# UNIVERSITY OF Southampton

University of Southampton Research Repository
ePrints Soton

# UNIVERSITY OF SOUTHAMPTON

## FACULTY OF SOCIAL, HUMAN AND MATHEMATICAL SCIENCES

# Stochastic Models of Stem Cell Dynamics

by

**Sonya Jane Ridden**

A thesis for the degree of Doctor of Philosophy

July 2016

UNIVERSITY OF SOUTHAMPTON

ABSTRACT

FACULTY OF SOCIAL, HUMAN AND MATHEMATICAL SCIENCES

Mathematics

Doctor of Philosophy

STOCHASTIC MODELS OF STEM CELL DYNAMICS

by Sonya Jane Ridden

There is a growing body of evidence to suggest that stem cell populations from both the embryo and the adult are heterogeneous in their gene expression patterns. However, the underlying mechanisms are not well understood. This thesis explores cell-to-cell variability in both multipotent and pluripotent stem cell populations using mathematical models to provide a theoretical framework to understand the collective dynamics of stem cell populations.

In the first part of the thesis we investigate the possibility that fluctuations in the transcription factor Nanog – which is central to the embryonic stem cell transcriptional regulatory network (ESC TRN) – regulate population variability by controlling important feedback mechanisms. Our analyses reveal the ESC TRN is rich in feedback, with global feedback structure critically dependent on Nanog, Oct4 and Sox2, which collectively participate in over two thirds of all feedback loops. Using a general measure of feedback centrality we show that removal of Nanog severely compromises the global feedback structure of the ESC TRN. These analyses indicate that Nanog fluctuations regulate population heterogeneity by transiently activating different regulatory subnetworks, driving transitions between a Nanog-expressing, feedback-rich, robust and self-perpetuating pluripotent state and a Nanog-diminished, feedback-sparse and differentiation-sensitive state.

The majority of studies characterising heterogeneity in Nanog expression have used live-cell fluorescent reporter strategies. However, recent evidence suggests that these reporters may not give a faithful reflection of endogenous Nanog expression because the introduction of the reporter construct can perturb the kinetics of the underlying regulatory network. To investigate the role of Nanog further we therefore sought to model in detail the dynamics of Nanog expression in heterozygous fluorescent knock-in reporter cell lines. We develop chemical master equation, chemical Langevin equation and reaction rate equation models of the reporter system to determine how this might disturb normal Nanog transcriptional control. Our analyses indicate that the reporter construct can weaken the strength of autoactivatory feedback loops that are central to Nanog regulation, and thereby qualitatively perturbs endogenous Nanog dynamics. These results question the efficacy of commonly used reporter strategies and therefore have important implications for the design and use of synthetic reporters in general, not just for Nanog.

In the second part of this thesis we consider the dynamics of populations of multipotent adult hematopoietic stem cells (HSCs). It is known that fluctuations within individual HSCs allow them to transit stochastically between functionally distinct metastable states, while the overall population distribution of expression remains stable. To investigate the relationship between single cell and population-level dynamics we propose a theoretical framework that views cellular multipotency as an instance of maximum entropy statistical inference, in which an underlying ergodic stochastic process gives rise to robust variability within the cell population. We illustrate this view by analysing expression fluctuations of the stem cell surface marker Sca1 in mouse HSCs and find that the observed dynamics naturally lie close to a critical state, thereby producing a diverse population that is able to respond rapidly to environmental changes. Although we focus on Sca1 dynamics, comparable expression fluctuations are known to generate functional diversity in other mammalian stem cell systems, including in pluripotent stem cells. Thus, the generation of ergodic expression fluctuations may be a generic way in which cell populations maintain robust multilineage differentiation potential under environmental uncertainty.

# Contents

# List of Tables

# List of Figures

# Declaration of Authorship

I, Sonya Jane Ridden, declare that the thesis entitled *Stochastic Models of Stem Cell Dynamics* and the work presented in the thesis are both my own, and have been generated by me as the result of my own original research. I confirm that:

- this work was done wholly while in candidature for a research degree at this University;

- where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated;

- where I have consulted the published work of others, this is always clearly attributed;

- where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work;

- I have acknowledged all main sources of help;

- where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself;

- parts of this work have been published as:

  [1] Sonya J. Ridden, Hannah H. Chang, Konstantinos C. Zygalakis, and Ben D. MacArthur. Entropy, ergodicity, and stem cell multipotency. *Physical Review Letters*, 115:208103, 2015

  [2] Sonya J. Ridden and Ben D. MacArthur. Chapter 2: Cell Fate Regulatory Networks, in *New frontiers of network analysis in systems biology*. Springer, New York, 2012

  [3] Ben D. MacArthur, Ana Sevilla, Michel Lenz, Franz-Josef Muller, Berhard M. Schuldt, Andreas A. Schuppert, Sonya J. Ridden, Patrick S. Stumpf, Miguel Fidalgo, Avi Ma'ayan, Jianlong Wang and Ihor R. Lemischka. Nanog-dependent feedback loops regulate murine embryonic stem cell heterogeneity. *Nature Cell Biology*, 14(11):1139 1147, 2012

  [4] Rosanna C. G. Smith, Patrick S. Stumpf, Sonya J. Ridden, Aaron Sim, Sarah Filippi, Heather Harrington, Ben D. MacArthur. Nanog fluctuations in ES cells highlight the problem of measurement in cell biology. *bioRxiv*, page 060558, 2016

Signed:........................................................ Date:...............................................

# Acknowledgements

I would like to thank my supervisors Dr. Ben MacArthur and Dr. Konstantinos (Kostas) Zygalakis for their relentless, active support, without which I could not have completed this thesis. In particular, Ben has provided an essential form of stability, and his innate ability to quell any sign of waning enthusiasm on my part had a positive effect on my momentum. Moreover, he has given me the priceless opportunity to contribute to respected academic publications. I consider myself very lucky to have had both Ben and Kostas as my supervisors.

