - a_{\nu}(\eta_i^{\rho'+1})) \frac{(\eta_j^{\rho\mu} - a_{\mu}(\eta_j^{\rho}))(\eta_k^{\rho'\nu} - a_{\nu}(\eta_k^{\rho'}))}{a^{\rho\mu}a^{\rho'\nu}(1 - a^{\rho\mu})(1 - a^{\rho'\nu})} \eta_j^{\bar{\rho}}\eta_k^{\bar{\rho}}
$$
$$
+ 2N^{-1} \sum_{\substack{\rho(\neq\bar{\rho}),\\\mu}} (\eta_i^{\rho+1} - a^{\rho+1})(\eta_i^{\rho+1,\mu} - a_{\mu}(\eta_i^{\rho+1})) \sum_{\substack{j(\neq i),\\k(\neq i)}} \frac{(\eta_j^{\rho} - a^{\rho})(\eta_k^{\rho\mu} - a_{\mu}(\eta_k^{\rho}))}{a^{\rho}a^{\rho\mu}(1 - a^{\rho})(1 - a^{\rho\mu})} \eta_j^{\bar{\rho}}\eta_k^{\bar{\rho}},
$$

$$(2.54)$$

which should be averaged to find the variance of the noise contribution to the local fields $\langle R^2 \rangle$. When averaging $R^2$, the cross term (bottom line of $R^2$ above) evaluates to zero because of independent gene expression levels in different cell cycle stages, and by first performing the conditional averaging over the daughter cell cycle gene expression levels. In fact, because of the inde-

pendent gene expression levels, the only contribution to $\langle R^2 \rangle$ from the first two lines comes from the scenario in which $\rho = \rho'$, $\mu = \nu$ and $j = k$, giving,

$$
\begin{aligned}
\langle R^2 \rangle =& N^{-2} \left\langle \sum_{\rho(\neq \bar\rho)} (\eta_i^{\rho+1} - a^{\rho+1})^2 \sum_{j(\neq i)} \frac{(\eta_j^\rho - a^\rho)^2}{(a^\rho(1-a^\rho))^2} (\eta_j^{\bar\rho})^2 \right\rangle \\
&+ N^{-2} \left\langle \sum_{\substack{\rho(\neq \bar\rho),\\ \mu}} (\eta_i^{\rho+1,\mu} - a_\mu(\eta_i^{\rho+1}))^2 \sum_{j(\neq i)} \frac{(\eta_j^{\rho\mu} - a_\mu(\eta_j^\rho))^2}{(a^{\rho\mu}(1-a^{\rho\mu}))^2} (\eta_j^{\bar\rho})^2 \right\rangle .
\end{aligned}
\tag{2.55}
$$

Then using the following averages,

$$
\langle (\eta^\rho - a^\rho)^2 \rangle = a^\rho(1 - a^\rho),
\tag{2.56}
$$

$$
\langle (\eta^\rho - a^\rho)^2 (\eta^\rho)^2 \rangle = a^\rho(1 - a^\rho)(1 - a^\rho),
\tag{2.57}
$$

$$
\begin{aligned}
\langle (\eta^{\rho\mu} - a_\mu(\eta^\rho))^2 \rangle &= \sum_{\eta^\rho, \eta^{\rho\mu}} p(\eta^\rho) W(\eta^{\rho\mu} | \eta^\rho) \left[ (\eta^{\rho\mu})^2 + a_\mu(\eta^\rho)^2 - 2 a_\mu(\eta^\rho) \eta^{\rho\mu} \right] \\
&= \sum_{\eta^\rho} p(\eta^\rho) \left[ a_\mu(\eta^\rho) + a_\mu(\eta^\rho)^2 - 2 a_\mu(\eta^\rho)^2 \right] \\
&= \sum_{\eta^\rho} p(\eta^\rho) (a_\mu(\eta^\rho) - a_\mu(\eta^\rho)^2) \\
&= a^{\rho\mu} - (1 - \gamma^{\rho\mu})^2 \frac{(a^{\rho\mu})^2}{a^\rho} - \frac{(\gamma^{\rho\mu} a^{\rho\mu})^2}{1 - a^\rho},
\end{aligned}
\tag{2.58}
$$

and

$$
\begin{aligned}
\langle (\eta^{\rho\mu} - a_\mu(\eta^\rho))^2 (\eta^\rho)^2 \rangle &= \sum_{\eta^\rho} (\eta^\rho)^2 \sum_{\eta^{\rho\mu}} (\eta^{\rho\mu} - a_\mu(\eta^\rho))^2 W(\eta^{\rho\mu}|\eta^\rho) \\
&= \sum_{\eta^\rho} \eta^\rho p(\eta^\rho) \sum_{\eta^{\rho\mu}} \left[ (\eta^\rho)^2 + a_\mu(\eta^\rho)^2 - 2\eta^{\rho\mu} a_\mu(\eta^\rho) \right] W(\eta^{\rho\mu}|\eta^\rho) \\
&= \sum_{\eta^\rho} \eta^\rho p(\eta^\rho) \left[ a_\mu(\eta^\rho) + a_\mu(\eta^\rho)^2 - 2a_\mu(\eta^\rho)^2 \right] \\
&= \sum_{\eta^\rho} \eta^\rho a_\mu(\eta^\rho)(1 - a_\mu(\eta^\rho)) p(\eta^\rho) \\
&= a^{\rho\mu}(1 - \gamma^{\rho\mu}) \left( 1 - \frac{(1 - \gamma^{\rho\mu}) a^{\rho\mu}}{a^\rho} \right) ,
\end{aligned}
\tag{2.59}
$$

one obtains

$$
\begin{aligned}
\langle R^2 \rangle =& \frac{(N-1)(C-1)}{N^2 a^\rho (1 - a^\rho)} \left[ a^{\bar{\rho}}(1 - a^{\bar{\rho}})(1 - a^{\bar{\rho}}) + a^{\bar{\rho}} a^{\rho+1}(1 - a^{\rho+1}) \right] \\
&+ \frac{M(N-1)(C-1)}{N^2 (a^{\rho\mu}(1 - a^{\rho\mu}))^2} \left[ a^{\rho\mu} - (1 - \gamma^{\rho\mu})^2 \frac{(a^{\rho\mu})^2}{a^\rho} - \frac{(\gamma^{\rho\mu} a^{\rho\mu})^2}{1 - a^\rho} \right] \\
&\times \left\{ a^{\bar{\rho}\mu}(1 - \gamma^{\bar{\rho}\mu}) \left( 1 - (1 - \gamma^{\bar{\rho}\mu}) \frac{a^{\bar{\rho}\mu}}{a^{\bar{\rho}}} \right) \right. \\
&\left. + a^{\bar{\rho}} \left( a^{\rho+1,\mu} - (1 - \gamma^{\rho+1,\mu})^2 \frac{(a^{\rho+1,\mu})^2}{a^{\rho+1}} - \frac{(\gamma^{\rho+1,\mu} a^{\rho+1,\mu})^2}{1 - a^{\rho+1}} \right) \right\} .
\end{aligned}
\tag{2.60}
$$

This is of the form

$$
\langle R^2 \rangle = \frac{(N-1)(C-1)}{N^2} A_1 + \frac{(N-1)(C-1)M}{N^2} A_2 ,
\tag{2.61}
$$

where $C$ is the number of cell cycle phases, $M$ is the number of daughter cells in the potency hierarchy and $A_1$ & $A_2$ are constants. For multicellular organisms, typically $M \gg C$ and we are in the limit of large $N$ where $(N-1) \approx N$. Thus, the dominant term in the variance of the noise

contribution to the local fields is of the form

$$\langle R^2 \rangle \sim \frac{M(C-1)}{N} , \tag{2.62}$$

which should be small to allow the cell cycle to proceed regularly. Therefore, to avoid a spin glass regime we require $N \gg M(C-1)$. This is in line with what is observed in multicellular organisms, e.g. for humans $N \sim 25,000$, $M \sim 300$ and $C = 5$.

Similarly one can perform the same analysis with a daughter cell cycle configuration as the condensed pattern $\boldsymbol{\eta}^{\bar{\rho}\bar{\mu}}$. In this case, the signal is

$$\begin{aligned}
S = {}& (\eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1}) \left\langle \frac{\eta^{\bar{\rho}} - a^{\bar{\rho}}}{a^{\bar{\rho}}(1 - a^{\bar{\rho}})} \eta^{\bar{\rho}\bar{\mu}} \right\rangle + (\eta_i^{\bar{\rho}+1,\bar{\mu}} - a_{\bar{\mu}}(\eta_i^{\bar{\rho}+1})) \left\langle \frac{\eta^{\bar{\rho}\bar{\mu}} - a_{\mu}(\eta^{\bar{\rho}})}{a^{\bar{\rho}\bar{\mu}}(1 - a^{\bar{\rho}\bar{\mu}})} \eta^{\bar{\rho}\bar{\mu}} \right\rangle \\
= {}& \frac{a^{\bar{\rho}\bar{\mu}}(1 - \gamma^{\bar{\rho}\bar{\mu}} - a^{\bar{\rho}})}{a^{\bar{\rho}}(1 - a^{\bar{\rho}})} (\eta_i^{\bar{\rho}+1} - a^{\bar{\rho}+1}) \\
& + (\eta_i^{\bar{\rho}+1,\bar{\mu}} - a_{\bar{\mu}}(\eta_i^{\bar{\rho}+1})) \left[ (1 - \gamma^{\bar{\rho}\bar{\mu}}) a^{\bar{\rho}\bar{\mu}} \left( 1 - (1 - \gamma^{\bar{\rho}\bar{\mu}}) \frac{a^{\bar{\rho}\bar{\mu}}}{a^{\bar{\rho}}} \right) \right. \\
& \left. + \gamma^{\bar{\rho}\bar{\mu}} a^{\bar{\rho}\bar{\mu}} \left( 1 - \frac{\gamma^{\bar{\rho}\bar{\mu}} a^{\bar{\rho}\bar{\mu}}}{1 - a^{\bar{\rho}}} \right) \right] , \tag{2.63}
\end{aligned}$$

and the sign of $S$ will depend on the choice of the probabilities $\gamma$ as well as the distributions of the gene expression levels. If the probabilities $\gamma$ are chosen appropriately the cell cycle of the daughter cell $\bar{\mu}$ will progress as expected. Applying the same reasoning to the noise contribution to the local field above, but with the gene expression configuration $\boldsymbol{\eta}^{\bar{\rho}\bar{\mu}}$ as the condensed pattern, one finds that $\langle R \rangle$ is a zero-mean random variable. The variance of $R$, is then given by

$$\begin{aligned}
\langle R^2 \rangle = {}& N^{-2} \left\langle \sum_{\rho(\neq\bar{\rho})} (\eta_i^{\rho+1} - a^{\rho+1})^2 \sum_{j(\neq i)} \frac{(\eta_j^{\rho} - a^{\rho})^2}{(a^{\rho}(1 - a^{\rho}))^2} (\eta_j^{\bar{\rho}\bar{\mu}})^2 \right\rangle \\
& + N^{-2} \left\langle \sum_{\substack{\rho(\neq\bar{\rho}), \\ \mu(\neq\bar{\mu})}} (\eta_i^{\rho+1,\mu} - a_{\mu}(\eta_i^{\rho+1}))^2 \sum_{j(\neq i)} \frac{(\eta_j^{\rho\mu} - a_{\mu}(\eta_j^{\rho}))^2}{(a^{\rho\mu}(1 - a^{\rho\mu}))^2} (\eta_j^{\bar{\rho}\bar{\mu}})^2 \right\rangle . \tag{2.64}
\end{aligned}$$

This can be evaluated using the averages,

$$\langle(\eta^\rho - a^\rho)^2(\eta^{\rho\mu})^2\rangle = \sum_{\eta^\rho}(\eta^\rho - a^\rho)^2 p(\eta^\rho)\sum_{\eta^{\rho\mu}}(\eta^{\rho\mu})^2 W(\eta^{\rho\mu}|\eta^\rho)$$

$$= \sum_{\eta^\rho}(\eta^\rho - a^\rho)^2 a_\mu(\eta^\rho)p(\eta^\rho)$$

$$= (1 - a^\rho)^2(1 - \gamma^{\rho\mu})a^{\rho\mu} + (a^\rho)^2\gamma^{\rho\mu}a^{\rho\mu}, \qquad (2.65)$$

and

$$\langle(\eta^{\rho\mu} - a_\mu(\eta^\rho))^2(\eta^{\rho\mu})^2\rangle = \sum_{\eta^\rho,\eta^{\rho\mu}} p(\eta^\rho)W(\eta^{\rho\mu}|\eta^\rho)\left[(\eta^{\rho\mu})^4 + (\eta^{\rho\mu})^2 a_\mu(\eta^\rho)^2 - 2(\eta^{\rho\mu})^3 a_\mu(\eta^\rho)\right]$$

$$= \sum_{\eta^\rho} p(\eta^\rho)\left[a_\mu(\eta^\rho) + a_\mu(\eta^\rho)^3 - 2a_\mu(\eta^\rho)^2\right]$$

$$= a^{\rho\mu} + (1 - \gamma^{\rho\mu})^2\frac{(a^{\rho\mu})^2}{a^\rho}\left[(1 - \gamma^{\rho\mu})\frac{a^{\rho\mu}}{a^\rho} - 2\right]$$

$$+ \frac{(\gamma^{\rho\mu}a^{\rho\mu})^2}{1 - a^\rho}\left[\frac{\gamma^{\rho\mu}a^{\rho\mu}}{1 - a^\rho} - 2\right],$$

$$(2.66)$$

to find that

$$\langle R^2\rangle = \frac{(C-1)(N-1)a^{\bar\rho\bar\mu}}{N^2 a^\rho(1 - a^\rho)}\left[1 + (1 - a^{\bar\rho})^2(1 - \gamma^{\bar\rho\bar\mu}) + (a^{\bar\rho})^2\gamma^{\bar\rho\bar\mu}\right]$$

$$+ \frac{M(C-1)(N-1)}{N^2 a^{\rho\mu}(1 - a^{\rho\mu})}\left[(1 - a^\rho)^2(1 - \gamma^{\rho\mu})a^{\rho\mu} + (a^\rho)^2\gamma^{\rho\mu}a^{\rho\mu}\right]$$

$$\times\left[1 + a^{\bar\rho\mu}\left[(1 - a^{\rho+1})^2(1 - \gamma^{\rho+1,\mu})a^{\rho+1,\mu} + (a^{\rho+1})^2\gamma^{\rho+1,\mu}a^{\rho+1,\mu}\right]\right].$$

$$(2.67)$$

Again this is of the form

$$\langle R^2\rangle = \frac{(N-1)(C-1)}{N^2}A_1 + \frac{(N-1)(C-1)M}{N^2}A_2, \qquad (2.68)$$

implying that, in the limit of large $N$ and $M$ the dominant term in the noise contribution to the local fields from all patterns other than the daughter cell cycle configuration of interest is

$$\langle R^2 \rangle \sim \frac{M(C-1)}{N} . \tag{2.69}$$

Thus, the daughter cell cycle will progress regularly provided that $N \gg M$ and an $\gamma$ is chosen such that the $S > 0$ at sites $\eta_i^{\rho\mu} = 1$ and $S > 0$ at sites $\eta_i^{\rho\mu} = 0$.

## APPENDIX 2.E   ORDER PARAMETERS FOR SPECIFIC CELL CYCLE CONFIGURATIONS

This appendix contains the calculation of the order parameters, $\widetilde{m}_\rho(\mathbf{n}(t))$ and $\widetilde{m}_{\rho\mu}(\mathbf{n}(t))$, when the system is in different levels of the cell hierarchy. First, $\widetilde{m}_\rho(\mathbf{n}(t))$ when the system is in the daughter cell cycle configuration. Equation (2.8) can be rewritten as follows by making use of the law of large numbers, $(N \to \infty)$ with $n_i = \eta_i^{\bar{\rho}\mu}$,

$$\widetilde{m}_\rho(\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}\mu}) = \left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1-a^\rho)} \eta^{\bar{\rho}\mu} \right\rangle . \tag{2.70}$$

If $\bar{\rho} \neq \rho$ then the daughter cell cycle phase is independent of the stem cell cycle phase and the expectation value factorises,

$$\widetilde{m}_\rho(\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}\mu}) = \left\langle \frac{\eta^\rho - a^\rho}{a^\rho(1-a^\rho)} \right\rangle \langle \eta^{\bar{\rho}\mu} \rangle , \tag{2.71}$$

and $\widetilde{m}_\rho(\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}\mu}) = 0$ since $\langle \eta^\rho - a^\rho \rangle = 0$.

However, if $\bar{\rho} = \rho$ then (2.70) can be written as,

$$
\begin{aligned}
\widetilde{m}_\rho(\mathbf{n} = \boldsymbol{\eta}^{\rho\mu}) &= \mathbb{E}\left[\mathbb{E}\left[\frac{\eta^\rho - a^\rho}{a^\rho(1 - a^\rho)}\eta^{\rho\mu}\middle|\eta^{\rho\mu}\right]\right] \\
&= \sum_{\eta^\rho}\frac{\eta^\rho - a^\rho}{a^\rho(1 - a^\rho)}p(\eta^\rho)\sum_{\eta^{\rho\mu}}\mathbb{W}(\eta^{\rho\mu}|\eta^\rho)\eta^{\rho\mu} \\
&= \frac{a^{\rho\mu}}{a^\rho}\left(\frac{1 - \gamma^{\rho\mu} - a^\rho}{1 - a^\rho}\right),
\end{aligned}
\tag{2.72}
$$

where $\mathbb{E}[x]$ and $\mathbb{E}[x|y]$ represents the expectation value of $x$ and of $x$ given $y$, respectively. Thus, provided $\gamma^{\rho\mu} < (1 - a^\rho)$, there is a positive (non-zero) value for the order parameter. This can be rewritten by noticing that the numerator is the covariance between $\eta^\rho$ and $\eta^{\rho\mu}$ and the denominator is the variance of $\eta^\rho$. Thus, for $n_i = \eta_i^{\rho\mu}\ \forall i$,

$$
\widetilde{m}_\rho(\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}\mu}) = \frac{\text{cov}\left[\eta^\rho, \eta^{\rho\mu}\right]}{\text{var}\left[\eta^\rho\right]}.
\tag{2.73}
$$

Next, $\widetilde{m}_{\rho\mu}(\mathbf{n}(t))$ when the system is in a stem cell cycle configuration. Following similar reasoning to the above, (2.9) can be written as,

$$
\widetilde{m}_{\rho\mu}(\mathbf{n} = \boldsymbol{\eta}^{\bar{\rho}}) = \left\langle\frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})}\eta^{\bar{\rho}}\right\rangle,
\tag{2.74}
$$

where again it is trivial that $m_{\rho\mu} = 0$ due to independence if $\bar{\rho} \neq \rho$.

For the case $\bar{\rho} = \rho$,

$$
\begin{aligned}
\widetilde{m}_{\rho\mu}(\mathbf{n} = \boldsymbol{\eta}^\rho) &= \mathbb{E}\left[\mathbb{E}\left[\frac{\eta^{\rho\mu} - a_\mu(\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})}\eta^\rho\middle|\eta^{\rho\mu}\right]\right] \\
&= \sum_{\eta^\rho}\eta^\rho p(\eta^\rho)\sum_{\eta^{\rho\mu}}\frac{(\eta^{\rho\mu} - a_\mu(\eta^\rho))\mathbb{W}(\eta^{\rho\mu}|\eta^\rho)}{a^{\rho\mu}(1 - a^{\rho\mu})}.
\end{aligned}
\tag{2.75}
$$

Then since $\sum_{\eta^{\rho\mu}}\mathbb{W}(\eta^{\rho\mu}|\eta^\rho) = 1$ and $\sum_{\eta^{\rho\mu}}\mathbb{W}(\eta^{\rho\mu}|\eta^\rho)\eta^{\rho\mu} = a_\mu(\eta^\rho)$ by their definitions, $\widetilde{m}_{\rho\mu}(\mathbf{n} = \boldsymbol{\eta}^\rho) = 0$. Therefore, during any phase of the stem cell cycle all of the order parameters

for each of the daughter cell cycle phases vanish.

As the overlap between a stem cell and its progeny increases the less distinguishable are their cellular identities, because their corresponding attractors become closer in the state space. Because of this, and that the activities in each cell cycle phase are fixed, binarized gene expression level data of known cell types could be used to determine the probabilities $\gamma^{\rho\mu}$ for differentiation events

$$\widetilde{m}_\rho(\mathbf{n}(t) = \boldsymbol{\eta}^{\rho\mu}) = \frac{a^{\rho\mu}(1 - \gamma^{\rho\mu} - a^\rho)}{a^\rho(1 - a^\rho)} \tag{2.76}$$

which can be rearranged to calculate the probability that gene is suppressed during differentiation,

$$\gamma^{\rho\mu} = (1 - a^\rho)\left(1 - \frac{a^\rho}{a^{\rho\mu}}\widetilde{m}_\rho(\boldsymbol{\eta}^{\rho\mu})\right). \tag{2.77}$$

## APPENDIX 2.F    STABILITY ANALYSIS OF EQUATIONS OF MOTION

The equations of motion (2.11) & (2.12) can be written in the general multivariate form,

$$\widetilde{\mathbf{m}}_{t+1} = \mathbf{F}(\widetilde{\mathbf{m}}_t), \tag{2.78}$$

where $\widetilde{\mathbf{m}}_t$ is a vector containing all of the order parameters for the stem and daughter cell cycle phases $\{\widetilde{\mathbf{m}}_\rho(t), \widetilde{\mathbf{m}}_{\rho\mu}(t)\}$. Furthermore, because of the periodicity of each cell cycle, applying the function $\mathbf{F}$ $C$ times to a state in one of the cell cycle configurations will return system to its initial phase, where $C$ is the number of stages in the cell cycle. Thus, for an attractor $\widetilde{\mathbf{m}}^*$ one has

$$\widetilde{\mathbf{m}}^* = \widetilde{\mathbf{m}}_t^* = \widetilde{\mathbf{m}}_{t+C}^* = \mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t), \tag{2.79}$$

where the notation $\mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t)$ indicates that the function $\mathbf{F}$ is applied $C$ times to the argument, i.e. $\mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t) = \mathbf{F}(\ldots \mathbf{F}(\mathbf{F}(\widetilde{\mathbf{m}}_t)))$. If the system is sufficiently close to an attractor $\widetilde{\mathbf{m}}^*$ of the dynamics, so at time t, $\widetilde{\mathbf{m}}_t = \widetilde{\mathbf{m}}_t^* + \boldsymbol{\delta}_t$, where $\boldsymbol{\delta}_t$ is a small, then we can Taylor expand the function

$\mathbf{F}^{(C)}$ about the deviation from the attractor. Starting from,

$$\widetilde{\mathbf{m}}_t^* + \boldsymbol{\delta}_t = \widetilde{\mathbf{m}}_{t+C}^* + \boldsymbol{\delta}_{t+C} = \mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t) \, , \qquad (2.80)$$

and Taylor expanding

$$\boldsymbol{\delta}_{t+C} = \triangledown\mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t^*)(\widetilde{\mathbf{m}}_t - \widetilde{\mathbf{m}}_t^*) + \dots \, , \qquad (2.81)$$

but the difference $\widetilde{\mathbf{m}}_t - \widetilde{\mathbf{m}}_t^*$ is just the distance from the attractor at time $t$, so this we can rewrite this as

$$\boldsymbol{\delta}_{t+C} = \triangledown\mathbf{F}^{(C)}(\widetilde{\mathbf{m}}_t^*)\boldsymbol{\delta}_t \, . \qquad (2.82)$$

Then using the progressive nature of the cell cycle attractors, the gradient can be rewritten in terms of the function $\mathbf{F}$ applied at the previous time step

$$\boldsymbol{\delta}_{t+C} = \triangledown\mathbf{F}\left(\mathbf{F}^{(C-1)}(\widetilde{\mathbf{m}}_t^*)\right)\boldsymbol{\delta}_t \, . \qquad (2.83)$$

Next using the chain rule one can evaluate the gradient of the function $\mathbf{F}$ to give

$$\boldsymbol{\delta}_{t+C} = \triangledown\mathbf{F}\left(\mathbf{F}^{(C-1)}(\widetilde{\mathbf{m}}_t^*)\right) \cdot \triangledown\mathbf{F}^{(C-1)}(\widetilde{\mathbf{m}}_t^*)\boldsymbol{\delta}_t \, , \qquad (2.84)$$

Repeating these last two steps of rewriting the function $\mathbf{F}^\tau(\dots)$ as $\mathbf{F}\left(\mathbf{F}^{(\tau-1)}(\dots)\right)$, and evaluating the gradient using the chain rule, one obtains

$$\boldsymbol{\delta}_{t+C} = \triangledown\mathbf{F}(\widetilde{\mathbf{m}}_{t+C-1}^*) \cdot \triangledown\mathbf{F}(\widetilde{\mathbf{m}}_{t+C-2}^*) \cdot \dots \cdot \triangledown\mathbf{F}(\widetilde{\mathbf{m}}_t^*)\boldsymbol{\delta}_t \, , \qquad (2.85)$$

which can be written as the following product

$$\delta_{t+C} = \prod_{\tau=C}^{0} \left[ \nabla \mathbf{F}(\widetilde{\mathbf{m}}_{t+\tau}^*) \right] \delta_t . \tag{2.86}$$

Thus the fixed point will be a stable attractor provided that the eigenvalues $\lambda$ of the matrix $\mathbf{B}$ are less than 1, where $\mathbf{B} = \prod_{\tau=C-1}^{0} \left[ \nabla \mathbf{F}(\widetilde{\mathbf{m}}_{t+\tau}^*) \right]$.