I would also like to thank the Institute for Complex Systems Simulation for their financial and emotional support, especially during the more uncomfortable periods of my life.

Finally, I cannot forget my lovely mum, whose unconditional love, and tea-to-desk delivery services have been vital to my survival (even though she has the tendency to put the milk in before removing the tea bag).

# Nomenclature

$\boldsymbol{x}(t)$       state vector

$M$       number of chemical reactions

$N$       number of molecular species

$R_j$       $j$th chemical reaction

$S_i$       $i$th molecular species

$\boldsymbol{\nu}_j$       state-change vector for reaction $j$

$a_j(\boldsymbol{x})$       propensity function for reaction $j$

$k_j$       reaction rate constant for reaction $j$

$p(\boldsymbol{x}, t)$       probability that the system is in state $\boldsymbol{x}$ at time $t$

$\xi(t)$       Gaussian white noise

$\boldsymbol{F}(\boldsymbol{x})$       $N \times 1$ drift vector

$\boldsymbol{S}$       $N \times M$ stoichiometric matrix

$\boldsymbol{a}(\boldsymbol{x})$       $M \times 1$ column vector of propensity functions

$\boldsymbol{D}(\boldsymbol{x})$       $N \times N$ diffusion matrix

$F(x)$       scalar drift function

$\sigma(x)$       scalar diffusion function

$J(x, t)$       probability current

$p_\infty(x)$       stationary probability density function

$\boldsymbol{A}$       network adjacency matrix

$\psi(x)$       scalar potential

CFPE       chemical Fokker-Planck equation

CLE       chemical Langevin equation

CME       chemical master equation

DNA       deoxyribonucleic acid

EML       erythroid, myeloid, and lymphocytic mouse hematopoietic progenitor cells

FI        fluorescence intensity

FPT       first passage time

ESC       embryonic stem cell

GFP       green fluorescent protein

HSC       hematopoietic stem cell

MFPT      mean first passage time

mRNA      messenger ribonucleic acid

ODE       ordinary differential equation

PPIN      protein-protein interaction network

RNAP      ribonucleic acid polymerase

RRE       reaction rate equation

Sca1      stem cell antigen 1 - surface marker in mouse hematopoietic progenitor cells

SDE       stochastic differential equation

SSA       stochastic simulation algorithm

TF        transcription factor

TRN       transcriptional regulatory network

# Chapter 1

# Introduction

It is now well documented that temporal gene expression fluctuations at the single-cell level give rise to heterogeneity in clonal stem cell populations [3], but the functional role of this diversity and the underlying molecular mechanisms are not well understood. In this thesis, we use mathematical models describing the regulation of gene expression to explore cell-to-cell variability in both adult and embryonic stem cell populations.

An important example of expression variability in embryonic stem cells is given by the protein Nanog, which is a key component of the mammalian embryonic stem cell regulatory network - a network of interacting proteins that governs embryonic stem cell fate decisions by influencing the expression of many genes associated with specific cell types. In Chapter 2, we explore models of gene expression in the core regulatory network to elucidate the underlying mechanisms that regulate variability. In Section 2.2, we investigate the role that feedback loops and the inherent stochasticity in transcription and translation have in controlling the variability of protein expression levels at the population level, and cell fate decisions at the single cell level. In Section 2.3, we focus on the experimentally observed heterogeneity of Nanog expression. Although many studies report a bimodal distribution of expression, new evidence suggests that this heterogeneity might be a result of the measurement method - a heterozygous knock-in reporter - interfering with the Nanog autoregulatory feedback loop, as opposed to a true reflection of the distribution of Nanog expression. We present a stochastic model of positive feedback to demonstrate how Nanog expression could be perturbed by the heterozygous knock-in strategy, thus supporting the notion that the observed bimodality is a reporter artefact.

In Chapter 3, we use tools from statistical mechanics and information theory to propose a theoretical framework for the functional role of the considerable cell-cell variability commonly exhibited by adult stem cells. In order to illustrate our perspective, we use a simple stochastic model to analyse the dynamics of a single protein in blood forming stem cells *in vitro*.

We begin by providing the biological background required to understand the systems modelled and analysed in this thesis. Starting with an introduction to genes and their expression, we then explain the role of transcription factors in the complex process of gene expression. Following an overview of cell types, we introduce genetic regulatory networks and their underlying biological processes.

In Section 1.2 of this chapter, we present the mathematical background common to all studies in this thesis. A review of the historical understanding of cell fate determination provides the foundation for the modelling methods used in the subsequent chapters. We then explain how mathematical models can provide the basis for theoretical description of genetic regulatory networks, and offer a deeper insight into the mechanisms that control heterogeneity in gene expression. Starting with the chemical master equation, we demonstrate how the progressive introduction of assumptions enables the transition from detailed stochastic models to deterministic reaction rate equations. Finally, in Section 1.3, we discuss noise in gene expression and highlight the importance of accounting for the stochasticity in this process in models of regulatory networks.

## 1.1  Molecular Biology

### 1.1.1  DNA and Genes

*Deoxyribonucleic acid (DNA)* is a molecule that contains information used to create proteins that are required for the development, function and reproduction of all known living organisms. DNA is stored in the cell nucleus, and consists of two long chains of small molecules called *nucleotides* that coil around each other to form a double helix. Each nucleotide is composed of one of the bases cytosine (C), guanine (G), adenine (A), or thymine (T), and other molecules that form the structure of the nucleotide. A *gene* is a section of DNA whose sequence of bases dictates the linear sequence of amino acids in the corresponding protein.