To find the eigenvalues of the matrix $\mathbf{B}$ one must first obtain the partial derivatives of the equations of motion in order to find the Jacobian at each time step, i.e. the matrix $\nabla F(\widetilde{\mathbf{m}}_{t+\tau}^*)$, to be evaluated at the attractor. The elements of the Jacobian are given by

$$\frac{\partial \widetilde{m}_\rho(t+1)}{\partial \widetilde{m}_{\bar{\rho}}(t)} = \frac{\beta}{4} \left\langle \frac{(\eta^\rho - a^\rho)(\eta^{\bar{\rho}+1} - a^{\bar{\rho}+1})}{a^\rho(1-a^\rho)} \operatorname{sech}^2 \left( \frac{\beta h(t)}{2} \bigg|_{\widetilde{\mathbf{m}}_t^*} \right) \right\rangle_{\eta^\rho, \eta^{\rho\mu}} , \tag{2.87}$$

$$\frac{\partial \widetilde{m}_\rho(t+1)}{\partial \widetilde{m}_{\bar{\rho}\bar{\mu}}(t)} = \frac{\beta}{4} \left\langle \frac{(\eta^\rho - a^\rho)(\eta^{\bar{\rho}+1,\bar{\mu}} - a_{\bar{\mu}}(\eta^{\bar{\rho}+1}))}{a^\rho(1-a^\rho)} \operatorname{sech}^2 \left( \frac{\beta h(t)}{2} \bigg|_{\widetilde{\mathbf{m}}_t^*} \right) \right\rangle_{\eta^\rho, \eta^{\rho\mu}} , \tag{2.88}$$

$$\frac{\partial \widetilde{m}_{\rho\mu}(t+1)}{\partial \widetilde{m}_{\bar{\rho}}(t)} = \frac{\beta}{4} \left\langle \frac{(\eta^{\rho\mu} - a_\mu(\eta^\rho))(\eta^{\bar{\rho}+1} - a^{\bar{\rho}+1})}{a^{\rho\mu}(1-a^{\rho\mu})} \operatorname{sech}^2 \left( \frac{\beta h(t)}{2} \bigg|_{\widetilde{\mathbf{m}}_t^*} \right) \right\rangle_{\eta^\rho, \eta^{\rho\mu}} , \tag{2.89}$$

$$\frac{\partial \widetilde{m}_{\rho\mu}(t+1)}{\partial \widetilde{m}_{\bar{\rho}\bar{\mu}}(t)} = \frac{\beta}{4} \left\langle \frac{(\eta^{\rho\mu} - a_\mu(\eta^\rho))(\eta^{\bar{\rho}+1,\bar{\mu}} - a_{\bar{\mu}}(\eta^{\bar{\rho}+1}))}{a^{\rho\mu}(1-a^{\rho\mu})} \operatorname{sech}^2 \left( \frac{\beta h(t)}{2} \bigg|_{\widetilde{\mathbf{m}}_t^*} \right) \right\rangle_{\eta^\rho, \eta^{\rho\mu}} .$$
$$\tag{2.90}$$

At each time step, i.e. phase of the cell cycle, only the diagonal part of these partial derivatives will contribute, i.e. $\rho = \bar{\rho}$ and $\mu = \bar{\mu}$, because the gene expression levels are assumed to be independent along the cell cycles. For each cell type, there is also only a single non-zero order
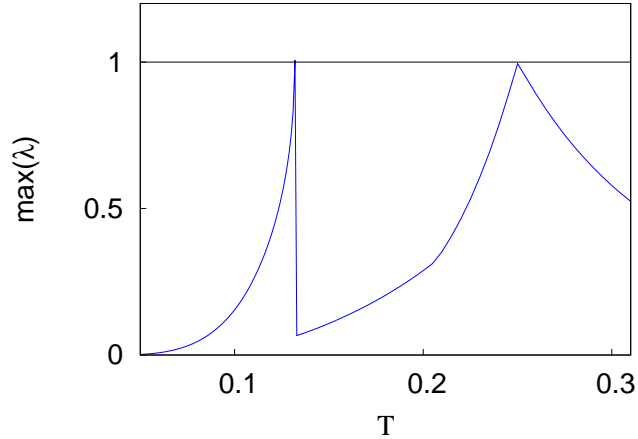
**Figure 2.F.1:** The maximum eigenvalue of the $\mathbf{B}$ as a function of the biological noise strength $T$. When $\max(\lambda) = 1$ there is a transition in the dynamics between solutions of the equations of motion for the dynamic order parameters (2.13) and (2.14). The same parameter values were used as in figure 2.4.3.

parameter, at each time and for each cell type, which is the order parameter of the cell cycle phase that the system is expected to be in given the initial conditions.

Using the partial derivatives above one can construct a Jacobian $\bigtriangledown\mathbf{F}(\widetilde{\mathbf{m}}_t)$, which needs be evaluated at each stage of a cell cycle to calculate $\mathbf{B} = \prod_{\tau=C}^{0} \left[\bigtriangledown\mathbf{F}(\widetilde{\mathbf{m}}_{t+\tau}^*)\right]$. Figure 2.F.1 shows the maximum eigenvalue of $\mathbf{B}$ for different values of the noise strength $T$, when the system was initialised in a daughter cell cycle. The maximum eigenvalue takes a value of 1 whenever a solution becomes unstable. Starting at low $T$, the system is initially stable with respect to the daughter cell cycle, until the critical $T$ for noise induced reprogramming is reached. Above this noise level there is a new solution that corresponds to the stem cell cycle attractor. This solution is then stable for intermediate values of $T$. If $T$ becomes too large the maximum eigenvalue reaches 1 again. This corresponds to the transition from the stem cell cycle to a high noise state in which the system has no overlap with any of the cell cycle patterns. The values of $T$ at which these transitions are found from the maximum eigenvalue are in exact agreement with those observed in figure 2.4.3.

APPENDIX 2.G    NORMALISED HAMMING DISTANCE

The Hamming distance between two vectors $\mathbf{x} = (x_1, x_2, \ldots x_N)$ and $\mathbf{y} = (y_1, y_2, \ldots y_N)$ is defined as follows,

$$\mathrm{d}\left[\mathbf{x}, \mathbf{y}\right] = \frac{1}{N} \sum_{i=1}^{N} |y_i - x_i| \ . \tag{2.91}$$

In the large $N$ limit we can replace the sum over $N$ using the law of large numbers to obtain,

$$\mathrm{d}\left[\mathbf{x}, \mathbf{y}\right] = \langle |y_i - x_i| \rangle \ . \tag{2.92}$$

Thus, for the same phase of the cell cycle the Hamming distance between the stem and daughter cell cycles is given by,

$$\mathrm{d}\left[\boldsymbol{\eta}^{\rho\mu}, \boldsymbol{\eta}^{\rho}\right] = a^{\rho} - a^{\rho\mu} + 2\gamma^{\rho\mu} a^{\rho\mu} \ , \tag{2.93}$$

where the averages were performed over the joint probability distribution, using the conditional and marginal distributions $p(\eta^{\rho}, \eta^{\rho\mu}) = W(\eta^{\rho\mu}|\eta^{\rho})p(\eta^{\rho})$. Similarly, the Hamming distance between a daughter cell cycle phase and the next phase of the stem cell cycle is given by,

$$\mathrm{d}\left[\boldsymbol{\eta}^{\rho\mu}, \boldsymbol{\eta}^{\rho+1}\right] = a^{\rho+1} - a^{\rho\mu} - 2a^{\rho\mu} a^{\rho+1} \ , \tag{2.94}$$

where the average was performed over the joint probability that factorises due to the independence of gene expression levels in different cell cycle phases, i.e. $p(\eta^{\rho+1}, \eta^{\rho\mu}) = p(\eta^{\rho+1})p(\eta^{\rho\mu})$

# 3

# MODEL VALIDATION

## 3.1. INTRODUCTION

Cell cycle specific gene expression profiles are, at the time of writing, incredibly rare, not least because of the difficulty involved in obtaining such data from experiments. The cell reprogramming model introduced in the previous chapter has already been shown to quantitatively and qualitatively capture aspects of cell reprogramming experiments. However, in this chapter, further attempts are made to assess the practical value of said model by using data from real experiments to test assumptions, parameters, and the key hypothesis of the model, i.e. there exists a single cell cycle phase that is maximally similar across different stages of the cell cycle.

Two data sets were used to achieve the results that follow. The first, that is publicly available, comes from The Human Protein Atlas and consists of RNA transcripts per million (TPM) for 64 different cell lines and 37 different tissue types [102, 103]. The unit TPM is a normalised unit which implies that for every 1 million RNA molecules sequenced from a sample $x$ of them came from that transcript/gene. The other data set, provided by the authors of [104], also includes transcript counts normalised to TPM, and comes from a study of the mouse enteric nervous system (ENS), i.e the system of neurons that surrounds the gut of a mouse. A summary of both datasets is given in table 3.1.1. Each of these datasets offers different qualities from a model validation perspective. In the Human data set each cell type is known and unique. The variation in cell types is great because of the diversity of 13 tissue types available in the database. This should give us a good estimate for ranges of individual parameters such as activities for the reprogramming model. Where the human data falls short is that it is not cell cycle specific and, therefore, could represent an average across multiple cell cycle phases.

On the other hand, the samples in the mouse data have each been labelled with a specific cell cycle phase. The cycle phases were determined using the maximum likelihood method available in the SCRAN package for the R programming language [105] by the authors of the study. This assignment allows the comparison of gene expression profiles of different samples in the same (or different) cell cycle phase. However, the mouse data is a lot less diverse, with all the cells sequenced coming from the ENS. This means that parameters such as the activity deduced from the data may be representative of only the ENS and not of the mouse as a whole. Furthermore, each cell in this dataset was only given an arbitrary label before the cell cycle phase assignment. Therefore, it is not possible to know if any two cells are of the same cell type and in different/the same cell cycle phases, or if they are different cell types altogether.

| Source | The Human Protein Atlas | R. Lasrado et al. [104] |
|---|---|---|
| Species | Human | Mouse |
| No. of sequenced genes | 19,628 | 9,628 |
| No. of sequenced cells | 56 | 120 |
| Cell cycle specific labels | N/A | $G_1$, $S$, $G_2/M$ |

**Table 3.1.1:** The transcript count data from RNA-sequencing experiments used in this chapter.

## 3.2. BINARY GENE EXPRESSION LEVELS

Before we can determine parameters, such as activities, the data first needs to be converted into the appropriate format. One of the assumptions that we make throughout this thesis, and a common one used in theoretical biology, is that genes have binary expression levels. The transcript count data can be converted to binary gene expression levels using a simple conditional statement, i.e. if the value of the transcript count is above a certain threshold the gene is expressed $\eta = 1$, otherwise, the gene is not expressed. Mathematically, this is equivalent to using a Heaviside step function $\eta_i = \Theta[x_i - \hat{x}]$, where $x_i$ is the transcript count corresponding to gene $i$ and $\hat{x}$ is the threshold value above which we say a gene is expressed. We choose the value of the threshold $\hat{x} = 0$ as this gives good agreement with reported values for the fraction of expressed genes in humans [22, 23]. However, due to the small number of molecules often involved in biological processes, even marginal changes to $\hat{x}$ can have a dramatic effect on the parameter values extracted from the data (see appendix 3.A).

The activity of each cell is determined by calculating the average of the binarized gene expression levels $a(\boldsymbol{\eta}) = \frac{1}{N_{genes}} \sum_i^{N_{genes}} \eta_i$, and is plotted in figure 3.2.1. The average activities are is $0.74$ and $0.72$ for the human and mouse cells respectively. These values are comparable to those reported in [22, 23] where different human tissues types had $\sim 70\%$ of genes expressed. The mouse cells have been organised into the three different cell cycle stages reported in the data set - $G_1$, $G_2/M$ & $S$. The average activity is lower for the $S$ and $G_2/M$ cell cycle stages. However, the

number of cells in these stages, $25$ and $15$ cells in $G_2/M$ & $S$ respectively, is markedly lower than the 80 cells in the $G_1$ phase. This could be due to the duration of the cell cycle phases. Because the long duration of the $G_1$ phase, cells that are sequenced at random and labelled by their cycle phase afterwards are much more likely to be in $G_1$ than any other phase. Hence, the bias of the population towards the $G_1$ phase. The difference in average activities across the phases could be a direct result of the smaller sample sizes for $G_2/M$ and $S$ phases, or a true characteristic that differentiates $G_1$ from the other two phases. If one believes the latter is true, then this suggests that the $G_1$ cell cycle phase could be the phase of maximal similarity across different cell types due to an increased fraction of expressed genes. This opposes the line of thought in chapter 2, where the common molecular machinery involved for processes occurring in the $S$ and $M$ phases were thought to result in a high similarity across different cell types.

The bimodal nature of the distribution of activities calculated from the gene expression levels of the mouse cells could be an artefact of the choice of the threshold used to binarize the data (see appendix 3.A). Increasing the threshold marginally above its current value $\hat{x} = 0$ removes the bimodal structure of the distribution (see appendix 3.A). However, we will continue to use $\hat{x} = 0$ because it gives a good agreement for the average fraction of expressed genes with other studies, as mentioned previously.

The activities from the human data show less structure than those of the mouse ENS. The activities of the human tissue data fluctuate around the average activity. On the other hand, the mouse activities have a bimodal distribution, with a sharp peak at $0.84$ representing the majority of the data and a broader peak at $0.6$. Although, for the $G_2/M$ and $S$ cells the minority have this high gene expression level, and the majority of activities are $a < 0.75$. As mentioned earlier this distinction could be down to a true difference caused by the processes involved at each stage of the cell cycle or it may be a finite size effect caused by the low number of samples in the $S$ and $G_2/M$ phases.

The variance of the gene expression levels also show that, whilst the organisms' gene expres-
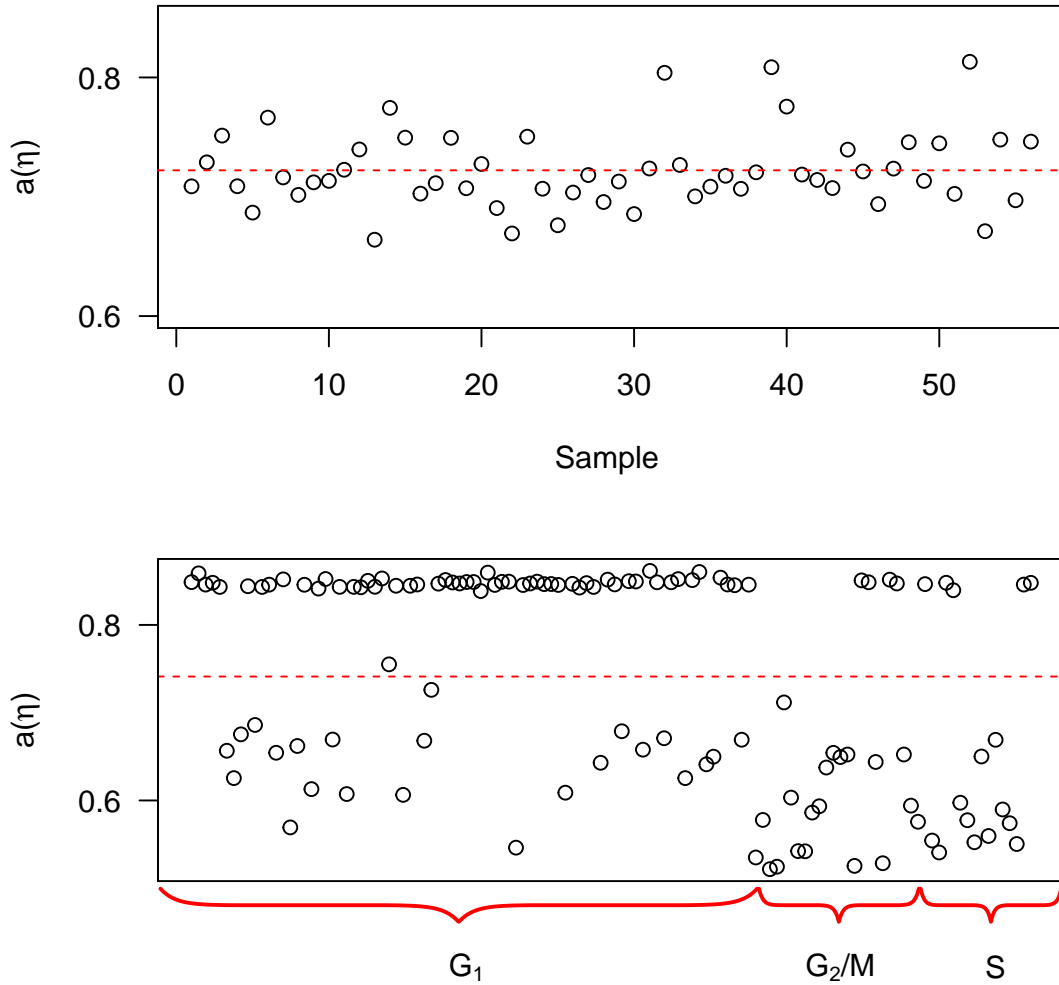
**Figure 3.2.1:** Activities for the different cell samples taken from human (top) and the mouse ENS (bottom). The mouse data are organised based on their maximum likelihood cell cycle phase. The red dashed lines in both plots give the average activity across all samples for each organism - $\langle a^{human} \rangle = 0.74$ and $\langle a^{mouse} \rangle = 0.72$.

sion levels fluctuate on a similar scale, the mouse data are much more localised about the average compared to the human data (figure 3.2.2).

## 3.3. CELL CYCLE SIMILARITIES

One of the key assumptions in the model presented in chapter 2 is that there is at least one cell cycle phase that is maximally similar across different cell types of an organism. The labelling of the mouse ENS data into 3 distinct cell cycle phases ($G_1$, $G_2/M$ & $S$) allows us to check this hypothesis. We have already seen that the $G_1$ cells of the mouse ENS have on average a higher activity than those in other cycle phases. In this section, we will study the gene expression levels further using similarity measures. The correlation between the gene expression levels of each cell sample of the mouse data gives a level of similarity between pairs of cells. A correlation of $\text{cor}\left[\eta, \eta'\right] = 1$ indicates that two cells have identical gene expression level patterns. Generally, the higher the correlation the greater level of similarity between any two gene expression profiles. A heat-map of the correlations between the gene expression patterns of each of the mouse ENS cells is shown in figure 3.3.3. (Correlations of gene expression patterns in the data from the Human Protein Atlas can be found in appendix 3.B.) The gene expression patterns of cells in the $G_1$ phase are highly correlated with those of the majority of other $G_1$ cells, whereas gene expression patterns of cells in $G_2/M$ and $S$ phases typically have lower levels of correlations with those of all other cells. This is further evidence in support that there is, at least for the mouse ENS data, a single cell cycle phase that is maximally similar across different cell types - the $G_1$ phase.

The covariance between the different gene expression level profiles reveals the same information. With the variances of the gene expression levels in each cell, they can also provide a predicted overlap between the different cell samples of the mouse ENS data using equation (2.20) or (2.21). The covariances and the resulting overlaps are plotted in figure 3.3.4. Unlike the plots of the correlation and covariance, the overlaps are not symmetric about the diagonal. This asym-
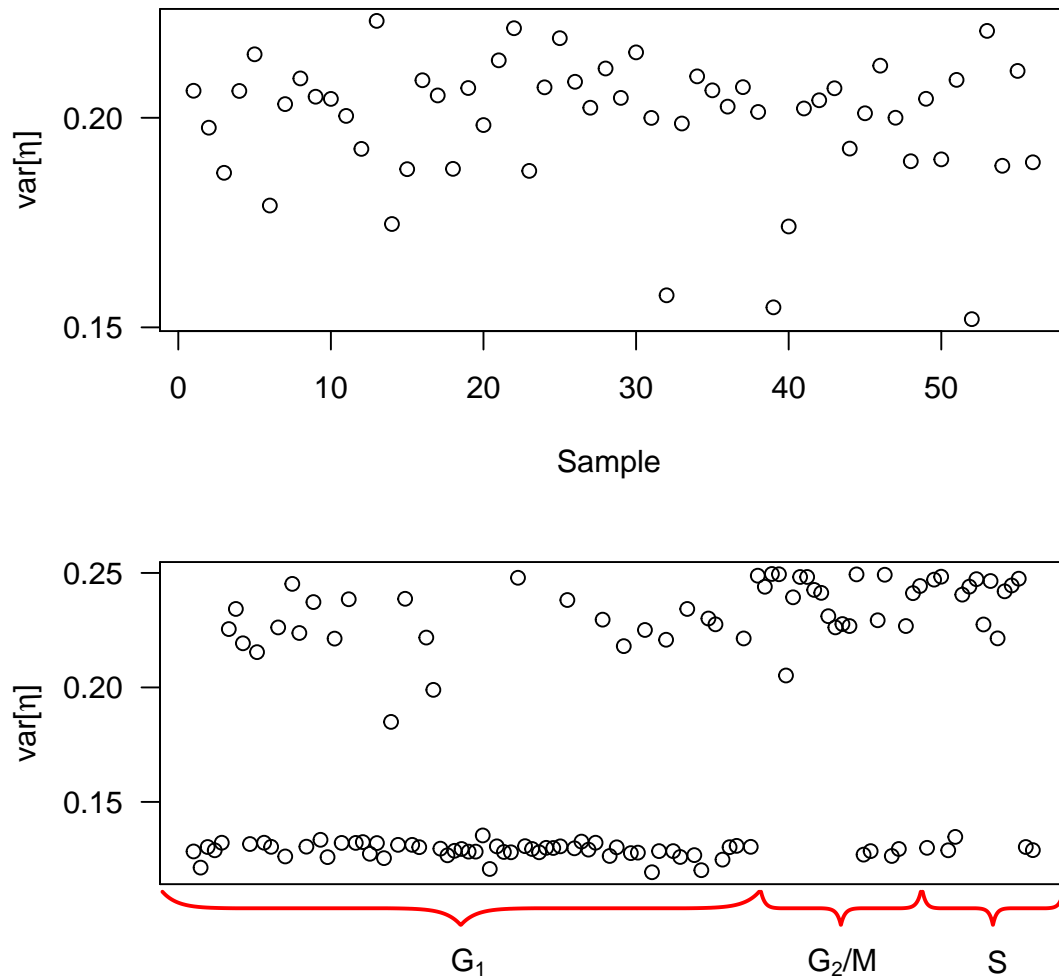
**Figure 3.2.2:** Variances of the gene expression levels for the different cell samples taken from the human (top) and mouse (bottom) data. The mouse samples are organised by their maximum likelihood cell cycle phase.

**Figure 3.3.3:** Correlations between the gene expression profiles of the different samples of the mouse ENS data, grouped according to cell cycle phase. The colour of each point in the plot gives the level of correlation between the two corresponding cell samples.

metry is caused by the choice of level of the hierarchy that the two cells are in when calculating the overlap between them - recall from chapter 2, (2.20) and (2.21) are valid for a two-level hierarchy where one cell is in a higher state of potency. Because, the samples in the mouse data are not labelled based on cell type/developmental time, all possible overlaps are plotted in figure 3.3.4. Thus, for a pair of gene expression patterns in cells $A$ and $B$, the overlap is calculated assuming $A$ and $B$ are the "stem" and "daughter" cells respectively, and then, assuming cell $B$ is the "daughter" cell and $A$ is the "stem" cell. In figure 3.3.4 the vertical and horizontal axis represents the choice of "stem" and "daughter" cells respectively, i.e. each "row" in the heat map corresponds to choosing that cell as the stem cell in a 2-level potency hierarchy (as in figure 2.2.3) and all other cells (including itself and cells labelled in different cycle phases) are then in the lower level of the hierarchy as daughter cells.

The overlaps between cells in different cell cycle phases in chapter 2 are zero. However, all overlaps calculated from the mouse data are non-zero. This difference arises as a direct result of

**Figure 3.3.4:** The covariance (top) between gene expression profiles in the mouse data and the overlap (bottom) predicted from them using equations (2.20) and (2.21). For the overlaps, the vertical and horizontal axes represent the choice of stem and daughter cells respectively. Cells in the $G_1$ phase have the highest level of similarity between different samples.

**Figure 3.3.5:** The probability that a gene is switched off during a hypothetical differentiation from one cell in the sample to another in the mouse ENS data set using (2.77). As with the bottom panel of figure 3.3.4, the vertical and horizontal axes represent the choice of stem and daughter cells used to calculate the probabilities respectively.

one of the assumptions in the hierarchy of cell cycles model, that is, the gene expression levels in different cell cycle phases are independent, and therefore, cell cycle phases uncorrelated. This is not the case for the mouse ENS data, in which all samples have some non-negligible degree of correlation (as shown in figure 3.3.3).

Supposing that the mouse data was compatible with the assumptions within our model, including the independent cell cycle phases, one would be able to infer the values of the probability that a gene is switched off during differentiation from stem to daughter cell, whilst keeping the activity constant, i.e. $\gamma^{\rho\mu}$, from the overlaps and activities. Using (2.77), the probabilities $\gamma^{\rho\mu}$ calculated from the activities of the mouse samples, and the overlaps between those samples, are plotted in figure 3.3.5.

## 3.4. DYNAMICS

With the binarized gene expression levels from the Human Protein Atlas and Mouse ENS datasets, the cell sample activities and $\gamma$ probabilities determined, it is possible to test the dynamics from chapter 2 on the data. The mouse data are labelled by 3 cell cycle phases, so it easy to store the gene expression levels into the cell cycle attractors of the model. However, because the cell types of each sample are unknown placing the stored patterns in levels of a potency hierarchy remains a challenge.