### 1.1.2 Gene Expression

Gene expression is the process of creating a protein from the information given by the DNA sequence of a gene. Many of the tasks involved in this process are carried out by *molecular machines*, which are assemblies of molecules that perform specific mechanical motions in response to external stimuli. *RNA polymerase II (RNAP II)* is a molecular machine that reads a gene sequence and transcribes it into messenger RNA (mRNA) molecules. This step - shown in Fig. 1.1, top left - is known as *transcription*, and produces a portable copy of the gene sequence. Transcription starts at a region of the gene called the *promoter*. Transcription stops at a terminator site at which point both the DNA and the completed mRNA molecule are released from the RNAP II.

The mRNA molecule is then transported outside the nucleus into the cell cytoplasm. Inside the cytoplasm are large molecular machines known as *ribosomes*, which bind to mRNA and translate it into a string of amino acids. This step, known as *translation*, is illustrated in Fig. 1.1, bottom left.

The string of amino acids becomes a protein when it folds into a specific three-dimensional structure. A sequence of three mRNA nucleotides codes for a single unique amino acid, and a complete mRNA molecule provides the template for the amino acid sequence of the protein. The way in which the protein folds, and therefore its resulting shape, depends on how the amino acids that make up the sequence are attracted or repelled by each other, and whether or not they are repelled by water (*hydrophobic*) or attracted to it (*hydrophilic*). The function of the resulting protein depends on its shape and the amino acids that make up the exterior, as this determines which other proteins may interact with it and to what extent. Interaction between proteins occurs due to an attraction between regions of their exterior amino acid sequences, and complementary shapes enable them to fit together like a lock and key. The coordinated behaviour of these proteins determines the physiological properties of the cell.

### 1.1.3 Transcription Factors

Some proteins become part of the structure of the cell, and others catalyse chemical reactions, such as the breakdown of a food source or toxin. In this thesis we are particularly interested in certain proteins called transcription factors (TFs). These proteins can enter the nucleus and interact with DNA, altering the rate of transcription of genes. Each gene contains a cluster

Figure 1.1: An illustration of the gene expression process. Left: The transcription process (top) and the translation process (bottom). Right: Transcription factors can bind to the promoter region of a gene and alter the rate of transcription by attracting or repelling the RNA polymerase II.

of binding sites, known as the promoter region, whose DNA sequences determine which TFs can control the activation or repression of the gene. TFs can bind to these promoter regions because they contain an external sequence of amino acids that has an affinity for a particular section of the DNA sequence in the promoter region. This enables them to control the expression of genes by promoting or blocking the recruitment of RNAP II to the promoter, thus initiating or preventing gene expression. Each TF is a single protein but they can form multiprotein complexes. They may work alone at the target site or together with other transcription factors. Fig. 1.1 (right) shows an illustration of a DNA-bound complex of TFs attracting an RNAP II to the promoter region of a gene.

In addition to the action of TFs, gene expression is also regulated by epigenetic regulatory mechanisms, such as histone acetylation and DNA methylation [5, 6], and signalling networks [7]. Proteins called *histones* lack DNA binding sequences, but can package DNA by acting as spools around which the DNA is wound. Together, they make up chromatin, whose structure is controlled by histone modifications and DNA methylation. Methylation (the addition of a methyl group to a molecule) of DNA tightens chromatin, and the acetylation of histones loosens the bindings [8–11]. Since tightly coiled chromatin is not accessible to the TFs and transcriptional machinery, whereas regions of chromatin that are less condensed allow active transcription, structure modification provides an additional mechanism for controlling gene expression.

Studies show that TFs control the opening of chromatin at specific sites [12–15], which

suggests that TFs and chromatic modification work together to activate or silence genes. It was previously thought that gene silencing was irreversible, however recent studies have shown that that histone modifications are reversible and dynamic [16, 17], which means that, in principle, each gene can be reversibly switched on and off.

### 1.1.4 Cell Types, Stem Cells and Development

Almost every cell in a multicellular organism contains the same set of DNA, but development (of an organism) gives rise to a wide range of cell types, each with very different physical characteristics. These differences are largely due to differences in gene expression i.e. which genes are turned "on" and which are "off". Therefore, cell types are characterised by their distinct morphology or function, and by distinct patterns of mRNA or protein expression.

Mammalian cells can be divided into three basic types: germ cells, somatic cells and stem cells. Germ cells give rise to gametes (sex cells - eggs or sperm), and somatic cells make up most of the body e.g. skin and muscle cells. Stem cells can give rise to somatic cells through the process of development in the embryo [18], and tissue maintenance and repair in the adult [19, 20].

There are two major categories of mammalian stem cells [21]:

1. Embryonic stem cells (ESCs), which are derived from the inner cell mass of the blastocyst [22, 23], eventually give rise to the structures of the fetus [24]. ESCs are *pluripotent*, which means that they can differentiate into all of the specialised embryonic cell types, and they can replicate indefinitely [25].

2. Adult (or somatic) stem cells can be found in small numbers in most adult tissues. They repair tissues in adult organisms by replenishing specialised cells and maintaining regenerative organs during normal cell turnover [19, 20]. Unlike ESCs, adult stem cells are usually *multipotent* [26], which means that their specialisation potential is limited to one or more specific lineages.

ESCs are characterised by two properties:

1. They can renew themselves indefinitely in culture through mitotic cell division [21, 25]. A stem cell can proliferate by dividing while maintaining the undifferentiated state

[27]. Cell division results in two daughter cells, and stem cells can divide symmetrically and produce two stem cells identical to the original, or they can divide symmetrically into two differentiated daughter cells. They can also divide asymmetrically to give rise to two different daughter cells, where one is a copy of the original stem cell and the other is differentiated [28]. Distinct daughter cells can be created because of an uneven distribution of regulatory molecules in the parent cell; the distinct cytoplasm that each daughter cell inherits results in a distinct pattern of differentiation for each daughter cell [28, 29].