This challenge is not the biggest to running the dynamics with the mouse ENS data stored in the interactions $J_{ij}$ though. As shown in figure 3.3.3, there are significant correlations in the gene expression levels between all pairs of the cells in the dataset. Because of the assumption in the model that the gene expression levels around the cell cycle are independent of one another, these correlations result in the breakdown of the model when it comes to its dynamics, even at low levels of the biological noise $T$, as can be seen in figure 3.4.6. The high correlations, and resulting high overlaps, between each stored gene expression level configuration result in the dynamics rapidly becoming trapped in a mixture state that has a high non-zero overlap with all stored configurations (figure 3.4.6). This behaviour is typical of the model when the binary gene expression levels of the mouse ENS data are stored as attractors in the interactions, regardless of the choice of which configurations are stem or daughter cells. The inability to distinguish between stored attractors could not only have been predicted by the high correlations between the mouse ENS data samples, but also from the probabilities $\gamma^{\rho\mu}$ which are typically small (figure 3.3.5). However, there are methods for storing correlated patterns in Hopfield-like neural network models. The most well known is the pseudoinverse learning rule [106]. This method adapts the storage prescription of the Hopfield model to include the inverse of the correlation between each to the stored patterns, i.e. $J_{ij} = \frac{1}{n} \sum_{\mu,\nu} \xi_i^\mu \xi_j^\nu (\mathbf{C}^{-1})_{\mu\nu}$ where $C_{\mu\nu} = \frac{1}{N} \sum_i \xi_i^\mu \xi_i^\nu$ is the correlation between patterns $\mu$ and $\nu$. However, to be compatible with our model would require the adaptation
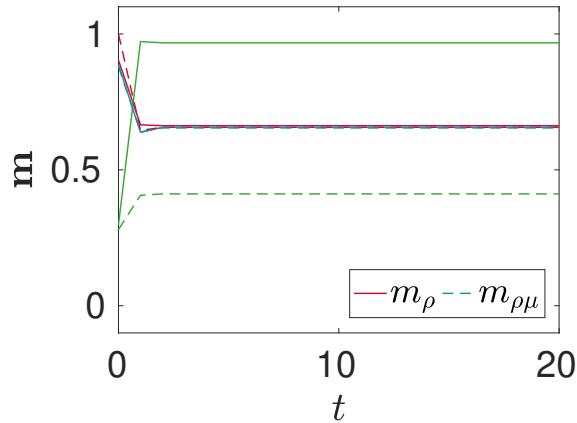
**Figure 3.4.6:** A typical evolution of the overlap between the gene expression levels from a simulation with mouse ENS cell cycle configurations stored in the interactions, $J_{ij}$. The system was initialised in a randomly selected mouse ENS daughter cell cycle phase and evolved deterministically $(T = 0)$. The system is rapidly trapped in a state that has a constant overlap with each configuration of the chosen stem and daughter cell cycle phases.

of the pseudoinverse rule to work with a hierarchy of cycles. The adaptation to cycles would simply involve replacing $\nu$ with $\mu + 1$ in $J_{ij}$ and has been studied for a variety of neural network structures [107], but further adapting the rule to store a hierarchy of these cycles would require careful considerations.

Working with the binary gene expression levels from the Human Protein Atlas dataset also presented challenges. Mainly, the cell cycle phase in which each cell was sampled from are unknown. One possibility would be to assign a cell cycle phase using a maximum likelihood method, similar to the preprocessing performed on the mouse ENS data by the authors of [104]. However, this is a large time-consuming task in itself, and thus, beyond the scope of this thesis. Instead, a single high potency cell line (HeLa) was selected as a template to generate gene expression levels for a synthetic stem cell cycle using a Markov transition matrix, similar to the method detailed in appendix 2.A. Then from this stem cell cycle, a set of daughter cell cycle specific gene expression levels was also generated using a Markov transition matrix. These synthetic cells were generated with transition matrices that maintain the activity of the HeLa gene expression profile in every configuration, but due to the manner in which they are generated are compatible with the dynam-

**Figure 3.4.7:** Simulation results using synthetic stem and daughter cells derived from the HeLa gene expression profile. The system is initialised in one of the daughter cell cycle phases. The parameters used for the simulation were taken from the HeLa cell, $N = 19,628$, $a^{\rho} = a^{\rho\mu} = 0.7$, apart from $\gamma^{\rho\mu} = 0.3 \; \forall \{\rho, \mu\}$ which was chosen to reduce the overlap between the stem and daughter cell cycles. Left - low noise evolution of the overlap of the system with the stem and daughter cells shows the persistence of the daughter cell cycle. Right - The envelope of the evolution of the dynamics at a high noise level ($T = 0.12$) showing the noise induced reprogramming. For both plots dashed and solid lines are the overlaps with the daughter and stem cell cycles respectively.

ics of the model in chapter 2.

When initialised in a cell cycle phase of the synthetic daughter cell, at low noise levels, the daughter cell cycle is sustained as expected. However, at high noise levels the synthetic stem cell cycle is retrieved from the daughter cells (see figure 3.4.7). Thus, the reprogramming model is compatible with gene expression patterns that have the same statistics as human cells. However, the model should be developed further to be compatible with correlations across the cell cycle in order to store real cell cycle specific gene expression levels, like those of the mouse ENS dataset, as opposed to working with synthetic ones.

## 3.5. Summary

In this chapter, gene expression data from human and mouse ENS cells from RNA sequencing experiments were converted into binary gene expression levels. These gene expression levels were

then used to determine typical values of parameters for the model presented in chapter 2. Parameters such as the average activities and the probability of silencing a gene on the differentiation between cells with the same activity in a two-level potency hierarchy are easily determined from the datasets and were found to lie in the ranges used in chapter 2.

The data was also used to test a number of assumptions made in the model. There is significant evidence in the mouse ENS data that there is indeed a single cell cycle phase, the $G_1$ phase, that is maximally similar across different cells. Cells in the $G_1$ phase not only have high correlations with other $G_1$ cells but have activities which are greater on average than cells in the $S$ and $G_2/M$ phases. The $S$ and $G_2/M$ cells have low correlations with other cells in the same and different phases. Interestingly, this goes against the original reasoning in chapter 2, where the cell cycle phase hypothesised to be maximally similar across different cell types was argued to be the $S$ or $M$ phases. This argument was based on the idea that the molecular machinery involved in the processes of duplicating the DNA in the $S$ phase, and separating chromatin $M$ in the phase, are consistent across different cell types. However, the $G_1$ phase is a crucial phase in which cells do not only grow, but they are prepared to enter the $S$ phase and have to pass multiple checkpoints. This, combined with the longer duration of the $G_1$ cell cycle phase may result in the higher correlations seen in gene expression levels of between different cell types in that phase. However, further studies are needed to understand the difference in correlations between the gene expression levels of different cells in the same cell cycle phases. Although the analysis presented in this chapter supports the hypothesis that there is a single maximally similar cell cycle phase across different cell types, due to the size of the cell cycle specific datasets that were studied, it is desirable that further investigation is carried out as such data becomes more readily available. This could reinforce the evidence for a single maximally similar cell cycle phase between different cells, as presented in this chapter, or assist in determining if the results presented are largely due to the limited size of the currently available data.

Trying to incorporate the human and mouse data into the dynamics of the model of the pre-

vious chapter has also further illuminated the need to develop that model to include correlations between the gene expression levels between phases of the cell cycle. The strong correlations between gene expression levels of different cycle phases and different cell types resulted in the attractors of a two-level hierarchy becoming indistinguishable. Often the dynamics becomes trapped in a mixture state and the model breaks down because of its assumption of independent cell cycle phases. Thus, whilst the model is successful in capturing some key aspects of the cell reprogramming process, it would be of great interest to build upon its current form. Methods exist for storing correlated patterns in Hopfield like neural networks. The most successful of which is the pseudoinverse method [106], in which the Hebbian learning prescription is adapted by including the inverse correlation matrix of the stored patterns in the interaction matrix. Adapting the interactions in chapter 2 using the reasoning of the pseudoinverse rule combined with the cycle and hierarchy storage prescriptions should allow the model to work with correlated gene expression patterns, such as those in the human and mouse datasets presented in this chapter.

# Appendices

## Appendix 3.A  Choice of threshold for binary gene expression levels

The data used in this chapter are in the form of transcription counts per million and have been converted into binary gene expression levels using a simple conversion,

$$\eta_i = \Theta[x_i - \hat{x}] \tag{3.1}$$

i.e. if the transcription count, $x_i$, is greater than some threshold value, $\hat{x}$, then gene $i$ is expressed $\eta_i = 1$, otherwise it is not $\eta_i = 0$. Thus, the choice of the threshold $\hat{x}$ plays an important role in determining which genes are expressed and, therefore, dictates values of parameters extracted from the data such as the activity of each cell sample. The average activity for each data set is plotted against the threshold $\hat{x}$ in figure 3.A.1.

The average activity of the mouse ENS samples is highly susceptible to the choice of $\hat{x}$, as $\langle a \rangle$ decays rapidly to zero over a small range of $\hat{x}$. Changing the threshold from $\hat{x} = 0$ to $\hat{x} = 1.5$ TPM approximately halves the average fraction of genes that are expressed in the mouse ENS samples. The change in $\langle a \rangle$ with $\hat{x}$ for the mouse data is non-linear, with the variance of the activities also decreasing with $\hat{x}$. On the other hand, the average activity of the human cells decreases nearly linearly, in the same range plotted, at a much slower rate. The variation in activities for the human

**Figure 3.A.1:** The average activity, $\langle a \rangle$, in the Human Protein Atlas (blue) and Mouse ENS (red) data sets for different values of the threshold used to convert the transcript counts into binary gene expression levels, $\hat{x}$. The bars show one standard deviation from the average activity.

data set remains roughly constant for the range of $\hat{x}$ shown in figure 3.A.1.

For the analysis presented in this chapter, the binary threshold of $\hat{x} = 0$ was used. This choice was made because it gives good agreement with [22, 23] for the Human data. If we assume that other organisms should have a similar fraction of expressed genes in their cells to human cells, this choice of threshold gives similar values of $\langle a \rangle$ for human and mouse cells. This would also be true for small non-zero values of $\hat{x}$, however choosing a non-zero value is harder to justify. The data is initially in the form of transcript counts from RNA-Seq. experiments, and thus naive, but at least justifiable, assumptions to make are that there was little experimental error and each transcript observed or not is directly due to the corresponding gene expression, or any experimental and measurement error has already been corrected for in the data during any preprocessing.

Interestingly, increasing the threshold slightly from zero to $\hat{x} = 0.1$ also removes the strong localisation of the the high activity mouse ENS cells around a single value as in figure 3.2.1 - see figure 3.A.2

**Figure 3.A.2:** The activity of each sample in the mouse ENS data set when a threshold of $\hat{x} = 0.1$ is used to binarize the gene expression levels.

## APPENDIX 3.B    HUMAN PROTEIN ATLAS DATA CELL-CELL SIMILARITIES

The correlations and covariances in the gene expression levels of the different cells taken from the Human Protein Atlas database are plotted in figures 3.B.1 and 3.B.2 respectively. There is less of an obvious structure in the correlation matrix of the human data, with the majority of pairs of cells having a strong correlation $cor[\eta, \eta'] \sim 0.6 - 0.7$, compared to the mouse ENS data whose correlation values are dependent on the cell cycle phase of the sample. Ordering the cell lines by the tissue type that they belong to also does not reveal any interesting patterns.

**Figure 3.B.1:** Correlations between the binary gene expression levels of each pair of cells taken from the Human Protein Atlas database.

**Figure 3.B.2:** Covariances between the binary gene expression levels of each pair of cells taken from the Human Protein Atlas database.

# 4

# Bipartite Gene Regulatory Networks

## 4.1. Introduction

In this chapter, the focus is moved away from the model of hierarchically stored cell cycles and turns to the interactions between genes that drive the expression levels. Along with morphological properties, cell types have long been characterised by the gene expression levels observed in experiments. Furthermore, cell types have recently been shown to be high dimensional attractors in the gene expression levels space [33] underlying their importance. However, this idea is not new and dates back to Waddington's metaphor of an epigenetic landscape [11] - see figure 1.1.1. In his analogy, developmental decisions are described by a ball rolling down a hill into different

valleys. These valleys represent different cell types that decrease in potency down the landscape. Waddington speculated that the shape of the landscape would be dictated by genes anchoring the landscape through their interactions [11]. The exact mechanisms that shaped the landscape were beyond the reach of the scientific methods of Waddington's time. However, the interactions between genes that drive expression levels has now been known for many decades - regulation via transcription factors (TFs). In the years since Waddington's original work, gene regulation has been studied meticulously using experimental and computational techniques [108, 109], such as gene editing [110], reporter genes and assay techniques [111], and gene regulatory network reconstruction [112, 113]. In parallel, there have been many attractor models that have attempted to describe cell fates, including that presented in chapter 2, although, there is still no universally accepted model explaining the mechanism behind the attractors.

Kauffman was the first to study cells as attractors of a dynamical system. He used a boolean network approach, in which the expression level of a gene is a random boolean function of its inputs, which are expression levels of other randomly chosen genes [12]. Whilst Kauffman networks have had some success in improving our understanding, because of the random nature of the interactions, they leave mechanistic details to our imagination. More recently, neural network models have had success in reproducing dynamics similar to those observed in experiments [81, 82, 114]. In these models the attractor structure is encoded in the gene interactions, or protein interactions, using a coupling matrix. However, the interactions in these types of models do not unveil any biological details beyond which genes should be co-expressed. Furthermore, the interactions in these Hopfield-like models are dense and therefore at odds with biological observations.

Models for cell fate decisions and reprogramming typically consist of small specific gene regulatory networks, because of experimental evidence for mutually exclusively expressed genes at branching points in development [34, 115–117] or, because they aim to model transitions between a select number of certain cell types [36, 118, 119]. Informaticians may study larger in-

teraction networks with the advent of high throughput experiments and big data. Furthermore, despite the fact they are deeply connected through gene expression and gene regulation, protein-protein interaction networks and gene regulatory networks (GRNs) are typically studied separately. It will be shown that studying a combined gene-TF network can provide better insight into the underlying molecular biology. Specifically, providing evidence that (i) TFs should be single proteins or small complexes that regulate many genes in order to maintain a steady state gene expression profile; (ii) multiple gene expression level attractors, or cell types, can be supported by a rewiring of the underlying gene regulatory network or specific prescriptions of the regulatory interactions between TFs and genes.

The rest of this chapter is organised as follows: First, a simple bipartite graph model for gene regulation is introduced along with a general model for gene expression dynamics. Next, the role of inhibition on the gene regulatory dynamics and the nature of gene expression level attractors supported by the dynamics is studied. Finally, the main findings are summarised and the implications of this model for future work on gene regulation and attractor models for cell types are discussed. Throughout this chapter, concepts and terminology from network science and graph theory will be used. For a comprehensive review of these ideas and any definitions, in terms of a biological setting, the reader is directed towards the comprehensive review by Pavlopoulos [120].

## 4.2. Proteins, Complexes and Transcription factors

Before the bipartite network formulation for gene regulation is introduced, the interplay between genes, proteins and protein complexes is considered here, to motivate the reduction to a bipartite graph and illustrate the simplicity of that formulation later on. Gene regulation is a complex biochemical feedback process that controls gene expression, and therefore, the concentration of gene products such as mRNA, proteins, and complexes, inside a cell. Gene regulation also allows a cell to react to external signals, such as morphogen gradients and cell stress. Thus, whilst this initial model will already seem complex, it is worth bearing in mind that it is a coarse-grained

version of reality as there are many more players involved in the game of gene regulation, such as multiple types of mRNA, tRNA, miRNA, enhancer regions of DNA, RNA polymerase, ribosomes, external signals, etc. These extra players could be included in the model depending on the level of detail one desires. However, because the system will be reduced down to a bipartite graph later on, they need not be included at this stage.

If it is assumed that every gene in a cell codes for a single protein, that could be a transcription factor (i.e. a protein that activates/inhibits a gene's expression level), and these proteins can also bind together to form protein complexes, that may also be transcription factors, then it is possible to construct a gene regulatory network (GRN) with multiple layers, like in figure 4.2.1, where each layer represents a different component involved in gene regulation. The network is complex, with some interactions, like gene expression, being directed, whilst others, e.g. complex formation/dissociation, may be undirected.

Reaction equations, in continuous time, can be written for the evolution of the concentration of protein molecules,

$$\dot{p}_i = n_i \eta_i - p_i \sum_j p_j \Pi_{ij}^+ + \sum_j c_{ij} \Pi_{ij}^- - \gamma_i p_i \,, \tag{4.1}$$

and complexes formed from 2 proteins,

$$\dot{c}_{ij} = p_i p_j \Pi_{ij}^+ - c_{ij} \Pi_{ij}^- - \gamma_{ij} c_{ij} \,. \tag{4.2}$$

Here $n_i$ is the binary gene expression level of gene $i$, $\boldsymbol{\eta}$ are the rates of protein synthesis, $\Pi^\pm$ are association/dissociation rates for protein/complexes and the $\boldsymbol{\gamma}$ variables are degradation rates. The protein-protein interaction network of a cell can be constructed from the non-zero values of $\Pi_{ij}^+$ that details which proteins interact with one another. Assuming that the dynamics of protein synthesis, dissociation and decay occurs at a much greater rate than that of gene expression, one

**Figure 4.2.1:** A network representation of gene regulation (not all nodes/edges are shown). A gene $n_i$ synthesises a protein $p_i$, which can reversibly bind (undirected edges) to form protein complexes $c_{ij}$. The proteins and complexes may act as transcription factors (TFs) for each of the genes. The activation/inhibition nature of each TF is labelled by $\xi \in \{0, \pm 1\}$ - the superscripts of $\xi$ denote which genes contribute to the transcription factors formation whilst the subscripts denote which gene it regulates.

can apply a separation of time scales resulting in stationarity in $\dot{\mathbf{p}}$ and $\dot{\mathbf{c}}$ for each time step in the gene expression dynamics. At stationarity (4.2) gives,

$$c_{ij} = \frac{\Pi_{ij}^+}{\Pi_{ij}^- + \gamma_{ij}} p_i p_j \,, \tag{4.3}$$

or

$$-p_i p_j \Pi_{ij}^+ + c_{ij} \Pi_{ij}^- = -\gamma_{ij} c_{ij} \,, \tag{4.4}$$

which allows us to simplify $(4.1)$ to

$$\dot{p}_i = n_i \eta_i - \sum_j \gamma_{ij} c_{ij} - \gamma_i p_i \, . \tag{4.5}$$

Now using the stationary solution of $(4.2)$, $\dot{p}_i$ can be written in terms of the protein concentrations only,

$$\dot{p}_i = n_i \eta_i - \left( \sum_j \frac{\gamma_{ij} \Pi_{ij}^+}{\Pi_{ij}^- + \gamma_{ij}} p_j + \gamma_i \right) p_i \, , \tag{4.6}$$

which has a stationary solution at,

$$p_i = \frac{n_i \eta_i}{\sum_j \dfrac{\gamma_{ij} \Pi_{ij}^+}{\Pi_{ij}^- + \gamma_{ij}} p_j + \gamma_i} \, . \tag{4.7}$$

When $n_i = 0$, protein molecules from gene $i$ are not synthesised and $(4.7)$ rightly gives $p_i = 0$. When $n_i = 1$ expanding the right hand side of $(4.7)$ in the limit of small protein concentration (i.e. $p_j = 0$) gives,

$$p_i \simeq \frac{\eta_i}{\gamma_i} \left( 1 - \frac{1}{\gamma_i} \sum_j \frac{\gamma_{ij} \Pi_{ij}^+}{\Pi_{ij}^- + \gamma_{ij}} p_j \right) \, , \tag{4.8}$$

which can be rearranged into the following form,

$$\gamma_i \simeq \sum_j \left[ \frac{\gamma_i^2}{\eta_i} \delta_{ij} + \frac{\gamma_{ij} \Pi_{ij}^+}{\Pi_{ij}^- + \gamma_{ij}} \right] p_j \, , \tag{4.9}$$

where the right hand side can be identified as the result of the multiplication of a matrix with the vector of protein concentrations, i.e. one has the matrix equation

$$\boldsymbol{\gamma} = \mathbf{M} \mathbf{p} \, . \tag{4.10}$$

Where $\mathbf{M}$ is a matrix with symmetric and diagonal parts, $\mathbf{M} = \mathbf{S} + \mathbf{D}$. By inverting this matrix

equation, the the stationary solution $(4.7)$ can be written as,

$$p_i = n_i \sum_j M_{ij}^{-1} \gamma_j \,. \tag{4.11}$$

Then, if the gene expression levels have discrete time dynamics on a slower scale, e.g. stages of the cell cycle, the general update rule for binary gene expression levels $(2.1)$ can be written in terms of the multilayer gene regulatory network parameters,

$$n_i(t+1) = \Theta \left[ \sum_j \xi_i^j p_j + \sum_{j,k} c_{jk} \xi_i^{jk} + \ldots + T z_i - \theta_i \right] \tag{4.12}$$

where $\xi_i^j \in \{0, \pm 1\}$ denotes the regulatory effect of the TFs. Rewriting this using the stationary solutions for the protein and complex dynamics gives,

$$n_i(t+1) = \Theta \left[ \sum_j \xi_i^j \underbrace{\sum_k M_{jk}^{-1} \gamma_k}_{J_{ij}} n_j + \sum_{j,k} \xi_i^{jk} \underbrace{\sum_{l,l'} \frac{\Pi_{jk}^+ \gamma_l \gamma_{l'} M_{jl}^{-1} M_{kl'}^{-1}}{\Pi_{jk}^- + \gamma_{jk}}}_{J_{ijk}} n_j n_k + \ldots \right] \,. \tag{4.13}$$

Note that if the regulatory interactions $\boldsymbol{\xi}$ are sparse, the effective interactions $J_{ij}$ and $J_{ijk}$ between genes are also sparse.

From the form of $J_{ij}$ and $J_{ijk}$, it can be seen that the interaction between two genes require regulation via single proteins and the interaction between three genes requires regulation through a complex of two proteins. Expanding this reasoning to interactions between many more genes, one sees that effective interactions between $m$ genes requires complexes that are transcription factors formed from the expression of $m-1$ genes. This information is important for the attractor nature of cell types in the gene regulatory network - for $300-400$ different cell types to be formed from $\sim 2,500$ regulatory genes (numbers here for humans) one would expect that interactions of higher order than pairs are needed to sustain the robust nature of many attractors in

a small network.

## 4.3. BIPARTITE GRAPH FORMULATION

The intricate dynamics of gene expression, protein synthesis, complex formation and activation & inhibition can be simplified using a directed bipartite graph. From this point on much of the notation used may be different from the previous chapters, however, all notation shall be explained on introduction. In this network representation, gene regulation is modelled using two sets of nodes (or vertices) representing $N$ individual regulatory genes and $\alpha N$ transcription factors respectively (see figure 4.3.2). Only regulatory genes are included in the network, i.e. those that contribute to TF formation, because, the expression of all genes that do not contribute to the synthesis of a TF are driven by this sub-network. Furthermore, because each gene does not typically regulate every other in a cell, the set of bipartite graphs studied are restricted to those with finite connectivity. The edges in the network are used to indicate the interactions between the genome and transcriptome. If the TF $\mu$ regulates the gene $i$, then they are connected by a directed edge $\xi_i^\mu \in \{\pm 1\}$, with the sign representing the activation/inhibition nature of the regulatory interaction. Similarly, if the gene $i$ expresses a protein that contributes to the formation of the TF $\mu$ there is a directed edge $\eta_i^\mu$, from node $i$ to $\mu$. The full connectivity of the gene regulatory network is jointly defined by the matrices $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$. This formulation is somewhat similar to the Wagner gene network model [69]. However, there is a key difference: a gene does not necessarily produce just a single species of transcriptional regulator (i.e. a TF), as assumed in Wagner's model, but may contribute to several different ones through the formation of complexes that contain the protein encoded by the gene in question.