2. They have the capacity to differentiate into any adult cell type [30]. Cellular differentiation is the process by which stem cells produce increasingly specialised progeny [31], and is generally accompanied by coordinated changes in gene expression and alterations in DNA structure that affect transcriptional access [5, 6].

Mammalian development begins when a sperm fertilises an egg and creates a single cell, called a zygote. The zygote (and subsequent blastomeres) are able to differentiate into all cell types, including the placental tissue, and are therefore *totipotent* [32]. The zygote divides into identical cells and after several cycles of cell division, these cells begin to specialise, forming a hollow sphere of cells, called a blastocyst [33]. The blastocyst has an outer layer of cells, and attached to the inside wall of this outer layer is a cluster of cells called the inner cell mass [33]. The cells of the inner cell mass go on to form nearly all of the tissues of the body [34, 35].

During the early stages of mammalian embryonic development, ESCs can give rise to three different groups of cells called *germ-layers*, and each group can generate a distinct set of tissue lineages [36–38]. The three germ layers are

- the ectoderm, which develops into skin and neural tissue;

- the mesoderm, which gives rise to blood, bone, cartilage, muscle, and fat;

- the endoderm, which generates tissues of the respiratory and digestive tracts.

Adult stem cells typically belong to one of these three germ-layers, to which their regeneration potential has been considered to be restricted, although recent experiments have challenged this notion and suggested that under certain circumstances these cells may convert from one

tissue lineage into a cell of an entirely distinct lineage (transdifferentiate), to contribute to a much wider range of specialised cell than previously anticipated [38]. Somatic cells are referred to as *terminally differentiated* [39] in that they are specialised cells that have reached the final stage of development after which no further specialisation is thought to occur. Somatic cells are generated from adult stem cells, which are not terminally differentiated; the job of an adult stem cell is to produce cells that carry out a specialised function [19, 20]. A stem cell differentiates into a *progenitor cell*, which is more committed to a particular lineage than a stem cell. A progenitor cell can differentiate to one or more type of cell, but can divide only a limited number of times [40, 41]. Some examples of stem and progenitor cells are:

- Hematopoietic stem cells give rise to red blood cells, white blood cells, and platelets, for example, and can be found in adult bone marrow [42, 43]. In Chapter 3, we focus on the expression of the protein Sca1 in hematopoietic stem cells.

- Mesenchymal stem cells give rise to stromal cells, fat cells, and types of bone cells, for example, and can also be found in adult bone marrow [44].

- Epithelial cells are progenitor cells that give rise to cells that line hollow organs and glands, and also make up the outer surface of the body [45].

- Muscle satellite cells are progenitor cells that contribute to the generation of adult muscle tissue [46].

### 1.1.5 Genetic Regulatory Networks

In mathematical terms a regulatory network is a set of interconnected components, called *nodes* – which represent the molecular entities involved (generally genes and proteins) – along with connections between them called *links*. Links represent interactions between molecular components and can be *directed* or *undirected*. For instance, much attention has been paid to using high-throughput experimental techniques to identify physical protein-protein interactions [47–51]. These data allow the inference of protein-protein interaction networks (PPINs), using a combination of experimental methods [49, 52, 53] and reverse engineering by the use of computational techniques [54–57]. In PPINs the nodes represent proteins and links represent physical interactions between proteins (i.e. binding). In this case, the links have no specific orientation: if protein A interacts with protein B then B also interacts with A. PPINs are therefore undirected. Although PPINs map out the physical

interactions between proteins – and thus the *possible* protein complexes which may form – they do not incorporate the consequences of these physical interactions, such as the induction or repression of gene expression by multi-protein complexes.

These interaction effects can be represented in the form of a transcriptional regulatory network



Figure 1.2: A transcriptional regulatory network with three TFs. Promotion of transcription is represented by an arrow, and inhibition by a T-bar.

(TRN) and much attention has also been paid to determining the structure and function of TRNs [58–63]. In TRNs the nodes are TFs and the links represent regulation of gene expression by upstream TFs. Unlike physical interactions between proteins there is a definite orientation to transcriptional regulatory interactions: if TF A regulates the transcription of TF B, it is not necessarily true that B regulates the transcription of A. Consequently, transcriptional regulatory networks are directed. Figure 1.2 shows an example of a transcriptional regulatory network consisting of three TFs labelled A, B and C. Promotion of transcription is represented by an arrow, and inhibition by a T-bar, thus A promotes the transcription of B, B promotes the transcription of C, and C represses that of A.

## 1.2 Mathematical Models of Regulatory Networks

The inherent complexity of the TRN makes it impossible to determine cell behaviour from the regulatory network architecture using experiment and intuition alone [64]. In this section we explain how mathematical models can provide the basis for theoretical description of TRNs, and offer a deeper insight into the mechanisms that control expression levels.

Mathematical models of regulatory networks convert known structural information (in the form of experimentally-derived protein-protein and protein-DNA interactions, for instance) into a set of equations that describe how molecular expression levels change over time as a result of the interactions between the components [65]. These equations can then be solved in order to reproduce observed dynamics and make novel predictions concerning cell behaviour [66].

There are numerous approaches to modelling regulatory networks, each of which has its own strengths and weaknesses, capturing the dynamic behaviour of the regulatory network

at different levels of detail [66]. Although closer to reality, detailed models involve a large number of parameters and therefore require a lot of information. In contrast, the simplest coarse-scale models only require knowledge of the architecture of the regulatory network, which can make them easier to construct and interpret [67].

In this thesis we make use of three different types of differential equation to model TRNs consisting of a small number of genes. Following a brief history of modelling cell fate, we describe these models in detail and explain how those that incorporate stochasticity can be simulated using computational methods. Starting with the chemical master equation (CME), we show how a series of approximations lead progressively to the stochastic chemical Langevin equation (CLE), the chemical Fokker-Planck equation (CFPE), and then to a set of ordinary differential equations (ODEs) known as the reaction rate equations (RREs).