Typically one does not know the structure of the vectors $\boldsymbol{\eta}$ & $\boldsymbol{\xi}$. Thus, they are treated here as independent random variables. The bipartite graph has 4 types of degrees: the TF in-degrees $c_\mu^{\text{in}}(\boldsymbol{\eta}) = \sum_{i=1}^N \eta_i^\mu$ are the number of genes that express a protein that contributes to the formation of the TF $\mu$, i.e. the size of the TF; the TF out-degrees $c_\mu^{\text{out}}(\boldsymbol{\xi}) = \sum_{i=1}^N |\xi_i^\mu|$ are the number

of regulatory genes that the TF $\mu$ regulates, i.e. the number of DNA binding domains on the TF multiplied by the number of occurrences of the response elements for those binding sites in the genome; the gene in-degrees $d_i^{\text{in}}(\boldsymbol{\xi}) = \sum_{\mu=1}^{\alpha N} |\xi_i^\mu|$ are the number of TFs that can regulate the gene; and the gene out-degrees $d_i^{\text{out}}(\boldsymbol{\eta}) = \sum_{\mu=1}^{\alpha N} \eta_i^\mu$ are the number of TFs that a protein expressed by a given gene can contribute to. For a given bipartite gene regulatory network the distributions of $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ can be defined as

$$p(\boldsymbol{\eta}) = \prod_{i,\mu} \left[ \frac{c_\mu^{\text{in}} d_i^{\text{out}}}{N\langle d^{\text{out}} \rangle} \delta_{\eta_i^\mu,1} + \left( 1 - \frac{c_\mu^{\text{in}} d_i^{\text{out}}}{N\langle d^{\text{out}} \rangle} \right) \delta_{\eta_i^\mu,0} \right] , \qquad (4.14)$$

$$p(|\boldsymbol{\xi}|) = \prod_{i,\mu} \left[ \frac{c_\mu^{\text{out}} d_i^{\text{in}}}{N\langle d^{\text{in}} \rangle} \delta_{|\xi_i^\mu|,1} + \left( 1 - \frac{c_\mu^{\text{out}} d_i^{\text{in}}}{N\langle d^{\text{in}} \rangle} \right) \delta_{|\xi_i^\mu|,0} \right] , \qquad (4.15)$$

where $\delta_{x,y}$ is the Kronecker delta function, defined as $1$ for $x = y$ and $0$ otherwise. These distributions for $\boldsymbol{\eta}$ & $\boldsymbol{\xi}$ assume that edges in the network are independent, hence the factorisation over $i$ and $\mu$, and the likelihood that a gene $i$ expresses a TF $\mu$ is governed by the size of the TF and the promiscuity of the protein that the gene expresses - i.e. the number of TFs that a protein expressed from a gene contributes to - and the likelihood that a TF $\mu$ regulates a gene $i$ is determined by the number of DNA binding sites of TF and the gene's binding site. For simplicity, if one assumes that $c_\mu^{\text{in}} = c_1 \; \forall \mu$, $c_\mu^{\text{out}} = c_2 \; \forall \mu$, $d_i^{\text{out}} = d_1 \; \forall i$ and $d_i^{\text{in}} = d_2 \; \forall i$, then, for large networks, $N \gg 1$, this gives Poisson distributions for the in- and out-degrees for both TFs and genes, $c_\mu^{\text{in}}(\boldsymbol{\eta}) \sim \Pi_{c_1}$, $c_\mu^{\text{out}}(\boldsymbol{\xi}) \sim \Pi_{c_2}$, $d_i^{\text{in}}(\boldsymbol{\xi}) \sim \Pi_{d_2}$, and $d_i^{\text{out}}(\boldsymbol{\eta}) \sim \Pi_{d_1}$.

Because genes are only connected to TFs by a directed edge, and vice versa, the following equalities must be true for in- and out-degrees,

$$\sum_{\mu=1}^{\alpha N} c_\mu^{\text{in}}(\boldsymbol{\eta}) = \sum_{i=1}^{N} d_i^{\text{out}}(\boldsymbol{\eta}) \qquad \forall \boldsymbol{\eta} , \qquad (4.16)$$

$$\sum_{\mu=1}^{\alpha N} c_\mu^{\text{out}}(\boldsymbol{\xi}) = \sum_{i=1}^{N} d_i^{\text{in}}(\boldsymbol{\xi}) \qquad \forall \boldsymbol{\xi} \,, \tag{4.17}$$

which leads to the following identities between the average in- and out-degrees,

$$d_1 = \alpha c_1 \,, \tag{4.18}$$

$$d_2 = \alpha c_2 \,. \tag{4.19}$$

Furthermore, because the degree distributions are Poisson, and thus defined completely by their first moment, it is not necessary to know the exact structure of $\boldsymbol{\eta}$ & $\boldsymbol{\xi}$ to study average properties of a gene regulatory network for a given system.

The inclusion of the transcriptome, in this bipartite image of a gene regulatory network, immediately highlights an interesting property of gene regulation: the interaction between $m$ unique genes requires a transcription factor that is a product of $m-1$ expressed genes. Also, any transcription factor that is a single protein molecule will have an in-degree $c_\mu^{\text{in}} = \sum_i \eta_i^\mu = 1$ and is thus synthesised by a single gene. These conclusions were not easily drawn in the previous section and help to highlight the explanatory power of the bipartite gene regulatory network.

Typically, experimental work focuses on the effective interactions between gene expression levels. For example, in knock-out experiments, a gene is silenced and the effect on other gene expression levels is monitored and compared to a control (and often the wild-type). It is possible to reduce the bipartite nature of a network down to a gene-gene interaction network to compare theoretical and numerical results with experiments, by simply integrating out the transcription factors (see Appendix 4.A).

**Figure 4.3.2:** A bipartite graph representation for gene-TF interactions in a system with $N$ genes and $\alpha N$ TFs (not all edges and nodes are shown). Edges $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ represent the directed interactions between genes-TF and TF-genes respectively. The average in- and out-degrees for the genes and TFs are given by $(c_1, c_2)$ and $(d_2, d_1)$ respectively. For conservation of degrees $\alpha c_1 = d_1$ and $\alpha c_2 = d_2$.

## 4.4. REGULATORY DYNAMICS

With the network structure defined in the previous section our attention is now turned to the gene regulatory dynamics. The common simplification that each gene expression level is given by a binary variable $\sigma_i \in \{0, 1\}$, with $i = 1 \ldots N$, is made. A gene expression pattern, or cell type, is then completely defined by the vector of gene expression levels $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \ldots, \sigma_N)$ and how this state changes in time is studied. The assumption of Boolean or "on/off" genes is used purely to simplify the mathematics, however, it may be relaxed if a more comprehensive description of the gene regulatory dynamics is required. The state of the transcriptome is also denoted with a vector of TF concentrations $\boldsymbol{\tau} = (\tau_1, \ldots, \tau_{\alpha N})$ with $\tau_\mu$ indicating the concentration of the TF $\mu$.

A general model for the dynamics of binary gene expression levels has the following form,

$$\sigma_i(t+1) = \Theta \left[ \sum_\mu \tau_\mu b_i^\mu \xi_i^\mu - \theta_i - Tz_i \right], \tag{4.20}$$

i.e. if the combined effect of all transcription factors that regulate a gene are greater than the level of noise in the system and a gene specific threshold $\theta_i$, the gene will be expressed in the next time

step, $\sigma_i(t+1) = 1$. Here, $z_i$ is a random variable that represents the fundamental stochastic nature of gene regulation; $T$ is a parameter that scales the strength of that noise level; and $b_i^\mu$ is the binding affinity of TF $\mu$ to its target gene $i$. The gene specific thresholds $\theta_i$ represent a barrier for which the regulatory interactions need to overcome to activate a gene and could describe, for example, the level of accessibility of a genes promoter site resulting from the chromatin structure. In order for a TF to be synthesised from a set of expressed genes, it is required that all of the genes that have an out-degree connected to that TF must be expressed at the same time. Thus, an order parameter $m_\mu(t)$, is introduced to keep track of this condition for each TF:

$$m_\mu(t) = \frac{\sum_i \eta_i^\mu \sigma_i(t)}{\sum_i \eta_i^\mu} \, ,$$ (4.21)

which takes a values $m_\mu(t) \in [0, 1]$. If none of the genes contributing to TF $\mu$ are expressed then $m_\mu(t) = 0$ and if they are all expressed $m_\mu(t) = 1$. The concentration of each TF then evolves according to the differential equation

$$\dot{\tau}_\mu = \pi_\mu^+ \delta_{m_\mu(t),1} - \pi_\mu^- \tau_\mu \, .$$ (4.22)

where $\pi_\mu^\pm$ are the production/degradation rates for the transcription factor $\mu$.

Assuming that the TF dynamics are much faster than the changes in gene expression levels allows for the application of a separation of time scales. This results in the TF concentration reaching a steady state before genes are regulated. Stationarity of the TF concentration is given by $\dot{\tau} = 0$, giving the steady state concentration

$$\tau_\mu = \frac{\pi_\mu^+}{\pi_\mu^-} \delta_{m_\mu(t),1} \, .$$ (4.23)

Substituting this steady state concentration into the dynamics results in the following evolution

of the gene expression levels,

$$\sigma_i(t+1) = \Theta \left[ \sum_\mu \xi_i^\mu \delta_{m_{\mu(t)},1} - \theta_i - Tz_i \right] ,$$

(4.24)

where the ratio of production to degradation rates and the binding affinities have been set to $\frac{\pi_\mu^+}{\pi_\mu^-} = b_i^\mu = 1$ for all TFs purely for simplicity. Because each gene expression level depends on all others, at every time step, through the order parameter $m_\mu(t)$, the dynamics for the gene expression levels are not only highly non-linear but they are coupled as well. Even when the dynamics is "linearised" by replacing the $\delta_{m_{\mu(t)},1}$ with $m_{\mu(t)}$, the dynamics remains coupled and non-trivial. However, this scenario is instructive as it will lead to a bound on the dynamics. The linearised version of the dynamics has Hebbian-type interactions but is both asymmetric and diluted:

$$\sigma_i(t+1) = \Theta \left[ \sum_j \sum_\mu \frac{\xi_i^\mu \eta_j^\mu}{c_\mu^{\text{in}}} \sigma_j(t) - \theta_i - Tz_i \right] .$$

(4.25)

Thus, in models with effective interactions of gene expression levels [82, 114] this dynamics highlights a structure for the interactions with the couplings $J_{ij} = \sum_\mu \frac{\xi_i^\mu \eta_i^\mu}{c_\mu^{in}}$.

The equations (4.24) & (4.25) shall be referred to as the "non-linear" and "linear" versions of the gene expression dynamics from now on. In both versions of the dynamics, it is possible to separate the interaction terms into a signal from a TF $\mu$ and the interference of all other TFs on a gene expression level. To see this, set $\sigma_i(t) = \eta_i^\mu \ \forall i$, giving $m_\mu(t) = 1$, and the dynamics evolves as $\sigma_i(t+1) = \Theta \left[ \xi_i^\mu + \sum_{\nu \neq \mu} \xi_i^\nu - \theta_i - Tz_i \right]$ - i.e. a gene will be activated or inhibited by a TF $\mu$ depending only on its regulatory nature, whilst all other TFs act (along with the noise in the system) to interfere with that signal. To simplify matters for the rest of this chapter, unless stated otherwise the noise and thresholds will remain fixed with $T = \theta_i = 0$ to reduce the number of parameters explored. However, they should not be forgotten as they play important roles in the dynamics of the model. For example, the gene regulatory dynamics is deterministic for $T = 0$
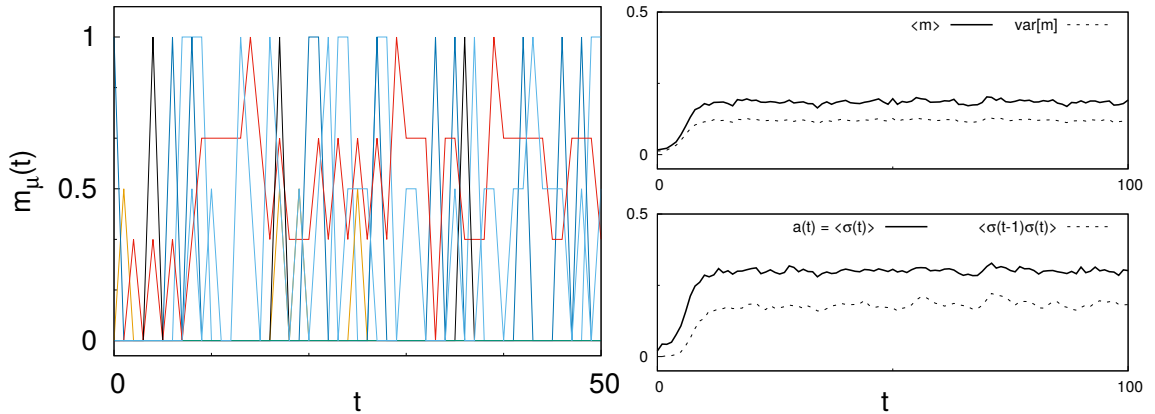
**Figure 4.4.3:** Sample trajectories of the order parameter for 10 randomly selected TFs in a simulation using the non-linear dynamics on bipartite GRN with Poisson degree statistics in a system with $N = \alpha N = 2,500$, $T = 0$ and $\theta_i = 0 \; \forall i$ (left). The average and variance of $m_\mu(t)$ (top right) and activity and autocorrelation for the gene expression levels (bottom right) throughout that same simulation. Each transcription factor in the network were chosen to be either an activator or inhibitor with equal probability.

but many biological processes are known to be stochastic.

Figure 4.4.3 shows the evolution of the order parameter $m_\mu(t)$, for a set of randomly selected TFs from a simulation using the non-linear dynamics (4.24), along with the average of the order parameter over all the TFs in that simulation and the activity, $a(t) = \langle \sigma(t) \rangle$, over a longer time window. In this example the dynamics is deterministic ($T = 0$), the degree distributions are Poisson ($c_1 = 1$ and $c_2 = 10$), and the TFs are activators or inhibitors with equal probability. Because of the absence of noise, the fluctuations seen in $m_\mu(t)$ and gene expression levels are governed solely by the regulatory interactions. In the sample trajectories plotted, there are many times at which a TF is synthesised, i.e. $m_\mu(t) = 1$. However for these TFs, genes contributing to them are silenced within a few time steps of its synthesis and the order parameter drops. This inhibition could be due to one of the finite fraction of TFs synthesised at each time step ($\langle m \rangle = \frac{1}{\alpha N} \sum_\mu m_\mu \approx 0.2$ in the steady state) or that the synthesised TF itself has a negative feedback loop with the gene(s) that express it. Furthermore, this simulation was set up such that all genes are initially not expressed $\sigma(t = 0) = 0$ and there exist just a small number of TFs (10) in the system. Despite the low initial number of TFs, the system quickly evolves to small fluctuations

about a steady state with approximately $1/3$ of the genes being expressed at every time step.



**Figure 4.4.4:** Average gene expression levels, $a(t)$, and fraction of transcription factor species, $a_{\mathrm{TF}}(t)$, observed in simulations of bipartite gene regulatory networks, with the non-linear (left column) and linear dynamics (right column), for both ferromagnetic (activating TFs only) and spin glass (with activators and inhibitors occurring with equal probability) regulatory interactions. The dynamics were simulated on different bipartite gene regulatory networks with Poisson degree distributions and $N = 2,500$, $\alpha = 1$, $c_1 = 1$, $T = 0$ and $\theta_i = 0$ $\forall i$. Each curve represents the average over 100 networks with the same connectivities $c_2$, ranging from $c_2 = 1, 2 \ldots 7$ (FM) and $c_2 = 1, 2 \ldots 20$ (SG). Each network reaches a steady state from an initially silenced configuration when a small number (10) TFs are introduced.

This behaviour is typical of the gene regulatory dynamics on bipartite networks with Poisson degree distributions (figure 4.4.4). The density of expressed genes, $a(t)$, and synthesised TFs, $a_{\mathrm{TF}}(t) = (\alpha N)^{-1} \sum_\mu \bar{\tau}_\mu(t)$ where $\bar{\tau}_\mu(t) = \Theta[\tau_\mu(t)]$, will evolve to reach a steady state that

depends on the interplay between activators and inhibitors; the initial conditions; and the structure of the network; when a small number of TFs are introduced to a network. When a network consists of only activators (i.e. the system has only ferromagnetic interactions) the system approaches a steady state in which approximately all the genes that are connected to the originally introduced TFs (by possibly several gene expression and regulatory steps) are expressed. However, if both activating and inhibiting regulatory interactions exist in the network (i.e. spin glass interactions), the dynamics plateaus to a lower net activity with fluctuations occurring around the steady state, due to competing TFs.

For a given network connectivity, the linear dynamics reaches a steady state in a fewer number of time steps than the non-linear dynamics. This is due to the constraint that requires all genes contributing to a TF to be activated simultaneously for that TF to be synthesised in (4.24). This is also reflected in the reduction in the fraction of unique TFs observed in simulations for the non-linear dynamics. However, for both linear and non linear dynamics; entirely activating and activator-inhibitor networks; there are certain network connectivities that remain in a silent state, $\sigma(t) = 0$, indefinitely even in the absence of noise (figure 4.4.4). Thus, in the next section, our attention turns to how the gene regulatory dynamics depends on the underlying structure of the bipartite network.

### 4.4.1. Percolation thresholds

A key property of any cell is its ability to maintain a gene expression pattern, allowing it to perform a specific function. Without sufficiently many expressed genes a cell would be unable to sustain itself and any biological function. A steady state gene expression level can be supported in two possible ways: (i) There exist many small disconnected clusters in the network, that do not interact, and the gene expression profile is the result of the sum of the activities in each cluster; (ii) The network is highly connected and the regulatory interactions give rise to a stable gene expression profile over the entire network. It is likely that the latter case is true. Cell reprogram-

ming experiments have demonstrated that nearly any adult cell can be transformed into a stem cell like state by introducing a small set of TFs (now known as the Yamanaka factors). Therefore, it is likely the targets of the Yamanaka factors are hubs (or are closely connected to hubs) in a large connected component. Furthermore, the steady state of gene expression of a cell must even be maintained in the early stages of development, before the maternal-zygotic transition when an embryo only translates maternally inherited mRNA. Across the transition, it is believed that two processes must occur: (i) an increase in zygotic gene expression and (ii) degradation of the maternal mRNA. For the former process to occur, a number of transcription factors, translated from maternal mRNA, must kick-start the zygotic gene expression dynamics of the GRN (which is even present before the transition).

This kind of phenomena can be studied by thinking in terms of percolation theory - which has been studied for directed random graphs [121] and undirected (scale free) bipartite graphs [122]. If a small number of transcription factors are introduced into an inactive GRN, the cell will sustain a non-trivial steady state level of gene expression only if a giant-cluster exists in the network. Consider, for example, a network with only activating TFs: if all gene expression levels in a cell are independent and the noise in the system activates and inhibits a gene with equal probability, then in the steady state one would expect that half of all genes are expressed. Thus, in order for a cell to sustain some non-trivial ordered state of gene expression levels the ergodicity in the system must be broken. Therefore, there must be a giant component in the underlying gene regulatory network. Hence, the percolation threshold of a bipartite GRN, i.e. the percolation threshold of a random directed bipartite graph, is calculated here using an adaptation of the cavity method [123]. To this end, the indicator variables $n_i$ and $n_\mu$ are introduced, and take value 1 if the gene $i$ or TF $\mu$ belong to the (out component of the) giant cluster, respectively, and are zero otherwise. If one assumes that all transcription factors are activators, i.e. $\xi_i^\mu \in \{0, 1\}$, these indicator variables can be written in terms of the corresponding indicator variables for their

neighbouring nodes in the gene regulatory network,

$$n_i = 1 - \prod_{\mu \in \partial_i^\xi} \left(1 - n_\mu^{(i)}\right) , \qquad (4.26)$$

$$n_\mu = \prod_{j \in \partial_\mu^\eta} n_j^{(\mu)} , \qquad (4.27)$$

where $\partial_i^\xi = \{\mu : \xi_i^\mu = 1\}$ represents the nodes that are the nearest neighbours of gene $i$ that are connected to it via one of its in-degrees, i.e. a $\xi_i^\mu$ edge, and $n_\mu^{(i)}$ is the indicator variable for TF $\mu$ in the cavity graph - a network with gene $i$ and all the edges connecting to it are removed. Similarly, $\partial_\mu^\eta$ and $n_j^{(\mu)}$ are the nearest neighbours of TF $\mu$ connected to one of its in-degrees, i.e. via an $\eta_j^\mu$ edge, and the indicator variable for gene $j$ on the cavity graph with TF $\mu$ removed respectively (see figure 4.4.5). These equations are constructed from the logic of the non-linear dynamics (4.24). In (4.27), the transcription factor $\mu$ belongs to the giant cluster if and only if all the genes contributing to its synthesis are on the giant cluster. Whereas, from (4.26), a gene $i$ belongs to the giant cluster if at least one of the TFs regulating it is also part of the giant cluster. This describes a bootstrap process on directed bipartite graphs. Bootstrap percolation [124] has been well studied on lattices [125–128], regular graphs [129, 130], trees [131], and complex networks [132, 133]. However, there are no exact results for bipartite graphs.

Similarly for the nearest-neighbours of $i$ and $\mu$, one has the cavity equations,

$$n_i^{(\mu)} = 1 - \prod_{\nu \in \partial_{i \backslash \mu}^\xi} \left(1 - n_\nu^{(i)}\right) , \qquad (4.28)$$

$$n_\mu^{(i)} = \prod_{k \in \partial_{\mu \backslash i}^\eta} n_k^{(\mu)} \qquad (4.29)$$

where the notation $\partial_{i \backslash \mu}$ is used to denote the set of nearest neighbours to $i$ excluding the node $\mu$.

**Figure 4.4.5:** Cavity graph (right) for the bipartite GRN (left) with the transcription factor $\mu$ and all its in-degrees removed. The separate branches of the network become independent if the graph is locally tree-like. The schematic shows only a sub-network and does not include self-regulatory interactions for clarity.

These equations are exact on tree-like graphs, and in the thermodynamic limit, will also be exact on graphs sampled from our ensemble, which are locally tree-like because of the sparsity of $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$. It is easy to see that one can continue constructing these equations for the neighbours of the $\mu$ and $i$, their neighbour's neighbours, and so on. In fact in the infinite system limit these equations become a set of stochastic recursion relations. In this limit, one can average (4.29) & (4.28) over the graph ensemble, and assuming that the edges $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ are independent, the following system of equations for the probabilities of being on the giant cluster and cavity probabilities are obtained,

$$g = \langle n_i \rangle = \sum_{d^{\text{in}}=1}^{\infty} P(d^{\text{in}}) \left[ 1 - (1 - \tilde{g})^{d^{\text{in}}} \right] , \tag{4.30}$$

$$\bar{g} = \langle n_\mu \rangle = \sum_{c^{\text{in}}=1}^{\infty} P(c^{\text{in}}) \hat{g}^{c^{\text{in}}} , \tag{4.31}$$

$$\hat{g} = \langle n_i^{(\mu)} \rangle = \sum_{d^{\text{in}}=1}^{\infty} P(d^{\text{in}}) \left[ 1 - (1 - \tilde{g})^{d^{\text{in}}} \right] , \tag{4.32}$$

$$\tilde{g} = \langle n_\mu^{(i)} \rangle = \sum_{c^{\text{in}}=1}^{\infty} P(c^{\text{in}}) \hat{g}^{c^{\text{in}}} . \tag{4.33}$$

The cavity probabilities $\hat{g}$ and $\tilde{g}$ are the probabilities that random edges from TFs to genes terminate at genes in the giant cluster, and similarly that random edges from genes to TFs lead to TFs in the giant cluster. That is the probability that a gene or TF will belong to a giant-cluster when one of its nearest neighbours, specifically successors, is removed from the graph. Here, the probabilities for a node to be on the giant cluster are equal to the corresponding cavity probabilities, $g = \hat{g}$ and $\bar{g} = \tilde{g}$, this is a result of the directedness, sparsity and independence of $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ as shown in appendix 4.B. The probabilities $P(d^{\text{out}})$ and $P(c^{\text{in}})$ are the gene and TF in-degree distributions respectively. Assuming genes and transcription factors have Poissonian degree distributions, the stability analysis of this system of equations gives rise to the following critical average transcription factor out-degree (see Appendix 4.B),

$$c_2^* = \frac{e^{c_1}}{\alpha c_1} \,, \tag{4.34}$$

above this critical value of $c_2$ a giant-cluster will exist in the network. This is a percolation threshold for a bipartite gene regulatory network with the non-linear regulatory dynamics (4.24).

One can construct a similar argument for the linear version of the dynamics (4.25) by relaxing the constraint in (4.27). Using the logic of the linear dynamics (4.25), a gene will be part of the giant cluster if it contributes to at least one TF on the giant cluster. Then, as before, a TF will belong to the giant cluster if it regulates at least one gene on the giant cluster. With this reasoning one obtains the following expressions for the indicator variables,

$$n_i = 1 - \prod_{\mu \in \partial_i^\xi} \left(1 - n_\mu^{(i)}\right) \,, \tag{4.35}$$

$$n_\mu = 1 - \prod_{j \in \partial_\mu^\eta} \left(1 - n_j^{(\mu)}\right) \,, \tag{4.36}$$

$$n_i^{(\mu)} = 1 - \prod_{\nu \in \partial_{i \backslash \mu}^{\xi}} \left(1 - n_\nu^{(i)}\right) , \tag{4.37}$$

$$n_\mu^{(i)} = 1 - \prod_{j \in \partial_{\mu \backslash i}^{\eta}} \left(1 - n_j^{(\mu)}\right) , \tag{4.38}$$

and similarly for the probabilities

$$g = \sum_{d^{\mathrm{in}}=1}^{\infty} P(d^{\mathrm{in}}) \left[1 - (1 - \tilde{g})^{d^{\mathrm{in}}}\right] , \tag{4.39}$$

$$\bar{g} = \sum_{c^{\mathrm{in}}=1}^{\infty} P(c^{\mathrm{in}}) \left[1 - (1 - \hat{g})^{c^{\mathrm{in}}}\right] , \tag{4.40}$$

$$\hat{g} = \sum_{d^{\mathrm{in}}=1}^{\infty} P(d^{\mathrm{in}}) \left[1 - (1 - \tilde{g})^{d^{\mathrm{in}}}\right] , \tag{4.41}$$

$$\tilde{g} = \sum_{c^{\mathrm{in}}=1}^{\infty} P(c^{\mathrm{in}}) \left[1 - (1 - \hat{g})^{c^{\mathrm{in}}}\right] . \tag{4.42}$$

The stability analysis of these probabilities, for Poisson degree distributions, result in the simpler form for $c_2^*$ (see appendix 4.B),

$$c_2^* = \frac{1}{\alpha c_1} . \tag{4.43}$$

This generalises the result obtained for undirected bipartite graphs, where $\boldsymbol{\eta} = \boldsymbol{\xi}$ and $c_1 = c_2$ [134]. Thus, this result could also have been achieved by marginalising the bipartite gene regulatory network over the TFs to obtain an effective gene-gene interaction network with average degree $\alpha c_1 c_2$ (as demonstrated in appendix 4.A), and then applying known results from directed graphs [121].

The percolation thresholds, $c_2^*$, for both the linear and non-linear dynamics are plotted in figure 4.4.6. In the non-linear dynamics, the exponential dominates, resulting in a scenario in which it
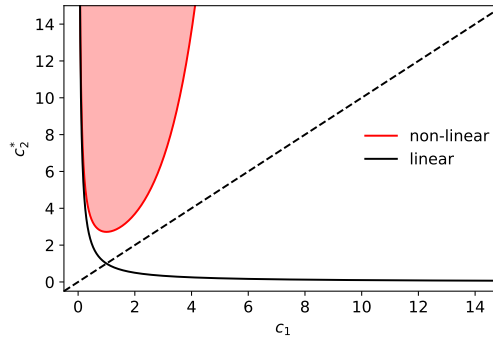
**Figure 4.4.6:** The critical average transcription factor out-degree, $c_2^*$, at which a giant component exists in the GRN as a function of the average transcription factor in-degree $c_1$ (for $\alpha = 1$). Solid lines represent the value of $c_2^*$ for a given $c_1$. In the non-linear dynamics a giant-component can only exist in a subset (shaded red) of the region above $c_2 = c_1$ (dashed line) due to the exponential in (4.34). Whereas, it is possible for a giant-cluster to be present in the network, for the linear version of the dynamics with $c_1 > c_2$.

is not possible to have a giant cluster in the network if $c_1 > c_2$ (for $\alpha = 1$). This is in line with the current understanding of molecular biology [100]. Transcription factors tend to be "promiscuous", regulating many genes, but are also simple complexes made up of fewer proteins in comparison. Contrary, for the linearised dynamics, it is possible for a giant-cluster to exist in the region $c_1 > c_2$.

Simulations of $N = 2,500$ genes and their transcription factors reveal a phase transition in both the steady state fraction of transcription factors synthesised $\langle a_{\text{TF}} \rangle$ and the steady state fraction of expressed genes $\langle a_{ss} \rangle$ (figure 4.4.7). These transitions occur at the values of $c_2^*$ predicted by (4.34) & (4.43). The simulations were performed by introducing a small fraction of transcription factors ($\sim 10$) to a GRN, with Poisson degree distributions, in which initially no other TFs exist and all genes are not expressed. It is worth noting that, because the networks generated for simulations were constructed using Poisson degree distributions, there are more nodes simulated than are actually a part of the Bipartite GRN. A gene is a regulatory gene if it is regulated by at least one TF and contributes to the synthesis of at least one TF, i.e. if $d_i^{\text{in}} \geq 1$ and $d_i^{\text{out}} \geq 1$. Similarly, for a protein or complex to truly be a TF it must be synthesised by at least one regulatory gene,

i.e $c_\mu^{\text{in}} \geq 1$, and then can regulate the expression any gene in the regulatory part of the network or the genes outside of this sub-network $c_\mu^{\text{out}} \geq 0$. In figure 4.4.7, both the analytic solutions and the simulation results take this into account and are normalised by the appropriate probabilities. For example $\bar{g}$ and $\langle a_{\text{TF}} \rangle$ are normalised by $P(c^{\text{in}} > 0)$.

Below the percolation threshold a gene regulatory network is strongly disconnected, and thus, introducing a small set of TFs will only activate a small number of genes. Above the percolation thresholds, introducing a small number of TFs results in an activation "avalanche", due to the presence of a giant-cluster in the network. If transcription factors have a large out-degree then they are able to activate many genes, which, in turn, increases the likelihood of new transcription factors being synthesised. Sufficiently far above the percolation threshold, the giant-cluster encompasses the majority, or entirety, of the network. So introducing a small number of TFs will result in the system equilibrating to a steady state in which all genes are expressed. Then, depending on the average transcription factor in-degree, $c_1$, the fraction of TFs that are synthesised in the steady states will vary. TFs with low in-degree are more likely to be activated for the non-linear dynamics, whilst they are less likely to be activated for the linear dynamics. All of these simulations were carried out deterministically $(T = 0)$ and with no gene-specific thresholds $\theta_i = 0$ $\forall i$. Introducing noise or gene-specific thresholds would alter $\langle a_{ss} \rangle$ and $\langle a_{\text{TF}} \rangle$. The simulations strongly agree with numerical solutions of $g$ and $\bar{g}$ for both the linear and non-linear dynamics.

It is worth reminding the reader that up until this point in our discussion only *activating* TFs have been included in the dynamics. The dynamics become much richer when inhibition plays a role too.

### 4.4.2. The effects of inhibition

Now consider how the dynamics are affected when TFs act as either *activators or inhibitors* with equal probability, i.e. spin glass interactions are used. The critical TF out-degree $c_2^*$ above which a net steady-state gene expression profile is supported by a network is, no longer the percolation
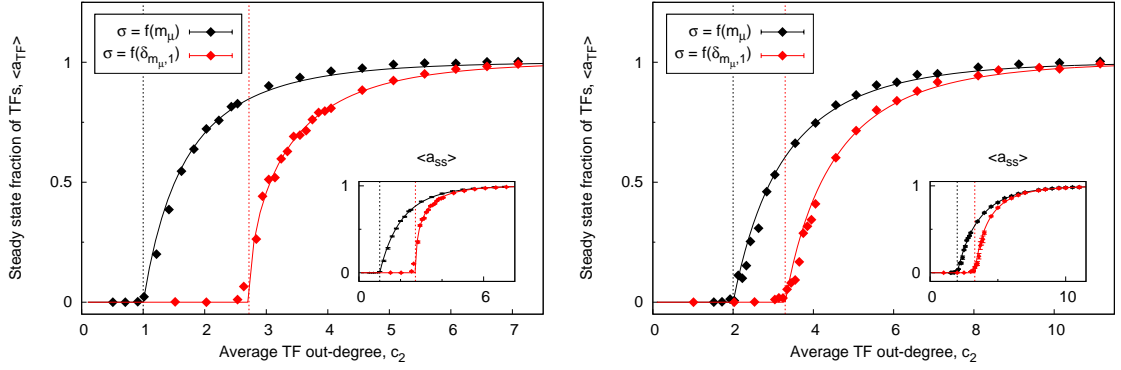
**Figure 4.4.7:** Steady state fraction of transcription factors for different average transcription factor out-degrees, $c_2$, in a bipartite GRN constructed with Poisson degree distributions, using the non-linear (black) and linear (red) versions of the gene regulatory dynamics - Left: $c_1 = 1.0$ and right: $c_1 = 0.5$. Insets show the average steady state gene expression levels for the same simulations. Average of over simulations (data points) were carried out with fixed $N = \alpha N = 2,500$, $\theta_i = 0$ $\forall i$ and $T = 0$. Curves are the probability that gene or TF belongs to the giant cluster of the GRN and are analytical solutions of (4.30), (4.31), (4.39) & (4.40) with $\alpha = 1$. The vertical dashed lines are the percolation thresholds predicted by (4.34) & (4.43).

threshold of the network, and cannot be calculated analytically. However, the dynamics can still be simulated allowing for a numerical study. Frustration in the system will result in more complex dynamics as TFs compete to regulate the same genes in a different manner. Thus, unlike in the activation only networks, it may be possible for a given network to support multiple stable gene expression level states, i.e. different $\boldsymbol{\sigma}$ without any rewiring of the GRN. Whilst histone modifications and chromatin markers could effectively rewire a GRN, the possibility of having multiple steady states in a gene regulation model due to just the interplay between activation and inhibition is a desirable one. When constructing a bipartite GRN, TFs will now activate or inhibit their target genes with a probability $P(\xi = 1) = \epsilon$ and $P(\xi = -1) = 1 - \epsilon$. As before, once the network is constructed its structure will remain fixed, i.e. $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$ are a form of quenched disorder.

In figure 4.4.8 (left panel), the fraction of transcription factors in the steady state is shown for simulations in which TFs are activators or inhibitors with equal probability, i.e. $\epsilon = 0.5$. Although these simulations are with the same number of genes ($N = 2,500$), and have the

same in- and out-degree statistics as those in figure 4.4.7, there is a notable decrease in both the densities of synthesised TFs and expressed genes in the steady state. This is a direct result of the presence of inhibitors. Unsurprisingly the value $c_2^*$ above which a non-zero $\langle a_{ss} \rangle$ and $\langle a_{\text{TF}} \rangle$ exist is also increased when inhibition is involved in the dynamics. This is because inhibitors act to silence genes even if they are part of a giant connected component of the GRN. Thus, above the percolation thresholds (4.34) & (4.43) a giant cluster will still exist, but it's existence is no longer sufficient to support a finite $\langle a_{ss} \rangle$ & $\langle a_{\text{TF}} \rangle$. Again, introducing noise or gene-specific thresholds to the dynamics would further alter the steady states of the system.

To determine whether or not multiple attractors (i.e. stable gene expression level patterns) exist for a given network, the distribution of overlaps (i.e. Pearson correlation coefficients) $q_{\alpha\beta} \in [-1, 1]$, of the steady state TF trajectories is studied after running simulations of the dynamics on the same network with different initial conditions. The difference in initial conditions is which small set of TFs (now 4 to parallel with reprogramming experiments) are introduced to the initially dormant gene regulatory network. The overlap is a measure of similarity between the steady states of two simulation runs $\alpha$ and $\beta$ - its formal definition is given in Appendix 4.C. The overlap $q_{\alpha\beta} = 1$ if the TFs synthesised in the steady state of two simulations are identical. Thus, the distribution of $q_{\alpha\beta}$ is studied to investigate whether or not multiple cell types can emerge from a given genetic network.

The overlap distribution $P(q_{\alpha\beta})$ (right panel of figure 4.4.8) was produced for an arbitrary point in the parameter space above $c_2^*$, such that the network will have non-zero $\langle a_{\text{TF}} \rangle$, using the non-linear (bottom) and linear (top) versions of the dynamics. The distribution does not show the self-overlaps ($q_{\alpha\alpha} = 1$), in order to focus on overlaps between different simulation runs. For both the linear and non-linear dynamics the distribution of overlaps has a single peak at $q_{\alpha\beta} = 1$. This implies, that for $\epsilon = 0.5$, each network supports only a single attractor regardless of the choice of dynamics.

Looking closer at the structure of the probability density function (pdf) of the overlap, near
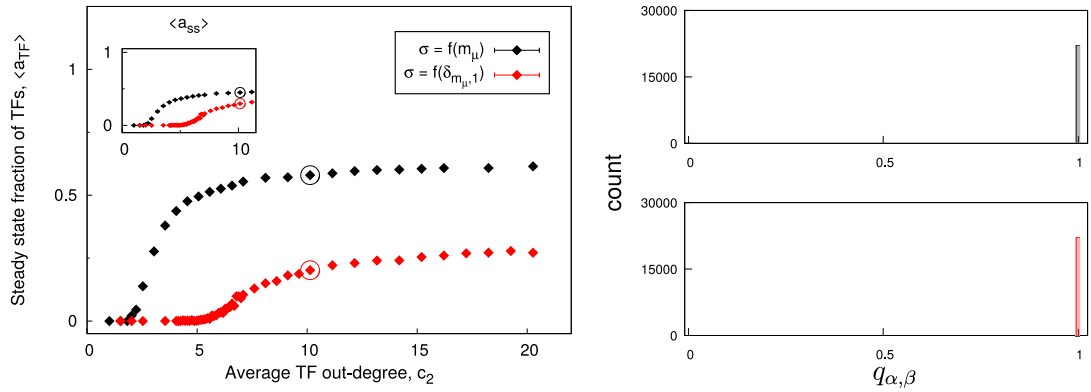
**Figure 4.4.8:** Effect of inhibition on steady state regulatory dynamics. Left - average fraction of TFs in the steady state, for different TF out-degrees, with $N = 2,500$, $\alpha = 1$, $c_1 = 1.0$ and $T = 0$, on a bipartite GRN with Poissonian degree distributions. The histograms of the overlap $q_{\alpha\beta}$ between $150$ simulation runs using the linear (top right) and non-linear (bottom right) dynamics for the points highlighted in the parameter space of the left panel. Self-overlaps $q_{\alpha\alpha}$ are not plotted.

$q_{\alpha\beta} = 1$, shows that the distribution is not a perfect $\delta$-function at exactly $q_{\alpha\beta} = 1$. However, increasing the time window over which the steady state average is performed moves the mass of the pdf towards $q_{\alpha\beta} = 1$, suggesting that there could be a single limit cycle attractor for the dynamics that either has a long period or a short period that is traversed many times (figure 4.4.9). In fact, it can be shown from the average deviations of values of the overlap from $q_{\alpha\beta} = 1$ that there is likely a single limit cycle with a short period (see appendix 4.D). This is supported by a recent study on the effects of dilution on the attractors of asymmetric neural network models, in which the effect of the sparsity and level of asymmetry of the interactions on the mean number and mean length of limit cycle attractors was explored numerically [135]. The authors suggest that increasing the asymmetry of interactions dramatically decreases the number of limit cycle attractors, whilst increasing the sparsity of the interactions decreases the length of the limit cycles. It is also shown that there is a dramatic decrease in the mean number of attractors as the interactions become very sparse. However, at high levels of dilution, this may be due to a lack of a giant component in the networks. The neural network models studied in [135] are analogous to the linear dynamics on a gene-gene network where the TFs have been integrated out of a bipartite

GRN.

One might believe that the existence of a single attractor is due to the lack of short loops in the effective gene-gene network for the choice of $(c_1, c_2) = (1, 10)$. However, when $c_1$ and $c_2$ are increased, with $c_2 \gg c_2^*$, to create a higher connectivity between genes (through TFs) a network still only supports a single attractor. Increasing the average in-degree of a TF, $c_1$, decreases the likelihood that each TF will be synthesised using the non-linear dynamics for any choice of $\epsilon$. This is because more genes need to be expressed simultaneously to produce a TF for increasing $c_1$. For the linear dynamics the opposite is true, increasing $c_1$ increases the probability a TF is expressed in the steady state. This is because, for the linear dynamics, TF synthesis requires at least one of its gene to be expressed. This behaviour can be seen in the top panel of figure 4.4.10 where empirical cumulative distribution functions (CDF) of the steady state frequency of TF synthesis, $\langle \bar{\tau}_\mu \rangle = \frac{1}{\Delta t} \sum_{t=t'}^{t'+\Delta t} \bar{\tau}_\mu(t)$ are plotted.

The average TF out-degree $c_2$ is equivalent to the average gene in-degree $d_2$ for $\alpha = 1$. Therefore, increasing $c_2$ increases the connectivity of the gene-gene network and the amount of TFs that compete to regulate each gene. Not only did increasing $c_2$ not effect the number of attractors observed, it also had no affect on $\langle \bar{\tau}_\mu \rangle$ (figure 4.4.10), because the network remains sparse compared to the total possible number of edges. Therefore, increasing $c_1$ and $c_2$ does not increase the number of attractors supported by a given network for either choice of the dynamics. Hence for $\epsilon = 0.5$ each network supports only a single attractor similar to the activation only networks. The difference in the attractors for different connectivities is in the frequency with which genes are expressed, and TFs are synthesised, in the steady state.

Increasing $\epsilon$ such that there is a bias in the network towards activation also does not affect the number of attractors supported for a given network. Instead, for fixed $(c_1, c_2)$ with a giant cluster in the network $(c_2 \gg c_2^*)$, as $\epsilon$ increases the ratio of the number of activating to inhibiting TFs that regulate a given gene increases on average. Thus, as $\epsilon$ increases the behaviour of the dynamics

**Figure 4.4.9:** The probability density function of the overlap between the same 150 simulations of the non-linear deterministic gene regulatory dynamics on a fixed network with $c_1 = 1$, $c_2 = 10$, $\epsilon = 0.5$. The different distributions arise from different time windows over which the steady state dynamics was averaged: $\Delta t = 500$ (top), $\Delta t = 2,000$ (middle), $\Delta t = 4,500$ (bottom)

**Figure 4.4.10:** Empirical cumulative distribution functions (CDFs) of the steady state frequency that a TF is synthesised in bipartite networks with different connectivities, using the linear (left column) and non-linear (right column) dynamics (at $T = 0$). The networks had fixed statistics of $c_2 = 100$, $\epsilon = 0.5$ (top row) and $c_1 = 4$, $\epsilon = 0.5$ (bottom row). The average TF in-degree dictates the number of TFs that are always synthesised, whilst the average TF out-degree $c_2$ has no significant effect on $F(\langle \tau_{mu} \rangle)$.

**Figure 4.4.11:** Empirical CDF of the steady state frequency with which TFs are synthesised, using the linear (left) and non-linear (right) dynamics on a network with $(c_1, c_2) = (1, 10)$ at $T = 0$. As the bias towards activation increases the number of TFs which are always expressed increases.

tends towards the behaviour of the FM interactions (figure 4.4.11). Increasing $\epsilon$ increases the frequency with which each gene is expressed, and in turn, each TF is synthesised in the steady state - resulting in an increase in $\langle a_{ss} \rangle$ and $\langle a_{\mathrm{TF}} \rangle$, not in the number of attractors for a given network.

The only instance in which a multiplicity of attractors is observed is when the graph has a symmetric structure, i.e. $\eta_i^\mu = \xi_i^\mu$. Although, this graph is no longer a description of a bipartite gene regulatory network, because, $\eta_i^\mu = -1$ has no biological meaning. Instead, the graph becomes a dilute Hopfield network - a type of dilute recurrent neural network model with symmetric interactions - when the linear dynamics are used. These graphs have a multiplicity of attractors, which can be seen from the bimodal overlap distributions in figure 4.4.12. Interestingly, both the linear and non-linear dynamics on dilute symmetric networks support multiple attractors. For the linear dynamics, there are many attractors with a high correlation and the system occasionally finds the same attractor. On the other hand, when the non-linear dynamics is used only $1$ out of $150$ simulations found a different attractor, which has a low correlation with all other simulation runs. This implies that the basin of attraction of a single attractor dominates the phase space of the dynamics, because of the harder constraint in the non-linear dynamics. Hence, whilst the robustness of attractors is amplified by the non-linear dynamics, the asymmetry of the interactions

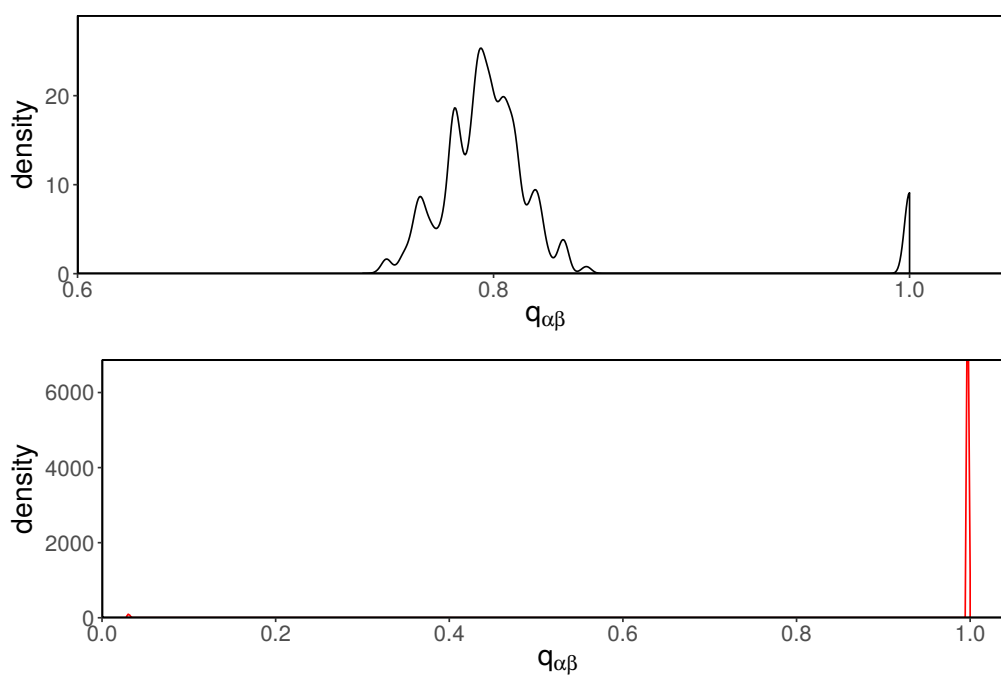**Figure 4.4.12:** Probability density functions for the overlap between simulations on a dilute symmetric network, with average degree $\langle c \rangle = \sqrt{10}$ and $\epsilon = 0.5$, using the linear (top) and non-linear (bottom) deterministic dynamics from the bipartite gene regulatory network model. The diversity of attractors is washed out by the non-linear dynamics.

in the bipartite gene regulatory network all but ensures it.

Therefore, provided a giant cluster exists in the bipartite gene regulatory network, a single steady state gene expression level profile exists for that network. The connectivity of the network, and the ratio of activators to inhibitors in the network, govern the frequency with which genes are expressed/TFs are synthesised in the steady state (for either choice of our dynamics). The overlap distribution of dynamics on different networks with the same statistics shows that there are different attractors for different networks (figure 4.4.13). However, the shape of the distribution shows that the attractors are either the same (or very highly correlated) or completely different. Thus, to have different attractors with a level of correlation between them, similar to those calculated from real data in chapter 3, would likely require the structure of the network to have a degree of similarity across different cell types. Hence, multicellular life could be achieved in one of three ways: (i) creating a specific network and then adding/removing edges to create the different cell types; (ii) using the same network across all cell types, but genes have different gene-specific thresholds $\theta_i$ for expression in different cell types; (iii) using the same network across all cell types with different rates $\pi_\mu^\pm$ and/or TF binding affinities, $b_i^\mu$ (both of these have so far been kept fixed with $\dfrac{\pi_\mu^+}{\pi_\mu^-} = b_i^\mu = 1$). The latter two of these choices would allow for a multiplicity of attractors without having to re-sample graphs in simulations. On the other hand, there is biological evidence that epigenetic markers, such as histone modifications, alter the chromatin structure to make certain genes more (in)accessible for regulation in different cell types. Enhancer regions of DNA are also thought to move promoter sites toward/away from genes altering the ability of TFs to regulate them. Furthermore, there is a significant body of literature linking the restructuring of chromatin to changes in cell types during cellular reprogramming experiments [10, 39, 84, 136–139]. Therefore, the more realistic approach may be to add/remove edges from a given network to create the diversity of cell types seen in multicellular organisms. Although, it is possible that the same behaviour could be captured using a distribution of gene-specific thresholds (and/or binding affinities) with high values of the $\theta_i$ (or low values of $b_i^\mu$) effectively removing edges to $i$

**Figure 4.4.13:** The probability density function of the overlap between 150 simulations of the non-linear dynamics each on different networks with $c_1 = 1$, $c_2 = 10$, $\epsilon = 0.5$ at $T = 0$.

from the network.

Alternatively, a different choice of the regulatory interactions used in the model could lead to multiple attractors in the dynamics on a single network. Previously, the linear dynamics has been governed by the local fields $h_i(t) = \sum_{j,\mu} \frac{\xi_i \eta_j^\mu}{c_\mu^{in}} \sigma_j(t) = \sum_\mu \xi_i^\mu m_\mu(t)$. If we alter these local fields to be of the form $h_i(t) = \sum_\mu \xi_i^\mu (m_\mu(t) - a(t))$, where $a(t) = \frac{1}{N} \sum_j \sigma_j(t)$ is the activity, and change our choice of the $\boldsymbol{\xi}$ such that $\xi_i^\mu = \frac{\eta_i^\mu - \frac{c_\mu^{in}}{N}}{1 - \frac{c_\mu^{in}}{N}}$, then the local fields can be written as,

$$
\begin{aligned}
h_i(t) &= \sum_\mu \frac{\eta_i^\mu - \frac{c_\mu^{in}}{N}}{1 - \frac{c_\mu^{in}}{N}} \left[ \sum_j \frac{\eta_j^\mu}{c_\mu^{in}} - \frac{1}{N} \sum_j \right] \sigma_j(t) \\
&= \sum_j \sum_\mu \frac{\eta_i^\mu - \frac{c_\mu^{in}}{N}}{1 - \frac{c_\mu^{in}}{N}} \frac{N}{c_\mu^{in}} \left[ \frac{\eta_j^\mu}{N} - \frac{c_\mu^{in}}{N^2} \right] \sigma_j(t) \\
&= \frac{1}{N} \sum_j \left\{ \sum_\mu \frac{(\eta_i^\mu - \frac{c_\mu^{in}}{N})(\eta_j^\mu - \frac{c_\mu^{in}}{N})}{\frac{c_\mu^{in}}{N} \left( 1 - \frac{c_\mu^{in}}{N} \right)} \right\} \sigma_j(t) \,.
\end{aligned}
\tag{4.44}
$$

This is now in the form $h_i(t) = \frac{1}{N}\sum_j J_{ij}\sigma_j(t)$, where $J_{ij}$ is the coupling prescription for multiple low activity configurations of $\boldsymbol{\eta}$, i.e. (2.6) with $\rho + 1 = \rho = \mu$ and $\langle\eta^\mu\rangle = \frac{c_\mu^{\text{in}}}{N}$. Hence, for this choice of $\boldsymbol{\eta}$-dependent $\boldsymbol{\xi}$, that is equivalent to picking $\xi_i^\mu = 1$ with probability $\frac{c_\mu^{\text{in}}}{N}$ and $\xi_i^\mu = -\frac{c_\mu^{\text{in}}}{N}$ with probability $\left(1 - \frac{c_\mu^{\text{in}}}{N}\right)$, there will exist multiple fixed point gene expression level attractors in the dynamics. However, this choice of $\boldsymbol{\xi}$ is more difficult to understand from a molecular biology perspective, and therefore, it could prove more challenging to construct the bipartite gene regulatory network from experimental data with this choice of $\boldsymbol{\xi}$.