### 1.2.1   History of Modelling Cell Fate

An early significant attempt to understand cell fate was presented in the 1950s by the developmental biologist Conrad Waddington when he introduced the notion of the 'epigenetic landscape' as a qualitative picture of development [68]. Waddington imagined the specification of different cell types occurring as a ball rolling down sloping channels in a landscape consisting of hills and valleys, whose geological structure is moulded by the genes that control development (see Fig. 1.3, left). As the ball (representing the cell) rolls down the hill, it reaches a point at which the channel splits in two, forcing the cell to chose between the different valleys. The downhill motion of the ball represents how the process moves inexorably forward in developmental time, while differentiation is controlled by the hills, which act as a barrier separating the landscape into distinct valleys. This intuitive metaphor for development and the discrete nature of cell fates was particularly insightful given that relatively little was known about protein synthesis prior to the discoveries of Watson and Crick, whose work on the structure of DNA was also published throughout the 1950s [69, 70].

A minor adaptation to Waddington's landscape was later proposed by Peter Andrews [71] who suggested that a rougher landscape featuring a series of dips in the valleys more accurately represents the progression through a series of relatively stable transitory cell types that occurs during differentiation (see Fig. 1.3, right). In this picture, the likelihood of movement from one state to another depends on the height of the terrain that surrounds the cell. Low

Figure 1.3: Left: Waddington's "Epigenetic landscape" was presented as a conceptualisation of development. The ball represents a cell whose development is driven by the slope of the landscape. As it rolls downhill, it must choose between discrete fates that are represented by the different valleys. Right: Andrews' adaptation of Waddington's epigenetic landscape, includes a series of dips in the valleys that represent the succession of metastable expression profiles that occur during development. Source: [71]

barriers between two adjacent dips would result in a high frequency of transitions, whereas high barriers would result in transitions occurring infrequently.

**Cell Types and Attractors of Complex Regulatory Networks**

The notion that cell fate decisions are regulated by complex networks was envisioned as early as the 1940s. The physicist Max Delbruck was an early proponent of the notion that distinct cell types correspond to dynamically stable states of underlying molecular regulatory networks [72, 73]. This notion was supported in 1960 by the molecular biologists Jaques Monod and Francois Jacob [74], who note that microbial regulatory elements "... could be connected into a wide variety of 'circuits' endowed with any desired degree of stability". Similarly (and also in the 1960s) the theoretical biologist Stuart Kauffman envisioned cell fates as arising from the dynamics of *complex* genetic regulatory networks [75, 76]. Since the structure of genetic regulatory networks was unknown in the 1960s, he used computational models to generate random networks in which each gene is a boolean variable (either "on" or "off"), and randomly assigned rules of interaction between them. Simulation of these models showed that large randomly structured regulatory networks obeying certain conditions give rise to some characteristics of a real differentiation process. In particular, he noted that they gave rise to relatively few attractor states, corresponding to the fact that the number of adult human cell types is considerably less than the number of possible genetic configurations. This led Kauffman to propose that mature cell types correspond to attractors of high-dimensional regulatory networks.

Experimental evidence supporting Kauffman's attractor hypothesis has been provided by Sui Huang and coworkers [77]. They showed that two biochemically distinct stimuli (the

solvent DMSO and the hormone ATRA) were both able to trigger neutrophil differentiation in human promyelocytic HL60 cells. Moreover, they demonstrated that these stimuli did not trigger differentiation in the same way: time-series data revealed that initially the two stimuli triggered divergent patterns of gene expression, but the two different trajectories eventually converged to the same differentiated neutrophil state, suggesting the presence of an attractor.

After 60 years of study, the molecular complexity of cell fate regulation is only just becoming clear and the dynamics of the cell fate regulatory networks that underlie development are still the subject of considerable research interest.

### 1.2.2 The Chemical Master Equation

As is the case for all the modelling approaches in this thesis, the master equation framework ignores spatial information and simply keeps track of the number of molecules of each type [78]. This simplification is valid for a well-stirred system, where molecules of each species are spread uniformly throughout the spatial domain in which they are confined [78]. Since we do not know the exact positions and velocities of the molecules, we think in terms of the probability that each reaction takes place [79]. It is also assumed that the system is in thermal equilibrium and that the volume of the spatial domain is fixed [80].

In a chemically reacting system there are $N$ molecular species $S_1, \ldots, S_N$ that take part in one or more of $M$ reactions $R_1, \ldots, R_M$ [80]. Reaction $R_j$ describes a single instantaneous physical event that changes the copy number of at least one species [81]. For example, a molecule of species A can bind to a molecule of species B to create a molecule of species C, thus decreasing the copy number of both species A and B by one, and increasing that of C by one.

The number of molecules of each species at time $t$ is quantified by the state vector

$$\boldsymbol{x}(t) = (x_1(t), \ldots x_N(t))^T$$

where $x_i(t)$ is a non-negative integer that represents the number of molecules of species $i$ present in the system at time $t$ [79]. The state vector, $\boldsymbol{x}(t)$, changes when one of the $M$ reactions takes place.

In the chemical master equation (CME) formulation, each reaction is characterised by two quantities [81, 82]:

1. The first is the resulting change in the copy numbers of the molecular species. The $j$th reaction has an associated *stoichiometric*, or *state-change* vector, $\boldsymbol{\nu}_j \in \mathbb{R}^N$, whose $i$th component, $\boldsymbol{\nu}_{ij}$ is the change in the number of molecules of species $i$ as a result of the $j$th reaction. Thus when reaction $j$ occurs, the state vector changes from $\boldsymbol{x}(t)$ to $\boldsymbol{x}(t) + \boldsymbol{\nu}_j$. These column vectors form the $N \times M$ stoichiometric matrix $[\boldsymbol{\nu}_1, \ldots, \boldsymbol{\nu}_M]$.

2. The second quantity is the *propensity function* $a_j(\boldsymbol{x})$. The quantity $a_j(\boldsymbol{x})dt$ is the probability that reaction $R_j$ occurs within the next infinitesimal time interval $[t, t+dt)$, given that the system is in state $\boldsymbol{x}$ at time $t$.

A chemical reaction is written in the form [83]

$$\text{reactants} \xrightarrow{reaction\ rate} \text{products}.$$

For example, consider a reaction where a molecule of A binds to a molecule of B, at rate $k$, to create a molecule of C. This is written as $A + B \xrightarrow{k} C$.

Every reaction may be decomposed into a series of *elemental* reactions that involve either one molecule (*unimolecular*) or two molecules (*bimolecular*) whose interaction yields one or more molecular products [81]. A unimolecular reaction is written as

$$S_i \xrightarrow{k} \text{products},$$

and a bimolecular reaction as

$$S_i + S_{i'} \xrightarrow{k} \text{products}.$$

The 'products' could be one or two molecules of any species, including those on the left hand side of the arrow. Reactions involving three or more molecules can be broken down into a series of elemental reactions [84]. For example, $S_1 + S_2 + S_3 \rightarrow S_4 + S_5$ could be written as $S_1 + S_2 \rightarrow S_{12}$ followed by $S_{12} + S_3 \rightarrow S_4 + S_5$. In the CME formalism it is assumed that reaction $R_j$ is a unimolecular or a bimolecular reaction.

Assuming the kinetics follow the law of mass action, the form of the function $a_j(\boldsymbol{x})$ for a bimolecular reaction is

$$a_j(\boldsymbol{x}) = k_j h_j(\boldsymbol{x}) \tag{1.1}$$

where $k_j$ is the *probability rate constant* or *reaction rate constant* for reaction $j$, defined such that the probability that a randomly chosen pair of molecules $S_i$ and $S_{i'}$ will collide *and* actually react according to $R_j$ in the next time period $dt$ [79]. The function $h_j(\boldsymbol{x})$ is the number of unique combinations of molecules $S_i$ and $S_{i'}$ present (the number of opportunities for the required molecules to collide). So for example, if $R_1$ is the reaction $S_1 + S_2 \xrightarrow{k_1} S_3$, then $a_1(\boldsymbol{x}) = k_1 x_1 x_2$, and if $R_2$ is the reaction $S_1 + S_1 \xrightarrow{k_2} S_4$, then $a_2(\boldsymbol{x}) = \dfrac{k_2}{2} x_1(x_1 - 1)$. If $R_j$ is a unimolecular, reaction Eq. (1.1) still applies, and the function $h_j(\boldsymbol{x})$ is just the copy number of the reactant species. For example, if $R_3$ is the reaction $S_1 \xrightarrow{k_3} S_2$ then $a_3(\boldsymbol{x}) = k_3 x_1$.

To model a constant rate of production of a species from an entity whose abundance we are not interested in capturing, we write the reaction in the form [85]

$$\phi \xrightarrow{k_j} \text{products,}$$

where $a_j(\boldsymbol{x}) = k_j$.

Given the probabilistic nature of the dynamics, we would like to study the evolution of the probability $p(\boldsymbol{x}, t)$ that the system is in state $\boldsymbol{x}$ at time $t$. A time evolution equation for $p(\boldsymbol{x}, t)$ can be deduced using the laws of probability. We do this by writing an equation for the probability of being in state $\boldsymbol{x}$ at time $t + \Delta t$, given that we know the probability of being in any possible state at time $t$ and assuming $\Delta t$ is small enough that at most one reaction can take place in the interval $[t, t + \Delta t)$ [79]. To be in state $\boldsymbol{x}$ at time $t + \Delta t$, the system must have either already been in state $\boldsymbol{x}$ at time $t$ and no reaction took place over $[t, t + \Delta t)$, or the system was in state $\boldsymbol{x} - \boldsymbol{\nu}_j$ for some $1 \leq j \leq M$ at time $t$ and the $j$th reaction took place in the interval $[t, t + \Delta t)$, causing the system to transition to state $\boldsymbol{x}$. So from the law of total probability we obtain [80]

$$\frac{p(\boldsymbol{x}, t + \Delta t) - p(\boldsymbol{x}, t)}{\Delta t} = \sum_{j=1}^{M} \left( a_j(\boldsymbol{x} - \boldsymbol{\nu}_j) p(\boldsymbol{x} - \boldsymbol{\nu}_j, t) - a_j(\boldsymbol{x}) p(\boldsymbol{x}, t) \right) + \mathcal{O}(\Delta t).$$

If we now let $\Delta t \to 0$ we get the CME [79]

$$\frac{dp(\boldsymbol{x},t)}{dt} = \sum_{j=1}^{M} a_j(\boldsymbol{x} - \boldsymbol{\nu}_j)p(\boldsymbol{x} - \boldsymbol{\nu}_j, t) - a_j(\boldsymbol{x})p(\boldsymbol{x}, t). \tag{1.2}$$

Since the state vector $\boldsymbol{x}$ can take a large number of possible values, the CME (1.2) is a very large system of coupled, linear ODEs. The $k$th equation gives the probability of the system being in the $k$th state at time $t$. Generally, the CME cannot be solved analytically, so we use the stochastic simulation algorithm (SSA, also known as Gillespie's algorithm) to obtain an evolving probability distribution by computing sample trajectories of the state vector $\boldsymbol{x}(t)$.