## 4.5. SUMMARY AND OUTLOOK

In this chapter, a general framework for modelling gene regulation using a bipartite network was constructed. The network structure integrates the genome and transcriptome with interactions between genes occurring only via TFs. However, due to the conservation of edges in the network, if the ratio of the number of regulatory genes to TFs $\alpha$ is known, a representative bipartite gene regulatory network can be constructed solely from the degree statistics of a gene-gene or protein-protein interaction network. It was shown that when TFs are constructed from the simultaneous synthesis of multiple gene products, they should regulate exponentially more genes on average than they are synthesised from to sustain a non-zero steady state density of expressed genes and synthesised TFs.

The degree distributions for genes and TFs govern the structure of the network, and ultimately, the steady state gene expression levels. For a fixed average TF in-degree, $c_1$, there is a critical average TF out-degree, $c_2^*$, above which a net steady state gene expression profile exists. This occurs only when a giant component exists in the bipartite gene regulatory network. The critical value $c_2^*$ can be calculated analytically for a network that contains no inhibitory interactions - in this case $c_2^*$ is also the percolation threshold of the network. When the constraint that all the genes contributing to a TF must be co-activated for it to be synthesised (i.e. the non-linear dynamics) is enforced, TFs must be constructed of a small number of proteins relative to the number of genes

they regulate, in order to sustain a steady state gene expression level. This is in line with what is currently observed from biological experiments, i.e. TF are typically small and promiscuous. However, for a less constrained version of the regulatory dynamics (i.e. the linear dynamics) this requirement vanishes.

The effects of inhibition on the dynamics was studied numerically with inhibitors shown to increase the value of $c_2^*$, because the presence of a giant cluster is no longer sufficient to sustain a steady state gene expression profile. Even with competing regulatory interactions, each network capable of sustaining a net gene expression pattern gives rise to a single attractor (or cell type), with the nature of that attractor governed by the structure of the network. The balance between activation and inhibition was found to affect the steady state density of expressed genes $a_{ss}$, as well as the frequency with which TFs are synthesised, by introducing frustration into the network. Changing the fraction of activators to inhibitors changes the amount of regulatory competition. Biasing the network towards activation increases the frequency with which genes are expressed, the number of genes that are always expressed and, therefore, increases $a_{ss}$. Hence the competition between activation and inhibition only controls the timing of gene expression, rather than facilitating multiple cell types.

Increasing the connectivity of the network also did not increase the number of attractors, but similarly altered the frequency with which genes are expressed by increasing the average number of target genes for each TF. Because any given network, with a giant cluster, was observed to support a single gene expression profile in the steady state, this supports the idea that multicellular life may require either different networks (or different rates for protein production, degradation and TF binding affinities) for each of their cell types or a different prescription for the regulatory interactions. The former could be done by changing the nature of the regulatory interactions (i.e. swapping the sign of the $\xi$-edges) or changing the accessibility of target genes (by adding or removing $\xi$-edges in the network or altering the gene-specific thresholds).

There are several pathways for future work. Firstly, only deterministic dynamics, i.e. gene reg-

ulation in the absence of noise, has been studied in this chapter. However, because the model allows for the inclusion of noise it can easily be extended to include this level of biological reality. One would expect noise to restrict the range in parameter space in which stable attractors are supported. Furthermore, the assumption that the connectivities $\eta$ and $\xi$ are statistically independent was made. This may not be true though, one might expect $\langle c^{\text{in}} \rangle$ and $\langle c^{\text{out}} \rangle$ to be correlated, with the number of DNA binding sites possibly increasing with the size of a TF. Thus, it could be worthwhile to extend this model to include correlations between the TF synthesis and gene regulation using data from knock-out experiments or known protein-protein and gene interaction networks. Also, the model does not take into account the effects of external signals (e.g. morphogen gradients and cell-cell interactions). Although, these could be included in the model with additional terms in the form of local or external fields. However, the most fruitful advancement of this model would be to investigate the nature and number of attractors supported when edges are added/removed from a network, or by using a distribution of values for the rates of protein synthesis and degradation, TF binding affinities and/or gene specific thresholds. Furthermore, if such values were dependent on gene expression profiles it may be possible for the dynamics to traverse from one attractor to another, encapsulating changes in cell state, for example, due to differentiation.

Despite the assumptions and limitations, the general framework presented in this chapter can be used to compare simulations with experimental data. However, the following data would be required to do so: (i) The distributions or statistics for in- and out-degrees of genes and/or TFs, (ii) the statistics of the number of activating and inhibiting TFs, (iii) the number of nodes in the regulatory network. Given these data, this model should be able to accurately predict the steady state density of expressed genes and frequency with which TFs are synthesised.

## APPENDIX 4.A    DEGREE DISTRIBUTIONS OF EFFECTIVE GENE-GENE INTERACTION NETWORK

Transcription factors act as intermediates in a GRN. If one would like to study the gene-gene interactions then we can ask what is the probability that they are connected by a directed edge? This involves integrating out the transcription factors of the bipartite network. If the degree sequence of the GRN are denoted by $\mathbf{k} = \{\mathbf{k}^{\text{in}}, \mathbf{k}^{\text{out}}\}$, the distribution of out-degrees in the GRN is given by,

$$p_{\text{out}}(k) = \left\langle \frac{1}{N} \sum_i \delta_{k, \sum_j A_{ij}} \right\rangle_{\eta, \xi} \tag{4.45}$$

where $A_{ij} = 1$ if a directed edge from $i$ to $j$ exists and is zero otherwise, i.e $A_{ij} = \Theta\left[\sum_\mu \eta_i^\mu |\xi_j^\mu|\right]$. In general one has,

$$p(\eta_i^\mu) = \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \delta_{\eta_i^\mu, 1} + \left(1 - \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle}\right) \delta_{\eta_i^\mu, 0}, \tag{4.46}$$

$$p(|\xi_i^\mu|) = \frac{d_i^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle} \delta_{|\xi_i^\mu|, 1} + \left(1 - \frac{d_i^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle}\right) \delta_{|\xi_i^\mu|, 0}. \tag{4.47}$$

Under the assumption that $\eta$ and $\xi$ are i.i.d random variables and rewriting the $\delta$-function using its Fourier representation, the general form of $p_{\text{out}}(k)$ is,

$$p_{\text{out}}(k) = \frac{1}{N} \sum_i \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega k} \langle e^{-\mathrm{i}\omega \sum_j A_{ij}} \rangle_{\eta,\xi} \,, \tag{4.48}$$

(where $\mathrm{i} = \sqrt{-1}$ is used to denote the imaginary number and $i$ denotes the index of the gene $i$ in the regulatory network.)

One can replace the $A_{ij}$ in the expression for $p_{\text{out}}(k)$ with $\tilde{A}_{ij} = \sum_\mu \eta_i^\mu \xi_i^\mu$. This is possible because the probability of having a connection in the network $\tilde{A}_{ij}$ is the same as in $A_{ij}$ to $\mathcal{O}(N^{-1})$, as shown below.

$$
\begin{aligned}
p(\tilde{A}_{ij}) &= \langle \delta_{\tilde{A}_{ij}, \sum_\mu \eta_i^\mu |\xi_j^\mu|} \rangle \\
&= \int \frac{d\omega}{2\pi} \langle e^{\mathrm{i}\omega\left(\tilde{A}_{ij} - \sum_\mu \eta_i^\mu |\xi_j^\mu|\right)} \rangle_{\eta,\xi} \\
&= \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega\tilde{A}_{ij}} \langle e^{-\mathrm{i}\omega \sum_\mu \eta_i^\mu |\xi_j^\mu|} \rangle_{\eta,\xi} \\
&= \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega\tilde{A}_{ij}} \prod_\mu \langle \eta_i^\mu |\xi_j^\mu| e^{-\mathrm{i}\omega} + 1 - \eta_i^\mu |\xi_j^\mu| \rangle_{\eta,\xi} \\
&= \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega\tilde{A}_{ij}} \prod_\mu \left[ \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle} e^{-\mathrm{i}\omega} + \left(1 - \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle}\right) \right] \\
&= \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega\tilde{A}_{ij}} \prod_\mu \left[ \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle} (e^{-\mathrm{i}\omega} - 1) + 1 \right] \\
&= \int \frac{d\omega}{2\pi} e^{\mathrm{i}\omega\tilde{A}_{ij}} e^{\sum_\mu \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N^2\langle d^{\text{out}}\rangle} \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{\langle d^{\text{in}}\rangle}(e^{-\mathrm{i}\omega}-1)} \,,
\end{aligned}
$$

where the sparse nature of GRN was used to exponentiate in the last line. Expanding the exponential and using the definition of the $\delta$-function one can see that,

$$p(\tilde{A}_{ij}) = \delta_{\tilde{A}_{ij},0}\left[1 - \sum_\mu \frac{d_i^{\text{out}} d_j^{\text{in}} c_\mu^{\text{out}} c_\mu^{\text{in}}}{N^2\langle d^{\text{out}}\rangle\langle d^{\text{in}}\rangle}\right] + \delta_{\tilde{A}_{ij},1}\left[\sum_\mu \frac{d_i^{\text{out}} d_j^{\text{in}} c_\mu^{\text{out}} c_\mu^{\text{in}}}{N^2\langle d^{\text{out}}\rangle\langle d^{\text{in}}\rangle}\right] + \mathcal{O}(N^{-2})\,. \tag{4.49}$$

Hence, $p(\tilde{A}_{ij} > 1)$ is $\mathcal{O}(N^{-2})$. Thus, to order $N^{-1}$, $p(\tilde{A}_{ij}) = p(A_{ij})$, so one can replace the averages over $A_{ij}$ with averages over the weighted edges $\tilde{A}_{ij}$. Therefore,

$$
\begin{aligned}
\left\langle e^{-i\omega \sum_j A_{ij}} \right\rangle_{\eta,\xi} &= \left\langle e^{-i\omega \sum_{j,\mu} \eta_i^\mu |\xi_j^\mu|} \right\rangle_{\eta,\xi} \\
&= \prod_\mu \left\langle e^{-i\omega \eta_i^\mu \sum_j |\xi_j^\mu|} \right\rangle_{\eta,\xi} \\
&= \prod_\mu \left[ 1 + \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \left( \left\langle e^{-i\omega \sum_j |\xi_j^\mu|} \right\rangle - 1 \right) \right] \\
&= \prod_\mu \left[ 1 + \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \left( \prod_j \left[ 1 + \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{N\langle d^{\text{in}}\rangle} \left( e^{-i\omega} - 1 \right) \right] - 1 \right) \right] \\
&= \prod_\mu \left[ 1 + \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \left( \exp\left\{ \frac{1}{N} \sum_j \frac{d_j^{\text{in}} c_\mu^{\text{out}}}{\langle d^{\text{in}}\rangle} \left( e^{-i\omega} - 1 \right) \right\} - 1 \right) \right] \\
&= \prod_\mu \left[ 1 + \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{N\langle d^{\text{out}}\rangle} \left( \exp\{ c_\mu^{\text{out}}(e^{-i\omega} - 1) \} - 1 \right) \right] \\
&= \exp\left\{ \frac{1}{N} \sum_\mu \frac{d_i^{\text{out}} c_\mu^{\text{in}}}{\langle d^{\text{out}}\rangle} \left( \exp\{ c_\mu^{\text{out}}(e^{-i\omega} - 1) \} - 1 \right) \right\} \\
&= e^{\frac{-d_i^{\text{out}} \alpha \langle c^{\text{in}}\rangle}{\langle d^{\text{out}}\rangle}} e^{\frac{\alpha d_i^{\text{out}}}{\langle d^{\text{out}}\rangle} \langle c^{\text{in}} e^{c^{\text{out}}(\exp(-i\omega)-1)} \rangle},
\end{aligned}
\tag{4.50}
$$

where, the last average $\langle c^{\text{in}} e^{c^{\text{out}}(\exp(-i\omega)-1)} \rangle$ is taken over the joint distribution $p(c^{\text{in}}, c^{\text{out}})$. Substituting this result into our expression for $p_{\text{out}}(k)$ and using the conservation of in- and out-degrees, gives

$$
\begin{aligned}
p_{\text{out}}(k) &= \frac{1}{N} \sum_i e^{-d_i^{\text{out}}} \int \frac{d\omega}{2\pi} e^{i\omega k} \exp\left\{ \frac{\alpha d_i^{\text{out}}}{\langle d^{\text{out}}\rangle} \langle c^{\text{in}} e^{c^{\text{out}}(e^{-i\omega}-1)} \rangle \right\} \\
&= \sum_d e^{-d} p_{\text{out}}(d) \int \frac{d\omega}{2\pi} e^{i\omega k} \exp\left\{ \frac{\alpha d}{\langle d^{\text{out}}\rangle} \langle c^{\text{in}} e^{c^{\text{out}}(e^{-i\omega}-1)} \rangle \right\},
\end{aligned}
\tag{4.51}
$$

where $p_{\text{out}}(d)$ is the out-degree distribution of the genes. For independent in- and out-degrees the degree distribution of the transcription factors factorises, $p(c^{\text{in}}, c^{\text{out}}) = p(c^{\text{in}})p(c^{\text{out}})$, and

making use of the conservation of edges in the bipartite network one has,

$$
\begin{aligned}
p_{\text{out}}(k) &= \sum_d e^{-d} p_{\text{out}}(d) \int \frac{d\omega}{2\pi} e^{i\omega k} \exp\left\{ d \left\langle \exp\left( c^{\text{out}}(e^{-i\omega} - 1) \right) \right\rangle \right\} \\
&= \sum_\lambda \frac{1}{\lambda!} \sum_d d^\lambda e^{-d} p_{\text{out}}(d) \int \frac{d\omega}{2\pi} e^{i\omega k} \left[ \sum_c p_{\text{out}}(c) \exp\left( c(e^{-i\omega} - 1) \right) \right]^\lambda \\
&= \sum_\lambda \frac{1}{\lambda!} \sum_d p_{\text{out}}(d) d^\lambda e^{-d} \int \frac{d\omega}{2\pi} e^{i\omega k} \int dx \delta\left( x - \sum_{r=1}^{\lambda} c_r \right) \\
&\qquad\qquad\qquad\qquad \times \left\{ \sum_{c_1,\ldots,c_\lambda} p_{\text{out}}(c_1) \ldots p_{\text{out}}(c_\lambda) e^{-x} \sum_s \frac{x^s}{s!} e^{-i\omega s} \right\} \\
&= \sum_\lambda \frac{1}{\lambda!} \sum_d p_{\text{out}}(d) d^\lambda e^{-d} \sum_{c_1,\ldots,c_\lambda} p_{\text{out}}(c_1) \ldots p_{\text{out}}(c_\lambda) \int dx \delta\left( x - \sum_{r=1}^{\lambda} c_r \right) e^{-x} \frac{x^k}{k!},
\end{aligned}
$$

$$(4.52)$$

Thus, the out-degree distribution in the effective gene-gene interaction network as,

$$
p_{\text{out}}(k) = \int dx e^{-x} \frac{x^k}{k!} P(x),
\tag{4.53}
$$

where

$$
P(x) = \sum_d p_{\text{out}}(d) e^{-d} \sum_{\lambda \geq 0} \frac{d^\lambda}{\lambda!} \sum_{c_1,\ldots,c_\lambda} p_{\text{out}}(c_1) \ldots p_{\text{out}}(c_\lambda) \delta\left( x - \sum_{r=1}^{\lambda} c_r \right).
\tag{4.54}
$$

Clearly, the out-degree distribution is normalised $\sum_{k\geq 0} p_{\text{out}}(k) = 1$. The average out-degree of the effective gene-gene interaction network would then by given by,

$$
\begin{aligned}
\langle k_{\text{out}} \rangle &= \sum_{k \geq 0} k p_{\text{out}}(k) \\
&= \int_0^\infty dy P(x) e^{-x} x \sum_{k \geq 0} \frac{x^{k-1}}{(k-1)!} \\
&= \int_0^\infty x P(x) dx \\
&= \sum_d p_{\text{out}}(d) e^{-d} \sum_{\lambda \geq 0} \frac{d^\lambda}{\lambda!} \sum_{c_1,\dots,c_\lambda} p_{\text{out}}(c_1) \dots p_{\text{out}}(c_\lambda) \sum_{r \leq \lambda} c_r \\
&= \sum_d p_{\text{out}}(d) e^{-d} \sum_{\lambda \geq 0} \frac{d^\lambda}{\lambda!} \lambda \sum_c p_{\text{out}}(c) c \\
&= \langle c^{\text{out}} \rangle \sum_d p_{\text{out}}(d) e^{-d} d \sum_{\lambda \geq 0} \frac{d^{\lambda-1}}{(\lambda-1)!} \\
&= \langle c^{\text{out}} \rangle \sum_d p_{\text{out}}(d) d \\
&= \langle c^{\text{out}} \rangle \langle d^{\text{out}} \rangle = \alpha c_1 c_2
\end{aligned}
\tag{4.55}
$$

It can be shown in a similar fashion that the in-degree distribution for the effective gene-gene interaction network is,

$$
p_{\text{in}}(k) = \int dy e^{-y} \frac{y^k}{k!} P(y) \,,
\tag{4.56}
$$

where

$$
P(y) = \sum_d p_{\text{in}}(d) e^{-d} \sum_{\lambda \geq 0} \frac{d^\lambda}{\lambda!} \sum_{c_1,\dots,c_\lambda} p_{\text{in}}(c_1) \dots p_{\text{in}}(c_\lambda) \delta \left( y - \sum_{r=1}^\lambda c_r \right) \,.
\tag{4.57}
$$

This is clearly normalised and gives the average in-degree in the effective gene-gene interaction network as $\langle k^{\text{in}} \rangle = \int dy y P(y) = \alpha c_1 c_2$.

## Appendix 4.B    Percolation thresholds

### 4.B.1    Non-linear dynamics

Here, the critical value of the the transcription factor out-degree $c_2^*$, above which a giant cluster will exist in the gene regulatory network, is calculated. The critical value can be found through a simple stability analysis.

Recall that we are using an adaptation of the cavity method to determine the probability that a gene or TF belongs to the giant cluster. In (4.26) the gene $i$ is connected to the TF $\mu$ via $\xi_i^\mu$. Then in the construction of the cavity graph for $n_\mu^{(i)}$, in (4.29), one removes all the genes connected to the TF $\mu$ via an $\eta$-edge. Due to the sparsity of the bipartite network, the likelihood that the gene $i$ contributes to the synthesis of the TF $\mu$ as well as being regulated by it is $\mathcal{O}(\frac{1}{N})$. Therefore, all of $c_\mu^{in}$ of the neighbours of $\mu$ must be in the giant cluster for $\mu$ to also belong to the giant cluster. This leads to, for the non-linear dynamics, the system of equations (4.32) & (4.33) which in their general form are:

$$\hat{g} = \langle n_i^{(\mu)} \rangle = \sum_{d^{\text{out}}} P(d^{\text{out}}) \frac{d^{\text{out}}}{\langle d^{\text{out}} \rangle} \sum_{d^{\text{in}}=1}^{\infty} P(d^{\text{in}}|d^{\text{out}}) \left[ 1 - (1 - \tilde{g})^{d^{\text{in}}} \right] , \qquad (4.58)$$

$$\tilde{g} = \langle n_\mu^{(i)} \rangle = \sum_{c^{\text{out}}} P(c^{\text{out}}) \frac{c^{\text{out}}}{\langle c^{\text{out}} \rangle} \sum_{c^{\text{in}}=1}^{\infty} P(c^{\text{in}}|c^{\text{out}}) \hat{g}^{c^{\text{in}}} . \qquad (4.59)$$

Recall, $n_i^{(\mu)}$ is the indicator variable for the gene $i$ on the cavity graph with the TF $\mu$ removed. Therefore, (4.58) arises from averaging (4.29) over all possible genes that have $\mu$ as a successor (i.e. have an out-degree that terminate at node $\mu$) and their predecessors (i.e. all TFs that have an incoming link to node $i$). This is done in (4.58) by taking the average over the degree distribution of the successors of $i$, $P(d^{\text{out}})d^{\text{out}}/\langle d^{\text{out}} \rangle$, and the average over the in-degree of $i$ on the cavity graph with $\mu$ and its in-degrees removed. A similar line of reasoning gives the general from of (4.59) by averaging on the cavity graph with $i$ and its in-degrees removed. Under the

assumption of independent edges $\eta$ and $\xi$, the in- and out-degrees for a node are independent. Combining this with the sparsity and directed nature of the bipartite gene regulatory network, the conditional distributions are equivalent to their corresponding marginals, $P(d^{\text{in}}|d^{\text{out}}) = P(d^{\text{in}})$ and $P(c^{\text{in}}|c^{\text{out}}) = P(c^{\text{in}})$. Then, the average over the successor degree distributions has no affect, and the cavity probabilities become equal to the probabilities of belonging to the giant cluster, i.e. $g = \hat{g}$ and $\bar{g} = \tilde{g}$, with

$$\hat{g} = \langle n_i^{(\mu)} \rangle = \sum_{d^{\text{in}}=1}^{\infty} P(d^{\text{in}}) \left[ 1 - (1 - \tilde{g})^{d^{\text{in}}} \right] = f_1(\tilde{g}, \hat{g}) \,, \tag{4.60}$$

$$\tilde{g} = \langle n_\mu^{(i)} \rangle = \sum_{c^{\text{in}}=1}^{\infty} P(c^{\text{in}})\hat{g}^{c^{\text{in}}} = f_2(\tilde{g}, \hat{g}) \,. \tag{4.61}$$

The point $(\tilde{g}, \hat{g}) = (0,0)$ is always a solution to these equations. This solution corresponds to the situation where there is no giant cluster in the network. Thus, the point at which this solution is no longer stable, will be the point at which a giant cluster emerges in the network. The solution $(\tilde{g}, \hat{g}) = (0,0)$ is stable provided that,

$$\left| \mathbf{J} |_{(\tilde{g},\hat{g})=(0,0)} \right| = \left| \left. \frac{\partial [f_1, f_2]}{\partial [\tilde{g}, \hat{g}]} \right|_{(0,0)} \right| < 1 \,.$$

Taking partial derivatives of (4.60) and (4.61) gives the Jacobian of the system evaluated at the point $(\tilde{g}, \hat{g}) = (0,0)$ as

$$\left| \mathbf{J} |_{(0,0)} \right| = \begin{vmatrix} 0 & \langle d^{\text{in}} \rangle \\ P(c^{\text{in}} = 1) & 0 \end{vmatrix} \,. \tag{4.62}$$

Hence, for Poisson degree distributions and the stability criteria $\left| \mathbf{J} |_{(\tilde{g},\hat{g})=(0,0)} \right| < 1$, a giant cluster will exist in the network, if and only if,

$$\alpha c_2 c_1 e^{-c_1} \geq 1 \,, \tag{4.63}$$

giving the percolation threshold for the bipartite gene regulatory network

$$c_2^* = \frac{e^{c_1}}{\alpha c_1} \ .$$

(4.64)

For Poisson degree distributions, $P(d^{\mathrm{in}}) = e^{-d_2} d_2^{d^{\mathrm{in}}}/d^{\mathrm{in}}!$ and $P(c^{\mathrm{in}}) = e^{-c_1} c_1^{c^{\mathrm{in}}}/c^{\mathrm{in}}!$, the averages in $(4.60)$ and $(4.61)$ can be performed exactly to find the probability that gene or TF belongs to the giant cluster, by making use of $\exp(x) = \sum_{k=0}^{\infty} \frac{x^k}{k!}$ and $d_2 = \alpha c_2$

$$g = \hat{g} = 1 - e^{-\alpha c_2 \tilde{g}}$$

(4.65)

$$\tilde{g} = \bar{g} = e^{c_1(\hat{g}-1)} - e^{-c_1}$$

(4.66)

These curves are plotted in figure 4.4.7 with $\bar{g}$ normalised by the $P(c^{\mathrm{in}} > 0) = 1 - P(c^{\mathrm{in}} = 0) = 1 - e^{-c_1}$.