**The SSA (Gillespie's Algorithm)**

The SSA [86, 87] enables us to build up an evolving probability distribution by computing sample trajectories of the state vector $\boldsymbol{x}(t)$. This involves successively advancing the system from its current state by exactly one reaction event, where the probability of a particular reaction occurring reflects the corresponding probability given by the CME.

At the heart of the SSA is the generation of two random numbers: the time $\tau$ to the next reaction, and the index $j$ of that reaction. Thus the key quantity is $p(\tau, j | \boldsymbol{x}, t)d\tau$, the probability that, given $\boldsymbol{x}$, the next reaction will occur in the time interval $[t + \tau, t + \tau + d\tau)$, and it will be the $j$th reaction. Since $\boldsymbol{x}(t)$ is a Markov process, we can begin to derive an expression for this probability by writing

$$\begin{aligned} p(\tau, j | \boldsymbol{x}, t).d\tau &= p\left(\text{no reaction takes place over } [t, t + \tau)\right) \\ &\quad \times p\left(j\text{th reaction takes place over } [t + \tau, t + \tau + d\tau)\right), \end{aligned} \tag{1.3}$$

where $d\tau$ is so small that at most one reaction can take place over that length of time. One way to derive the first probability on the RHS of this equation is to divide $[t, t + \tau)$ into $n$ intervals, as illustrated below, and determine the limit as $n \to \infty$. Now, the probability of

no reactions occurring in a time interval of length $\tau/n$ is

$$1 - a_0(\boldsymbol{x})\frac{\tau}{n},$$

where

$$a_0(\boldsymbol{x}) = \sum_{j=1}^{M} a_j(\boldsymbol{x}),$$

so the probability of no reactions occurring in all $n$ consecutive intervals is

$$\lim_{n\to\infty} \left(1 - a_0(\boldsymbol{x})\frac{\tau}{n}\right)^n = e^{-a_0(\boldsymbol{x})\tau}.$$

The second term on the RHS of Eq. (1.3) – the probability that the $j$th reaction takes place over the interval $[t + \tau, t + \tau + d\tau)$ – is simply $a_j(\boldsymbol{x})d\tau$, thus

$$p(\tau, j|\boldsymbol{x}, t)d\tau = e^{-a_0(\boldsymbol{x})\tau}a_j(\boldsymbol{x})d\tau,$$

$$\Rightarrow \quad p(\tau, j|\boldsymbol{x}, t) = e^{-a_0(\boldsymbol{x})\tau}a_j(\boldsymbol{x}). \tag{1.4}$$

This can be re-written as

$$p(\tau, j|\boldsymbol{x}, t) = \frac{a_j(\boldsymbol{x})}{a_0(\boldsymbol{x})}a_0(\boldsymbol{x})e^{-a_0(\boldsymbol{x})\tau}. \tag{1.5}$$

So the joint probability density function $p(\tau, j|\boldsymbol{x}, t)$ can be written as the product of two separate density functions:

1. $j$ is the integer random variable with probability mass $\frac{a_j(\boldsymbol{x})}{a_0(\boldsymbol{x})}$.

2. $\tau$ is the continuous random variable with an exponential probability density with mean $\frac{1}{a_0(\boldsymbol{x})}$.

The fact that $\tau$ and $j$ are statistically independent is very useful from a computational perspective, since it allows us to simulate a reaction index and a reaction time independently, using a sample from the uniform distribution on the interval $(0, 1)$ and the cumulative distribution functions.

The "direct" method" (we will see some approximations later) of implementing the SSA is therefore as follows:

1. Given that the system is in state $\boldsymbol{x}$ at time $t$, evaluate $a_1(\boldsymbol{x}), \ldots, a_M(\boldsymbol{x})$ and $a_0(\boldsymbol{x}) = \sum_{j=1}^{M} a_j(\boldsymbol{x})$.

2. Draw two random numbers from the uniform distribution on the interval $(0, 1)$, $r_1$ and $r_2$, and calculate $\tau$ and $j$ according to:

$$\tau = \frac{1}{a_0(\boldsymbol{x})} \ln \left( \frac{1}{r_1} \right)$$

$$j = \text{the smallest integer satisfying } \sum_{k=1}^{j} a_k(\boldsymbol{x}) > r_2 a_0(\boldsymbol{x})$$

3. Update $t \to t + \tau$ and $\boldsymbol{x} \to \boldsymbol{x} + \boldsymbol{\nu}_j$

4. Return to step 1 or end the simulation.

The SSA is an exact stochastic method to simulate a chemical reaction system, in the sense that the statistical properties that underlie the CME are reproduced precisely [78]. It is easy to implement, but can be very slow when there are large numbers of molecules in the system because reactions occur frequently, resulting in many iterations within a given period of time. An alternative strategy called the tau-leaping approximation speeds up simulations of the CME by making approximations to the probabilities assigned to the successive states [81], and as we will see, features in the transition from the CME to the chemical Langevin equation.

### 1.2.3    The Chemical Fokker-Planck Equation

The Chemical Fokker-Planck Equation (CFPE) is a multivariate partial differential equation that describes the time evolution of the joint probability distribution of the state of a chemically reacting system. We now show how the multivariate CFPE can be derived from the CME.