### 4.B.2   LINEAR DYNAMICS

The percolation threshold for the linear dynamics can be found following a similar line of reasoning to that for the non-linear dynamics. The key difference is that there is no longer the hard constraint that all genes contributing to a TF are required to belong to the giant cluster for it to be as well. Instead it is sufficient that at least one the genes contributing to a TF belongs to the giant cluster. Thus, with this constraint relaxed, the equations for the indicators become symmetric and one obtains the cavity probabilities,

$$\hat{g} = \langle n_i^{(\mu)} \rangle = \sum_{d^{\mathrm{out}}} P(d^{\mathrm{out}}) \frac{d^{\mathrm{out}}}{\langle d^{\mathrm{out}} \rangle} \sum_{d^{\mathrm{in}}=1}^{\infty} P(d^{\mathrm{in}}|d^{\mathrm{out}}) \left[ 1 - (1 - \tilde{g})^{d^{\mathrm{in}}} \right] ,$$

(4.67)

$$\tilde{g} = \langle n_\mu^{(i)} \rangle = \sum_{c^{\mathrm{out}}} P(c^{\mathrm{out}}) \frac{c^{\mathrm{out}}}{\langle c^{\mathrm{out}} \rangle} \sum_{c^{\mathrm{in}}=1}^{\infty} P(c^{\mathrm{in}}|c^{\mathrm{out}}) \left[ 1 - (1 - \hat{g})^{c^{\mathrm{in}}} \right] ,$$

(4.68)

and therefore the analysis is much simpler following the same steps as for the non-linear dynamics. The point $(\tilde{g}, \hat{g}) = (0, 0)$ is still a solution with the Jacobian evaluated at this point given by,

$$\mathbf{J}|_{(\tilde{g},\hat{g})=(0,0)} = \begin{pmatrix} 0 & \langle d^{\text{in}} \rangle \\ \langle c^{\text{in}} \rangle & 0 \end{pmatrix} . \tag{4.69}$$

Therefore, the criteria for a giant component to exist in the bipartite gene regulatory network is $\alpha c_2 c_1 \geq 1$ for the Poisson degree distributions. This gives rise to the percolation threshold for the linear dynamics,

$$c_2^* = \frac{1}{\alpha c_1} . \tag{4.70}$$

Then the probabilities that a gene or TF belong to the giant cluster, for the case of Poisson degree distributions, are

$$g = \hat{g} = 1 - e^{-\alpha c_2 \tilde{g}} , \tag{4.71}$$

$$\bar{g} = \tilde{g} = 1 - e^{-c_1 \hat{g}} . \tag{4.72}$$

Again, these probabilities are plotted in figure 4.4.7.


APPENDIX 4.C    OVERLAP BETWEEN SIMULATIONS

The overlap $q_{\alpha\beta}$ is used as a measure of similarity of the steady state gene expression levels between two simulation runs, by comparing the TF profiles between different simulations. Formally it is defined as,

$$q_{\alpha,\beta} = \frac{\tilde{q}_{\alpha\beta}}{\sqrt{\tilde{q}_{\alpha\alpha}\tilde{q}_{\beta\beta}}} , \tag{4.73}$$

with

$$\tilde{q}_{\alpha\beta} = \frac{1}{\alpha N} \sum_\mu \left( \langle \bar{\tau}_\mu \rangle_\alpha - \langle \bar{\tau} \rangle_\alpha \right) \left( \langle \bar{\tau}_\mu \rangle_\beta - \langle \bar{\tau} \rangle_\beta \right) , \tag{4.74}$$

$$\langle \bar{\tau} \rangle_\alpha = \frac{1}{\alpha N} \sum_\mu \langle \bar{\tau}_\mu \rangle_\alpha \,, \tag{4.75}$$

where $\langle \ldots \rangle_\alpha$ denotes the time average in the steady state of simulation run $\alpha$; and $\tau_\mu$ indicates whether a TF has been synthesised or not with

$bar\tau_\mu(t) = \Theta\left[\tau_\mu(t)\right]$. The overlap is strictly defined in the same manner as a Pearson correlation coefficient with $q_{\alpha\beta} \in [-1, 1]$. When $\beta = \alpha$, $q_{\alpha\alpha} = 1$ and the gene expression levels are identical. Any value of $q_{\alpha\beta} \neq 1$ indicates that there is a difference in the steady state gene expression profiles. The distribution of $q_{\alpha\beta}$ is studied to determine whether or not multiple attractors exist for the each choice of the dynamics (linear and non-linear) on a given network.

## APPENDIX 4.D    INFERRED ATTRACTOR LENGTH FROM DEVIATIONS IN THE OVER-LAP DISTRIBUTION

In this appendix, it is demonstrated how one can infer the possible length of a limit cycle attractor of the dynamics in a bipartite gene regulatory network. For this example the overlap distributions in figure 4.4.9 are used along with the statistics of the network the dynamics are simulated on for those distributions, i.e. $(c_1, c_2, \epsilon) = (1, 10, 0.5)$.

If we assume that the transient to the steady state is short on this network, as generally demonstrated for increasing connectivity in figure 4.4.4, then the dynamics will converge to the limit cycle attractor in $\mathcal{O}(1)$ time steps. Next, if the time window $\Delta t$ used to perform the average over the steady state is much greater than the length of the limit cycle $\ell$, the limit cycle will be fully traversed $N_{\Delta t}$ times. There can also be a fraction of the limit cycle length $x$ traversed in each $\Delta t$ as well as the $N_{\Delta t}$ full cycles giving $\Delta t = \ell\left(N_{\Delta t} + x\right)$. Over many simulation runs the average of $x$ would be expected to be $\langle x \rangle = 0.5$. If the probability that any gene in the network is expressed is given by $a$ then the probability that a gene is expressed differently in simulations runs is given by $2a(1 - a)$. That is, a gene $i$ could be expressed in simulation $\alpha$ and not in $\beta$ or vice versa. Then

| $\Delta t$ | $\approx \langle \Delta q \rangle$ | $\ell$ |
|---|---|---|
| 500 | $2.3 \times 10^{-3}$ | 4.6 |
| $2,000$ | $6.3 \times 10^{-4}$ | 5.0 |
| $4,500$ | $2.6 \times 10^{-4}$ | 4.7 |

**Table 4.D.1:** Length of limit cycle attractors inferred from mean deviation from $q_{\alpha\beta} = 1$ in distributions of the overlap using different time windows for the steady state averaging.

the average deviation in the overlap from $q_{\alpha\beta} = 1$ is given by

$$\langle \Delta q \rangle = \frac{2a(1-a)x\ell}{\Delta t} \, , \tag{4.76}$$

where $x\ell$ is the number of sites on which there is a difference in gene expression between simulation runs. Using this expression for $\langle \Delta q \rangle$ and the mean values of the deviation from $q_{\alpha\beta} = 1$ in figure 4.4.9, one can infer the length of the limit cycle attractor for that network shown in the table 4.D.1. There is also an $\mathcal{O}(1/N)$ effect on the overlap of the initially introduced TFs hitting finite clusters that are able to sustain a net steady state gene expression level alongside the contribution from the giant cluster. Therefore, the network used to construct the overlap distributions in figure 4.4.9 likely has an attractor that has is a limit cycle of length $\ell \simeq 5$. However, this periodicity is difficult to identify directly from trajectories of the dynamics.

# 5

## Summary & outlook

In this thesis, tools from the field of statistical physics, in particular, the theories of neural networks and directed random graphs, and Monte Carlo simulation methods, have been employed to explore how the expression of genes may lead to multicellular life and the known transitions between different cell types. The motivation behind this thesis has been twofold. Firstly, to create predictive models that are more biologically grounded than those currently available to understand one of the biggest advancements in molecular biology - cellular reprogramming. Secondly, to build a powerful mathematical framework that encourages communication and collaboration between experimentalists and theoreticians interested in this field of research.

In chapter 2, a model for cell reprogramming was constructed that builds on key features of cellular biology - mainly cell cycles and potencies. This level of biological realism does not currently exist in other neural network type models for the interaction between gene expression levels. In the model, cell types are hierarchically related dynamical attractors of the effective interactions between gene expression levels. Stages of the cell cycle are fully characterised by the configuration of gene expression levels, and reprogramming corresponds to triggering transitions between such configurations. Two possible mechanisms for reprogramming were found: cycle specific perturbations and a noise-induced switching. The former corresponds to a directed perturbation that induces a transition into a cycle-state of a different cell type in the potency hierarchy (mainly a stem cell), whilst the latter is a priori undirected and could be induced, e.g. by a (stochastic) change in the cellular environment. These reprogramming protocols were found to be effective in large regimes of the parameter space and make specific predictions concerning reprogramming dynamics that are broadly in line with experimental findings, including the number of genes that need to be perturbed to reprogram a cell to an induced pluripotent stem cell.

More specifically, two critical points were found numerically in the phase space of the model. The first, a critical noise level $T$ above which the dynamics of a system converges towards the gene expression profiles of a stem cell from an initial state that has a strong correlation with a differentiated cell. Next, there were a critical fraction of gene expression levels $q_r$ which could be altered to achieve the same kind of reprogramming. These perturbations mimic Yamanaka's line of reasoning for the original reprogramming experiment, by altering the gene expression level of a small number of genes to be the same as in the cell type that one desires upon reprogramming, namely a stem cell. At low levels of noise, this fraction of genes is in-line with the number expected to be affected by retroviral transduction of the Yamanaka reprogramming factors. Combining these two results it was shown how the number of perturbed genes required to transition to a stem cell like state decreases with increasing noise. Perturbations, through direct changes in gene expression levels, are most effective when they do not act to disrupt the progression of the cell cycle

and the desired final state is the cell cycle phase with the shortest Hamming distance between the initial and desired resultant cell type.

In chapter 3, data from RNA sequencing experiments were used to further test the hypotheses of the cell reprogramming model introduced in chapter 2. It was shown that there is evidence in support of a single cell cycle phase with maximal similarity across different cell types of the mouse enteric nervous system. Interestingly, this cycle phase appears to be the $G_1$ phase rather than the $S$ or $M$ phases as hypothesised in chapter 2. The higher level of similarity in this phase is manifested in an increase in the number of expressed genes during $G_1$ leading to strong correlations in gene expression profiles. Data from the Human Protein Atlas project was used to find the average density of expressed genes and the correlation between gene expression levels from different tissue types. The 70% of expressed genes in human cells is in agreement with previous studies and provided a useful benchmark for calibration. However, whilst validating parts of the reprogramming model, the current lack of available cell cycle specific data, especially for human cells, means that conclusions drawn from the analysis in chapter 3 may evolve with experimental technologies and increasing efforts in this area.

In chapter 4, a general framework for gene regulation was constructed as the dynamics on a directed bipartite graph. The integration of the transcriptome and genome into a single bipartite gene regulatory network allowed for a greater understanding of the structure and dynamics of gene regulation compared to previous Boolean or neural network models. Specifically, it demonstrated that transcription factors should typically regulate exponentially many more genes than those that contribute to their synthesis. Using percolation theory and an adaptation of the cavity method analytical expressions were derived for the average out-degree of a TF required to create a giant cluster in the gene regulatory network, when TFs are either fully formed from the products of expressed genes and when they may function with errors in their production. Later in this chapter, numerical analysis showed that the competition between regulatory interactions controls the frequency with which genes are expressed, rather than the number of cell types that arise from a

gene regulatory network. With any given gene regulatory network shown to support just a single gene expression profile, the choice of interactions in the model provides support for the hypothesis that a requisite for multicellular life is a rewiring of the underlying gene regulatory network in different cell types. This is most likely achieved in nature through chromatin (de)condensation, or alternative TF behaviour in different cell types, e.g. different rates of formation and/or binding affinities for their target genes, of which there is already significant evidence for in the experimental literature. The bipartite gene regulatory network can easily be adapted to study either of these effects due to its general ground-up construction.

Fruitful directions for future research include improving the understanding, and quantification of, the level of biological noise strength $T$ in the gene expression level dynamics. This should not only illustrate the meaning of the noise-induced reprogramming mechanism found in chapter 2, but it would also allow for a meaningful level of stochasticity to be introduced to the currently deterministic bipartite gene regulatory dynamics studied in chapter 4. Although, developing either model to include the levels of correlation between the gene expression profiles of an organism seen in chapter 3, would be the most desirable advancement of this work. Suggestions on how to achieve this are detailed in the summaries of the respective chapters.

Finally, as with any theoretical model, continued attempts should be made to fit parameter values, test and scrutinise the hypotheses and assumptions made in the cell reprogramming and bipartite gene regulatory network models contained in this thesis. For these models, like any, are only interesting from a mathematical perspective unless they can accurately capture and predict the behaviour of the systems that they are based on. This will require communication and collaboration between theorists and experimentalists, as with any interdisciplinary research. It is the hope of the author that this thesis will facilitate such work, regardless of any scrutiny that this may lead to in the accuracy of this thesis.

# References

[1] J. Frisén, U. Lendahl, and T. Perlmann. The 2012 Nobel Prize in Physiology or Medicine - Advanced Information. Technical report, 2012. URL `http://www.nobelprize.org/nobel{_}prizes/medicine/laureates/2012/advanced.html`.

[2] A. B. C. Cherry and G. Q. Daley. Reprogramming Cellular Identity for Regenerative Medicine. *Cell*, 148:1110–1122, 2012. doi: 10.1016/j.cell.2012.02.031.

[3] A. B. C. Cherry and G. Q. Daley. Reprogrammed cells for disease modeling and regenerative medicine. *Annu. Rev. Med.*, 64:277–290, 2013. doi: 10.1146/annurev-med-050311-163324. URL `http://www.ncbi.nlm.nih.gov/pubmed/23327523http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3629705`.

[4] S. J. Engle and D. Puppala. Perspective Integrating Human Pluripotent Stem Cells into Drug Development. *Cell Stem Cell*, 12:669–677, 2013. doi: 10.1016/j.stem.2013.05.011. URL `http://dx.doi.org/10.1016/j.stem.2013.05.011`.

[5] R. R. Kanherkar, et al. Cellular reprogramming for understanding and treating human disease. *Front. Cell Dev. Biol.*, 2(67):1–21, nov 2014. doi: 10.3389/fcell.2014.00067.

URL `http://journal.frontiersin.org/article/10.3389/fcell.2014.00067/abstract`.

[6] Y. Avior, I. Sagi, and N. Benvenisty. Pluripotent stem cells in disease modelling and drug discovery. *Nat. Rev. Mol. Cell Biol.*, 17(3):170–182, jan 2016. doi: 10.1038/nrm.2015.27. URL `http://www.nature.com/doifinder/10.1038/nrm.2015.27`.

[7] R. Morris, et al. Mathematical approaches to modeling development and reprogramming. *Proc. Natl. Acad. Sci.*, 111(14):5076–5082, apr 2014. doi: 10.1073/pnas.1317150111. URL `http://www.pnas.org/cgi/doi/10.1073/pnas.1317150111`.

[8] I. Sancho-Martinez, et al. Reprogramming by lineage specifiers: Blurring the lines between pluripotency and differentiation. *Curr. Opin. Genet. Dev.*, 28:57–63, 2014. doi: 10.1016/j.gde.2014.09.009. URL `http://dx.doi.org/10.1016/j.gde.2014.09.009`.

[9] K. Takahashi and S. Yamanaka. A developmental framework for induced pluripotency. *Development*, 142(19):3274–3285, 2015. doi: 10.1242/dev.114249. URL `http://www.ncbi.nlm.nih.gov/pubmed/26443632`.

[10] Z. D. Smith, C. Sindhu, and A. Meissner. Molecular features of cellular reprogramming and development. *Nat. Rev. Mol. Cell Biol.*, 17:139–154, 2016. doi: 10.1038/nrm.2016.6.

[11] C. H. Waddington. The Strategy of the Genes: A Discussion of Some Aspects of Theoretical Biology. Routledge, Taylor & Francis Group, New York, routledge edition, 1957.

[12] S. Kauffman. Metabolic stability and epigenesis in randomly constructed genetic nets. *J. Theor. Biol.*, 22(3):437–467, 1969. doi: 10.1016/0022-5193(69)90015-0.

[13] L. Wolpert and J. H. Lewis. Towards a Theory of Development. *Fed. Proc.*, 34(1):14–20, 1975.

[14] J. B. Gurdon. The Developmental Capacity of Nuclei taken from Intestinal Epithelium Cells of Feeding Tadpoles. *J. Embryol. Exp. Morphol.*, 10(4):622–640, 1962.

[15] H. M. Blau, et al. Cytoplasmic activation of human nuclear genes in stable heterocaryons. *Cell*, 32(4):1171–80, apr 1983. doi: 10.1016/0092-8674(83)90300-8. URL http://www.ncbi.nlm.nih.gov/pubmed/6839359.

[16] N. Takagi, et al. Reversal of X-inactivation in female mouse somatic cells hybridized with murine teratocarcinoma stem cells in vitro. *Cell*, 34(3):1053–62, oct 1983. doi: 10.1016/0092-8674(83)90563-9. URL http://www.ncbi.nlm.nih.gov/pubmed/6627391.

[17] K. Takahashi, et al. Induction of Pluripotent Stem Cells from Mouse Embryonic and Adult Fibroblast Cultures by Defined Factors. *Cell*, 126(4):663–676, aug 2006. doi: 10.1016/j.cell.2006.07.024. URL http://linkinghub.elsevier.com/retrieve/pii/S0092867406009767.

[18] K. Takahashi, et al. Induction of Pluripotent Stem Cells from Adult Human Fibroblasts by Defined Factors. *Cell*, 131(5):861–872, nov 2007. doi: 10.1016/j.cell.2007.11.019. URL http://linkinghub.elsevier.com/retrieve/pii/S0092867407014717.

[19] J. Hanna, et al. Direct cell reprogramming is a stochastic process amenable to acceleration. *Nature*, 462:595–601, 2009. doi: 10.1038/nature08592.

[20] J. Woong Han and Y.-s. Yoon. Induced Pluripotent Stem Cells: Emerging Techniques for Nuclear Reprogramming. *Antioxi Redox Signal.*, 15:1799–1813, 2011. doi: 10.1089/ars.2010.3814.

[21] P. Hou, et al. Pluripotent Stem Cells Induced from Mouse Somatic Cells by Small-Molecule Compounds. *Science (80-. ).*, 341:651–654, 2013.

[22] D. Ramsköld, et al. An Abundance of Ubiquitously Expressed Genes Revealed by Tissue Transcriptome Sequence Data. *PLoS Comput Biol*, 5(12):1–11, 2009. doi: 10.1371/journal.pcbi.1000598.

[23] F. Ponten, et al. A global view of protein expression in human cells, tissues, and organs. *Mol. Syst. Biol.*, 5(337):1–9, 2009. doi: 10.1038/msb.2009.93. URL http://www.ncbi.nlm.nih.gov/pubmed/20029370.

[24] C. Trapnell. Defining cell types and states with single-cell genomics. *Genome Res.*, 25(10):1491–1498, oct 2015. doi: 10.1101/gr.190595.115. URL http://genome.cshlp.org/lookup/doi/10.1101/gr.190595.115.

[25] F. Buettner, et al. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat. Biotechnol.*, 33(2):155–162, 2015. doi: 10.1038/nbt.3102.

[26] C. Marr, J. X. Zhou, and S. Huang. Single-cell gene expression profiling and cell state dynamics: collecting data, correlating data points and connecting the dots. *Curr. Opin. Biotechnol.*, 39:207–214, 2016. doi: 10.1016/j.copbio.2016.04.015. URL http://dx.doi.org/10.1016/j.copbio.2016.04.015.

[27] V. Moignard and B. Göttgens. Transcriptional mechanisms of cell fate decisions revealed by single cell expression profiling. *BioEssays*, 36(4):419–426, apr 2014. doi: 10.1002/bies.201300102. URL http://doi.wiley.com/10.1002/bies.201300102.

[28] B. Treutlein, et al. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature*, 509:371–375, 2014. doi: 10.1038/nature13173.

[29] A. C. D'Alessio, et al. A Systematic Approach to Identify Candidate Transcription Factors that Control Cell Identity. *Stem Cell Reports*, 5:763–775, nov 2015. doi: 10.1016/j.stemcr.2015.09.016. URL `http://linkinghub.elsevier.com/retrieve/pii/S2213671115002787`.

[30] A. Olsson, et al. Single-cell analysis of mixed-lineage states leading to a binary cell fate choice. *Nature*, 537:698–702, 2016.

[31] S. Hormoz, et al. Inferring Cell-State Transition Dynamics from Lineage Trees and Endpoint Single-Cell Measurements. *Cell Syst.*, 3(5):419–433, nov 2016. doi: 10.1016/j.cels.2016.10.015. URL `http://linkinghub.elsevier.com/retrieve/pii/S2405471216303349`.

[32] A. M. Molinaro and B. J. Pearson. In silico lineage tracing through single cell transcriptomics identifies a neural stem cell population in planarians. *Genome Biol.*, 17(87):1–17, 2016. doi: 10.1186/s13059-016-0937-9.

[33] S. Huang, et al. Cell Fates as High-Dimensional Attractor States of a Complex Gene Regulatory Network. *Phys. Rev. Lett.*, 94(12):128701–4, 2005. doi: 10.1103/PhysRevLett.94.128701.

[34] S. Huang. Reprogramming cell fates: Reconciling rarity with robustness. *BioEssays*, 31:546–560, 2009. doi: 10.1002/bies.200800189.

[35] D. V. Foster, et al. A model of sequential branching in hierarchical cell fate determination. *J. Theor. Biol.*, 260:589–597, 2009. doi: 10.1016/j.jtbi.2009.07.005. URL `http://ac.els-cdn.com/S002251930900318X/1-s2.0-S002251930900318X-main.pdf?{_}tid=5941f78a-0324-11e7-ba01-00000aacb35f{&}acdnat=1488884211{_}ddcd643580f0557697100c7f1c94bc27`.

[36] J. X. Zhou, L. Brusch, and S. Huang. Predicting Pancreas Cell Fate Decisions and Reprogramming with a Hierarchical Multi-Attractor Model. *PLoS One*, 6(3):e14752, mar 2011. doi: 10.1371/journal.pone.0014752. URL `http://dx.plos.org/10.1371/journal.pone.0014752`.

[37] J. X. Zhou and S. Huang. Understanding gene circuits at cell-fate branch points for rational cell reprogramming. *Trends Genet.*, 27(2):55–62, 2011. doi: 10.1016/j.tig.2010.11.002. URL `http://ac.els-cdn.com/ S0168952510002222/1-s2.0-S0168952510002222-main.pdf?{_}tid= cf7b000a-0323-11e7-90c3-00000aacb362{&}acdnat= 1488883979{_}7253e6b88ab4f5a6619d655d96004121`.

[38] A. Soufi, et al. Pioneer Transcription Factors Target Partial DNA Motifs on Nucleosomes to Initiate Reprogramming. *Cell*, 161:1–14, 2014. doi: 10.1016/j.cell.2015.03.017. URL `http://dx.doi.org/10.1016/j.cell.2015.03.017`.

[39] M. N. Krause, I. Sancho-Martinez, and J. C. Izpisua Belmonte. Understanding the molecular mechanisms of reprogramming. *Biochem. Biophys. Res. Commun.*, 473(3):693–697, 2016. doi: 10.1016/j.bbrc.2015.11.120.

[40] M. N. Artyomov, A. Meissner, and A. K. Chakraborty. A Model for Genetic and Epigenetic Regulatory Networks Identifies Rare Pathways for Transcription Factor Induced Pluripotency. *PLoS Comput. Biol.*, 6(5):e1000785, may 2010. doi: 10.1371/journal.pcbi.1000785. URL `http://dx.plos.org/10.1371/journal.pcbi.1000785`.

[41] T. Ly, et al. A proteomic chronology of gene expression through the cell cycle in human myeloid leukemia cells. *Elife*, 3:1–36, 2014. doi: 10.7554/eLife.01630.

[42] A. Scialdone, et al. Computational assignment of cell-cycle stage from single-cell transcriptome data. *Methods*, 85:54–61, 2015. doi: 10.1016/j.ymeth.2015.06.021.

[43] R. Grima, et al. Steady-state fluctuations of a genetic feedback loop: An exact solution. *J. Chem. Phys.*, 137:35104, 2012. doi: 10.1063/1.4736721. URL `http://dx.doi.org/10.1063/1.4736721]`.

[44] S. Smith and R. Grima. Plasticity of the truth table of low-leakage genetic logic gates. *Phys. Rev. E*, 98(062410), 2018. doi: 10.1103/PhysRevE.98.062410. URL `http://grimagroup.bio.ed.ac.uk/documents/SmithGrimaLogic.pdf`.

[45] B. D. MacArthur, C. P. Please, and R. O. C. Oreffo. Stochasticity and the molecular mechanisms of induced pluripotency. *PLoS One*, 3(8):e3086, 2008. doi: 10.1371/journal.pone.0003086.

[46] S. J. Ridden, et al. Entropy, Ergodicity, and Stem Cell Multipotency. *Phys. Rev. Lett.*, 115(20):208103, nov 2015. doi: 10.1103/PhysRevLett.115.208103. URL `https://link.aps.org/doi/10.1103/PhysRevLett.115.208103`.

[47] A. Trounson and N. D. DeWitt. Pluripotent stem cells progressing to the clinic. *Nat. Rev. Mol. Cell Biol.*, 17(3):194–200, feb 2016. doi: 10.1038/nrm.2016.10. URL `http://www.nature.com/doifinder/10.1038/nrm.2016.10`.

[48] B. D. MacArthur, A. Ma'ayan, and I. R. Lemischka. Systems biology of stem cell fate and cellular reprogramming. *Nat. Rev. Mol. Cell Biol.*, 10:672–681, sep 2009. doi: 10.1038/nrm2766. URL `http://www.nature.com/doifinder/10.1038/nrm2766`.

[49] M. Herberg and I. Roeder. Computational modelling of embryonic stem-cell fate control. *Development*, 142:2250–2260, 2015. doi: 10.1242/dev.116343.

[50] M. Santillán, et al. Modeling the epigenetic attractors landscape: toward a post-genomic mechanistic understanding of development. *Front. Genet.*, 6(160):1–14, 2015. doi: 10.3389/fgene.2015.00160.

[51] B. Ebrahimi. Biological computational approaches: new hopes to improve (re)programming robustness, regenerative medicine and cancer therapeutics. *Differentiation*, pp. 1–6, 2016. doi: 10.1016/j.diff.2016.03.001. URL `http://dx.doi.org/10.1016/j.diff.2016.03.001`.

[52] K. Takahashi and S. Yamanaka. A decade of transcription factor-mediated reprogramming to pluripotency. *Nat. Rev. Mol. Cell Biol.*, 17:183–193, 2016. doi: 10.1038/nrm.2016.8.

[53] C. R. S. Banerji, et al. Cellular network entropy as the energy potential in Waddington's differentiation landscape. *Sci. Rep.*, 2013. doi: 10.1038/srep03039.

[54] F. Nicol-Benoit, P. Le Goff, and D. Michel. Drawing a Waddington landscape to capture dynamic epigenetics. *Biol. Cell*, 105:576–584, 2013. doi: 10.1111/boc.201300029.

[55] J. Wang, et al. The Potential Landscape of Genetic Circuits Imposes the Arrow of Time in Stem Cell Differentiation. *Biophys. J.*, 99:29–39, 2010. doi: 10.1016/j.bpj.2010.03.058. URL `http://ac.els-cdn.com/ S0006349510004248/1-s2.0-S0006349510004248-main.pdf?{_}tid= 30768ffa-0324-11e7-9e09-00000aab0f02{&}acdnat= 1488884142{_}519d660a7d1d75b407c579bfd52e21b7`.

[56] O. J. L Rackham, et al. A predictive computational framework for direct reprogramming between human cell types. *Nat. Genet.*, 48(3):331–335, 2016. doi: 10.1038/ng.3487.

[57] J. L. Fierst and P. C. Phillips. Modeling the Evolution of Complex Genetic Systems: The Gene Network Family Tree. *J. Exp. Zool. (Mol. Dev. Evol.)*, 324(B):1–12, 2015. doi:

10.1002/jez.b.22597. URL

`https://onlinelibrary.wiley.com/doi/pdf/10.1002/jez.b.22597`.

[58] S. G. Brush. History of the Lenz-Ising Model. *Rev. Mod. Phys.*, 39(4):883–893, oct 1967.

doi: 10.1103/RevModPhys.39.883. URL

`https://link.aps.org/doi/10.1103/RevModPhys.39.883`.

[59] E. Ising. Beitrag zur Theorie des Ferromagnetismus. *Z.Phys*, 31:253–258, 1925. URL

`https:`

`//link.springer.com/content/pdf/10.1007{%}2FBF02980577.pdf`.

[60] W. S. Mcculloch and W. Pitts. A Logical Calculus of the Ideas Immanent in Nervous

Activity. *Bull. Math. Biophys.*, 5:115–133, 1943. URL `https:`

`//link.springer.com/content/pdf/10.1007{%}2FBF02478259.pdf`.

[61] J. J. Tyson, T. Laomettachit, and P. Kraikivski. Modeling the dynamic behavior of

biochemical regulatory networks. *J. Theor. Biol.*, 462:514–527, feb 2019. doi:

10.1016/j.jtbi.2018.11.034. URL

`https://linkinghub.elsevier.com/retrieve/pii/S0022519318305873`.

[62] T. Castellani and A. Cavagna. Spin-glass theory for pedestrians. *J. Stat. Mech. Theory

Exp.*, P05012, 2005. doi: 10.1088/1742-5468/2005/05/P05012. URL

`http://arxiv.org/abs/cond-mat/0505032{%}5Cnhttp:`

`//stacks.iop.org/1742-5468/2005/i=05/a=P05012?key=crossref.`

`65b9c18afe91f699a0b084ed4cfb1a81`.

[63] S. F. Edwards and P. W. Anderson. Theory of spin glasses. *J. Phys. F Met. Phys.*, 5:965,

1975. URL `https:`

`//iopscience.iop.org/article/10.1088/0305-4608/5/5/017/pdf`.

[64] D. Sherrington and S. Kirkpatrick. Solvable Model of a Spin-Glass. *Phys. Rev. Lett.*, 35(26):1792–1796, dec 1975. doi: 10.1103/PhysRevLett.35.1792. URL `https://link.aps.org/doi/10.1103/PhysRevLett.35.1792`.

[65] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Spin-glass models of neural networks. *Phys. Rev. A*, 32(2):1007–1018, 1985. doi: 10.1103/PhysRevA.32.1007. URL `http://pra.aps.org/abstract/PRA/v32/i2/p1007{_}1`.

[66] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci.*, 79(8):2554–2558, apr 1982. URL `http://www.ncbi.nlm.nih.gov/pubmed/6953413http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC346238`.

[67] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Storing Infinite Numbers of Patterns in a Spin-Glass Model of Neural Networks. *Phys. Rev. Lett.*, 55(14):1530–1533, sep 1985. doi: 10.1103/PhysRevLett.55.1530. URL `https://link.aps.org/doi/10.1103/PhysRevLett.55.1530`.

[68] D. J. Amit, H. Gutfreund, and H. Sompolinsky. Statistical Mechanics of Neural Networks near Saturation. *Ann. Phys. (N. Y).*, 173:30–67, 1987. URL `https://ac.els-cdn.com/0003491687900923/1-s2.0-0003491687900923-main.pdf?{_}tid=30a35faf-93e9-4ad8-ae93-55a96854f9b2{&}acdnat=1549300368{_}153945871a825f1dd566bd46cb3e2503`.

[69] A. Wagner. Evolution of gene networks by gene duplications: A mathematical model and its implications on genome organization. *Proc. Natl. Acad. Sci.*, 91:4387–4391, 1994. URL `https://www.jstor.org/stable/pdf/2364661.pdf?refreqid=excelsior{%}3Ab8999dc60d0a3e0116d9a443f96418b6`.

[70] G. P. Wagner and L. Altenberg. Perspective: Complex Adaptations and the Evolution of

Evolvability. *Evolution (N. Y).*, 50(3):967–976, jun 1996. doi:

10.1111/j.1558-5646.1996.tb02339.x. URL

`http://doi.wiley.com/10.1111/j.1558-5646.1996.tb02339.x`.

[71] M. L. Siegal and A. Bergman. Waddington's canalization revisited: Developmental

stability and evolution. *PNAS*, 99(16):10528–10532, 2002. URL

`www.pnas.orgcgidoi10.1073pnas.102303999`.

[72] A. Bergman and M. Siegal. Evolutionary capacitance as a general feature of complex

gene networks. *Nature*, 424, 2003. doi: 10.1038/nature01765. URL

`www.nature.com/nature`.

[73] R. B. R. Azevedo, et al. Sexual reproduction selects for robustness and negative epistasis

in artificial gene networks. *Nat. Lett.*, 440(7080):87–90, mar 2006. doi:

10.1038/nature04488. URL `http://www.nature.com/articles/nature04488`.

[74] T. MacCarthy and A. Bergman. Coevolution of robustness, epistasis, and recombination

favors asexual reproduction. *Proc. Natl. Acad. Sci.*, 104(31):12801–12806, aug 2007.

doi: 10.1073/pnas.77.8.4838. URL

`http://www.ncbi.nlm.nih.gov/pubmed/16592864http:`

`//www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC349943`.

[75] E. Borenstein and D. C. Krakauer. An End to Endless Forms: Epistasis, Phenotype

Distribution Bias, and Nonuniform Evolution. *PLoS Comput. Biol.*, 4(10):e1000202, oct

2008. doi: 10.1371/journal.pcbi.1000202. URL

`https://dx.plos.org/10.1371/journal.pcbi.1000202`.

[76] J. Draghi and G. P. Wagner. The evolutionary dynamics of evolvability in a gene network

model. *J. Evol. Biol.*, 22(3):599–611, mar 2009. doi:

10.1111/j.1420-9101.2008.01663.x. URL

http://doi.wiley.com/10.1111/j.1420-9101.2008.01663.x.

[77] J. L. Fierst. A history of phenotypic plasticity accelerates adaptation to a new

environment. *J. Evol. Biol.*, 24(9):1992–2001, sep 2011. doi:

10.1111/j.1420-9101.2011.02333.x. URL

http://doi.wiley.com/10.1111/j.1420-9101.2011.02333.x.

[78] C. Espinosa-Soto, O. C. Martin, and A. Wagner. Phenotypic plasticity can facilitate

adaptive evolution in gene regulatory circuits. *BMC Evol. Biol.*, 11(5):1–14, dec 2011.

doi: 10.1186/1471-2148-11-5. URL http:

//bmcevolbiol.biomedcentral.com/articles/10.1186/1471-2148-11-5.

[79] Y. Le Cunff and K. Pakdaman. Phenotype–genotype relation in Wagner's canalization

model. *J. Theor. Biol.*, 314:69–83, dec 2012. doi: 10.1016/J.JTBI.2012.08.020. URL

https://www.sciencedirect.com/science/article/pii/

S0022519312004420?via{%}3Dihub.

[80] B. Derrida and H. Flyvbjerg. Multivalley structure in Kauffman's model: analogy with

spin glasses. *J. Phys. A. Math. Gen.*, 19(16):L1003–L1008, 1986. doi:

10.1088/0305-4470/19/16/010. URL

http://iopscience.iop.org/0305-4470/19/16/010{%}5Cnhttp:

//iopscience.iop.org/0305-4470/19/16/010/pdf/

0305-4470{_}19{_}16{_}010.pdf.

[81] A. H. Lang, et al. Epigenetic landscapes explain partially reprogrammed cells and

identify key reprogramming genes. *PLoS Comput. Biol.*, 10(8):1–13, 2014. doi:

10.1371/journal.pcbi.1003734.

[82] A. Szedlak, et al. Cell cycle time series gene expression data encoded as cyclic attractors

in Hopfield systems. *PLOS Comput. Biol.*, 13(11):e1005849, nov 2017. doi: 10.1371/journal.pcbi.1005849. URL http://dx.plos.org/10.1371/journal.pcbi.1005849.

[83] A. Szedlak, G. Paternostro, and C. Piermarocchi. Control of asymmetric Hopfield networks and application to cancer attractors. *PLoS One*, 9(8):e105842, 2014. doi: 10.1371/journal.pone.0105842. URL http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0105842.

[84] T. Vierbuchen and M. Wernig. Molecular Roadblocks for Cellular Reprogramming. *Mol. Cell*, 47(6):827–838, 2012. doi: 10.1016/j.molcel.2012.09.008. URL http://dx.doi.org/10.1016/j.molcel.2012.09.008.

[85] S. F. Gilbert. Epigenetic Landscaping: Waddington's Use of Cell Fate Bifurcation Diagrams. *Biol. Philos.*, 6:135–154, 1991.

[86] P. Mazzarello. A unifying concept: the history of cell theory. *Nat. Cell Biol.*, 1:E13–E15, 1999.

[87] K. A. Johnson and R. S. Goody. The Original Michaelis Constant: Translation of the 1913 Michaelis– Menten Paper. *Biochemistry*, 50:8264–8269, 2011. doi: 10.1021/bi201284u.

[88] S. Bös, R. Kühn, and J. L. van Hemmen. Martingale approach to neural networks with hierarchically structured information. *Zeitschrift für Phys. B Condens. Matter*, 71(2):261–271, jun 1988. doi: 10.1007/BF01312798. URL http://link.springer.com/10.1007/BF01312798.

[89] R. J. Cho, et al. A Genome-Wide Transcriptional Analysis of the Mitotic Cell Cycle. *Mol. Cell*, 2:65–73, jul 1998. doi: 10.1016/S1097-2765(00)80114-8. URL http://linkinghub.elsevier.com/retrieve/pii/S1097276500801148.

[90] H. Sompolinsky and I. Kanter. Temporal Association in Asymmetric Neural Networks. *Phys. Rev. Lett.*, 57(22):2861–2864, 1986.

[91] H. Gutfreund and M. Mezard. Processing of Temporal Sequences in Neural Networks. *Phys. Rev. Lett.*, 61(2):235–238, 1988. doi: 10.+e.

[92] C. Cortes, A. Krogh, and J. A. Hertz. Hierarchical associative networks. *J. Phys. A. Math. Gen.*, 20(13):4449–4455, sep 1987. doi: 10.1088/0305-4470/20/13/044. URL `http://stacks.iop.org/0305-4470/20/i=13/a=044?key=crossref.80669a1daea667a0b36926ccf44f6baa`.

[93] A. Krogh and J. A. Hertz. Mean-field analysis of hierarchical associative networks with 'magnetisation'. *J. Phys. A Math. Gen*, 21:2211–2224, 1988. URL `http://iopscience.iop.org/0305-4470/21/9/033`.

[94] N. Parga and M. A. Virasoro. The ultrametric organization of memories in a neural network. *J. Phys.*, 47(11):1857–1864, 1986. doi: 10.1051/jphys:01986004701101857000>. URL `https://hal.archives-ouvertes.fr/jpa-00210382`.

[95] D. Bollé, J. Busquets Blanco, and T. Verbeiren. The signal-to-noise analysis of the Little-Hopfield model revisited. *J. Phys. A Math. Gen. J. Phys. A Math. Gen*, 37:1951–1969, 2004. doi: 10.1088/0305-4470/37/6/001. URL `http://iopscience.iop.org/article/10.1088/0305-4470/37/6/001/pdf`.

[96] S. M. Chambers, et al. Hematopoietic Fingerprints: An Expression Database of Stem Cells and Their Progeny. *Cell Stem Cell*, 1:578–591, 2007. doi: 10.1016/j.stem.2007.10.003.

[97] E. S. Lander and et Al. Initial sequencing and analysis of the human genome. *Nature*, 409:860–921, 2001.

[98] J. C. Venter, et al. The Sequence of the Human Genome. *Science (80-. ).*, 291:1304–1351, 2001.

[99] J. M. Vaquerizas, et al. A census of human transcription factors: function, expression and evolution. *Nat. Rev. Genet.*, 10(4):252–263, 2009. doi: 10.1038/nrg2538. URL `http://www.ncbi.nlm.nih.gov/pubmed/19274049`.

[100] V. Narang, et al. Automated Identification of Core Regulatory Genes in Human Gene Regulatory Networks. *PLoS Comput. Biol.*, 11(9):e1004504, 2015. doi: 10.1371/journal.pcbi.1004504. URL `http://www.ncbi.nlm.nih.gov/pubmed/26393364http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4578944`.

[101] A. C. C. Coolen, R. Kuehn, and P. Sollich. Theory of Neural Information Processing Systems. Oxford University Press, 2005.

[102] M. Uhlén, et al. Tissue-based map of the human proteome. *Science (80-. ).*, 347(6220), 2015. doi: 10.1126/science.1260419. URL `http://science.sciencemag.org/`.

[103] The Human Protein Atlas. URL `https://www.proteinatlas.org`.

[104] R. Lasrado, et al. Lineage-dependent spatial and functional organization of the mammalian enteric nervous system. *Science (80-. ).*, 356:722–726, 2017. URL `http://science.sciencemag.org/content/sci/356/6339/722.full.pdf`.

[105] A. T. Lun, D. J. McCarthy, and J. C. Marioni. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Research*, 5(2122):1–68, oct 2016. doi: 10.12688/f1000research.9501.2. URL `https://f1000research.com/articles/5-2122/v2`.

[106] I. Kanter and H. Sompolinsky. Associative recall of memory without errors. *Phys. Rev. A,*

35(1), 1987. URL https://pdfs.semanticscholar.org/f3e2/
cfad939b7e4f0ea7783883afa997e48920c1.pdf.

[107] C. Zhang, G. Dangelmayr, and I. Oprea. Storing cycles in Hopfield-type networks with
pseudoinverse learning rule: Admissibility and network topology. *Neural Networks*,
46:283–298, 2013. doi: 10.1016/j.neunet.2013.06.008. URL
http://dx.doi.org/10.1016/j.neunet.2013.06.008.

[108] L. Narlikar and I. Ovcharenko. Identifying regulatory elements in eukaryotic genomes.
*Briefings Funct. genomics proteomics*, 8(4):215–230, 2009. doi: 10.1093/bfgp/elp014.
URL https:
//www.ncbi.nlm.nih.gov/pmc/articles/PMC2764519/pdf/elp014.pdf.

[109] D. Kleftogiannis, et al. Where we stand, where we are moving: Surveying computational
techniques for identifying miRNA genes and uncovering their regulatory role. *J. Biomed.
Inform.*, 46:563–573, 2013. doi: 10.1016/j.jbi.2013.02.002. URL
http://dx.doi.org/10.1016/j.jbi.2013.02.002.

[110] T. Anton, E. Karg, and S. Bultmann. Applications of the CRISPR/Cas system beyond
gene editing. *Biol. Methods Protoc.*, 3(1):1–10, 2018. doi:
10.1093/biomethods/bpy002. URL https://academic.oup.com/biomethods/
article-abstract/3/1/bpy002/4995153.

[111] E. Schenborn and D. Groskreutz. Reporter Gene Vectors and Assays. *Mol. Biotechnol.*,
13(1):29–44, 1999. URL https://link.springer.com/content/pdf/10.
1385{%}2FMB{%}3A13{%}3A1{%}3A29.pdf.

[112] L. En Chai, et al. A review on the computational approaches for gene regulatory network
construction. *Comput. Biol. Med.*, 48:55–65, 2014. doi:

10.1016/j.compbiomed.2014.02.011. URL

`http://dx.doi.org/10.1016/j.compbiomed.2014.02.011`.

[113] D. Thompson, A. Regev, and S. Roy. Comparative Analysis of Gene Regulatory
Networks: From Network Reconstruction to Evolution. *Annu. Rev. Cell Dev. Biol,*
31:399–428, 2015. doi: 10.1146/annurev-cellbio-100913-012908. URL
`www.annualreviews.org`.

[114] R. Hannam, A. Annibale, and R. Kühn. Cell reprogramming modelled as transitions in a
hierarchy of cell cycles. *J. Phys. A Math. Theor.,* 50(425601):1–23, 2017. doi:
10.1088/1751-8121/aa89a2.

[115] D. V. Foster, et al. A model of sequential branching in hierarchical cell fate determination.
*J. Theor. Biol.,* 260:589–597, 2009. doi: 10.1016/j.jtbi.2009.07.005. URL `http://ac.`
`els-cdn.com/S002251930900318X/1-s2.0-S002251930900318X-main.`
`pdf?{_}tid=5941f78a-0324-11e7-ba01-00000aacb35f{&}acdnat=`
`1488884211{_}ddcd643580f0557697100c7f1c94bc27`.

[116] J. X. Zhou and S. Huang. Understanding gene circuits at cell-fate branch points for
rational cell reprogramming. *Trends Genet.,* 27(2):55–62, 2010. doi:
10.1016/j.tig.2010.11.002. URL `http://ac.els-cdn.com/`
`S0168952510002222/1-s2.0-S0168952510002222-main.pdf?{_}tid=`
`cf7b000a-0323-11e7-90c3-00000aacb362{&}acdnat=`
`1488883979{_}7253e6b88ab4f5a6619d655d96004121`.

[117] X. Fang, et al. Cell fate potentials and switching kinetics uncovered in a classic bistable
genetic switch. *Nat. Commun.,* 9(2787):1–9, 2018. doi: 10.1038/s41467-018-05071-1.
URL `www.nature.com/naturecommunications`.

[118] C. Li and J. Wang. Quantifying Cell Fate Decisions for Differentiation and

Reprogramming of a Human Stem Cell Network: Landscape and Biological Paths. *PLoS Comput. Biol.*, 9(8):e1003165, 2013. doi: 10.1371/journal.pcbi.1003165. URL `www.ploscompbiol.org`.

[119] S.-J. Dunn, et al. Defining an essential transcription factor program for naïve pluripotency. *Science (80-. ).*, 344(6188):1156–1160, jun 2014. doi: 10.1126/science.1248882. URL `http://www.ncbi.nlm.nih.gov/pubmed/24904165{%}5Cnhttp://www.sciencemag.org/cgi/doi/10.1126/science.1248882`.

[120] G. A. Pavlopoulos, et al. Using graph theory to analyze biological networks. *BioData Min.*, 4(10):1–27, 2011. doi: 10.1186/1756-0381-4-10. URL `https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3101653/pdf/1756-0381-4-10.pdf`.

[121] M. Boguñá and M. Á. Serrano. Generalized percolation in random directed networks. *Phys. Rev. E*, 72(1):016106, jul 2005. doi: 10.1103/PhysRevE.72.016106. URL `https://link.aps.org/doi/10.1103/PhysRevE.72.016106`.

[122] H. Hooyberghs, B. Van Schaeybroeck, and J. Indekeu. Percolation on bipartite scale-free networks. *Phys. A Stat. Mech. its Appl.*, 389(15):2920–2929, aug 2010. doi: 10.1016/J.PHYSA.2009.12.068. URL `https://www.sciencedirect.com/science/article/pii/S0378437110000336?via{%}3Dihub`.

[123] M. Mézard and G. Parisi. The Bethe lattice spin glass revisited. *Eur. Phys. J. B*, 20(2):217–233, mar 2001. doi: 10.1007/PL00011099. URL `http://link.springer.com/10.1007/PL00011099`.

[124] P. M. Kogut and P. L. Leath. Bootstrap percolation transitions on real lattices. *J. Phys. C Solid State Phys*, 14:3187–3194, 1981. URL `http://iopscience.iop.org/article/10.1088/0022-3719/14/22/013/pdf`.

[125] M. Aizenman and J. L. Lebowitz. Metastability effects in bootstrap percolation. *J. Phys. A. Math. Gen.*, 21:3801–3813, 1988. URL `http://iopscience.iop.org/article/10.1088/0305-4470/21/19/017/pdf`.

[126] R. H. Schonmann. On the Behavior of Some Cellular Automata Related to Bootstrap Percolation. *Ann. Probab.*, 20(1):174–193, 1992. URL `https://www.jstor.org/stable/pdf/2244552.pdf?refreqid=excelsior{%}3Ab751edc32f8956b09dbe194f413b4c0c`.

[127] J. Gravner and A. E. Holroyd. Slow Convergence in Bootstrap Percolation. *Ann. Appl. Probab.*, 18(3):909–928, 2008. doi: 10.1214/07-AAP473. URL `https://www.jstor.org/stable/pdf/25442654.pdf?refreqid=excelsior{%}3Aae42600dc3ad6ba0744f4cdee4e93dfe`.

[128] J. Balogh and B. Bollobás. Bootstrap percolation on the hypercube. *Probab. Theory Relat. Fields*, 134:624–648, 2006. doi: 10.1007/s00440-005-0451-6. URL `https://link.springer.com/content/pdf/10.1007{%}2Fs00440-005-0451-6.pdf`.

[129] J. Balogh and B. G. Pittel. Bootstrap Percolation on the Random Regular Graph. *Random Struct. Alg*, 30:257–286, 2006. doi: 10.1002/rsa.20158. URL `www.interscience.wiley.com`.

[130] L. R. G. Fontes and ·. R. H. Schonmann. Bootstrap Percolation on Homogeneous Trees Has 2 Phase Transitions. *J Stat Phys*, 132:839–861, 2008. doi: 10.1007/s10955-008-9583-2. URL `https://link.springer.com/content/pdf/10.1007{%}2Fs10955-008-9583-2.pdf`.

[131] J. Balogh, Y. Peres, and G. Pete. Bootstrap Percolation on Infinite Trees and Non-Amenable Groups. *Comb. Probab. Comput.*, 15(05):715, sep 2006. doi:

10.1017/S0963548306007619. URL

`http://www.journals.cambridge.org/abstract{_}S0963548306007619`.

[132] G. J. Baxter, et al. Bootstrap percolation on complex networks. *Phys. Rev. E,* 82(1):011103, jul 2010. doi: 10.1103/PhysRevE.82.011103. URL

`https://link.aps.org/doi/10.1103/PhysRevE.82.011103`.

[133] J. Gao, T. Zhou, and Y. Hu. Bootstrap percolation on spatial networks. *Sci. Rep.,* 5(14662):1–10, dec 2015. doi: 10.1038/srep14662. URL

`http://www.nature.com/articles/srep14662`.

[134] A. Annibale, A. C. C. Coolen, and N. Planell-Morell. Quantifying noise in mass spectrometry and yeast two-hybrid protein interaction detection experiments. *J. R. Soc. Interface,* 12(20150573):1–29, 2015. doi: 10.1098/rsif.2015.0573. URL

`http://dx.doi.org/10.1098/rsif.2015.0573`.

[135] V. Folli, et al. Effect of dilution in asymmetric recurrent neural networks. *Neural Networks,* 104:50–59, aug 2018. doi: 10.1016/j.neunet.2018.04.003. URL

`https://linkinghub.elsevier.com/retrieve/pii/S0893608018301230`.

[136] S. Orkin and K. Hochedlinger. Chromatin Connections to Pluripotency and Cellular Reprogramming. *Cell,* 145(6):835–850, jun 2011. doi: 10.1016/J.CELL.2011.05.019. URL `https: //www.sciencedirect.com/science/article/pii/S0092867411005769`.

[137] S. A. Morris and G. Q. Daley. A blueprint for engineering cell fate: current technologies to reprogram cell identity. *Cell Res.,* 23(1):33–48, jan 2013. doi: 10.1038/cr.2013.1. URL `http://www.nature.com/articles/cr20131`.

[138] L. David and J. M. Polo. Phases of reprogramming. *Stem Cell Res.,* 12:754–761, 2014. doi: 10.1016/j.scr.2014.03.007.

[139] B. Nashun, P. W. Hill, and P. Hajkova. Reprogramming of cell fate: epigenetic memory and the erasure of memories past. *EMBO J.*, 34:1296–1308, 2015. doi: 10.15252/embj.201490649. URL

`http://emboj.embopress.org/content/embojnl/34/10/1296.full.pdf.`