First, we relax the condition that the components of $\boldsymbol{x}$ are integers, and allow them to take real values (which is reasonable for a continuous, large number approximation). Assuming that the function $f_j(\boldsymbol{x}) = a_j(\boldsymbol{x})p(\boldsymbol{x}, t)$ is analytic (infinitely differentiable) in the real variable $\boldsymbol{x}$, we can use Taylor's theorem to expand the first term on the right hand side of Eq. (1.2). The Taylor series expansion of a general function $f(\boldsymbol{x})$ is

$$f(\boldsymbol{x} - \boldsymbol{\nu}) = f(\boldsymbol{x}) - \boldsymbol{\nu}^T D f(\boldsymbol{x}) + \frac{1}{2} \boldsymbol{\nu}^T D^2 f(\boldsymbol{x}) \, \boldsymbol{\nu} + \ldots,$$

where $Df(\boldsymbol{x})$ is the Jacobian of $f(\boldsymbol{x})$ and $D^2 f(\boldsymbol{x})$ is the Hessian matrix. Substituting $f(\boldsymbol{x})$ with $a_j(\boldsymbol{x})p(\boldsymbol{x},t)$ and cancelling the first term of the Taylor Series expansion with the last term on the RHS of Eq. (1.2), the CME becomes:

$$
\begin{aligned}
\frac{\partial p(\boldsymbol{x},t)}{\partial t} &= -\sum_{j=1}^{M} \boldsymbol{\nu}_j^T D[a_j(\boldsymbol{x})p(\boldsymbol{x},t)] + \frac{1}{2}\sum_{j=1}^{M} \boldsymbol{\nu}_j^T D^2[a_j(\boldsymbol{x})p(\boldsymbol{x},t)]\boldsymbol{\nu}_j + \dots, \\
&= -\sum_{i=1}^{N} \frac{\partial}{\partial x_i}\left[\sum_{j=1}^{M}\nu_{ij}a_j(\boldsymbol{x})p(\boldsymbol{x},t)\right] + \frac{1}{2}\sum_{i=1}^{N}\sum_{k=1}^{N}\frac{\partial^2}{\partial x_i x_k}\left[\sum_{j=1}^{M}\nu_{ij}\nu_{kj}a_j(\boldsymbol{x})p(\boldsymbol{x},t)\right] + \dots \\
&= -\sum_{i=1}^{N}\frac{\partial}{\partial x_i}[F_i(\boldsymbol{x})p(\boldsymbol{x},t)] + \frac{1}{2}\sum_{i=1}^{N}\sum_{k=1}^{N}\frac{\partial^2}{\partial x_i x_k}[D_{ik}(\boldsymbol{x})p(\boldsymbol{x},t)] + \dots,
\end{aligned}
\tag{1.6}
$$

where $F_i(\boldsymbol{x}) = \sum_{j=1}^{M}\nu_{ij}a_j(\boldsymbol{x})$ and $D_{ik}(\boldsymbol{x}) = \sum_{j=1}^{M}\nu_{ij}\nu_{kj}a_j(\boldsymbol{x})$. Eq. (1.6) is the Kramers-Moyal expansion of the CME. Truncation after the second order term leads to the multivariate CFPE. The CFPE is a partial differential equation that describes the time evolution of the probability distribution function, $p(\boldsymbol{x},t)$. It can be written in the form

$$
\frac{\partial p(\boldsymbol{x},t)}{\partial t} = -\nabla \cdot \left[\mathbf{F}(\boldsymbol{x})p(\boldsymbol{x},t) - \frac{1}{2}\nabla \cdot [\mathbf{D}(\boldsymbol{x})p(\boldsymbol{x},t)]\right],
\tag{1.7}
$$

where $\boldsymbol{F}(\boldsymbol{x}) = \boldsymbol{S}\,\boldsymbol{a}(\boldsymbol{x})$ is the $N \times 1$ drift vector, $\boldsymbol{S}$ is the $N \times M$ stoichiometric matrix, $\boldsymbol{a}(\boldsymbol{x})$ is the $M \times 1$ column vector of propensity functions, and $\boldsymbol{D}(\boldsymbol{x}) = \sum_{j=1}^{M} \boldsymbol{\nu}_j\,\boldsymbol{\nu}_j^T a_j(\boldsymbol{x})$ is the $N \times N$ diffusion matrix.

In one dimension the CFPE reads:

$$
\frac{\partial p(x,t)}{\partial t} = -\frac{\partial}{\partial x}\left[\sum_{j=1}^{M}\nu_j a_j p(x,t)\right] + \frac{1}{2}\frac{\partial^2}{\partial x^2}\left[\sum_{j=1}^{M}\nu_j^2 a_j p(x,t)\right],
\tag{1.8}
$$

which is commonly written as

$$
\frac{\partial p(x,t)}{\partial t} = -\frac{\partial}{\partial x}[F(x)p(x,t)] + \frac{1}{2}\frac{\partial^2}{\partial x^2}\left[\sigma^2(x)p(x,t)\right],
\tag{1.9}
$$

where $F(x) = \sum_{j=1}^{M}\nu_j a_j$ and $\sigma^2(x) = \sum_{j=1}^{M}\nu_j^2 a_j$, or

$$
\frac{\partial p(x,t)}{\partial t} = -\frac{\partial}{\partial x}J(x,t),
\tag{1.10}
$$

where

$$
J(x,t) = F(x)p(x,t) - \frac{1}{2}\frac{\partial}{\partial x}\left[\sigma^2(x)p(x,t)\right].
\tag{1.11}
$$

At the steady-state, the probability current $J(x,t) = c$, where $c$ is a constant. Since there is no probability flow through $x = 0$ (i.e. negative numbers of molecules are not possible), $c$ must equal zero, and so the stationary probability density function, $p_\infty(x)$ satisfies

$$\frac{dp_\infty(x)}{dx} + \left( \frac{2}{\sigma^2(x)} \frac{d\sigma^2(x)}{dx} - \frac{F(x)}{\sigma^2(x)} \right) p_\infty(x) = 0.$$

Using the integration factor

$$I(x) = \exp \left( 2 \ln \left( \sigma^2(x) \right) - \int_0^x \frac{F(y)}{\sigma^2(y)} dx \right),$$

we obtain

$$p_\infty(x) = \frac{Z^{-1}}{\sigma^2(x)} \exp \left[ 2 \int_0^x \frac{F(y)}{\sigma^2(y)} dy \right], \tag{1.12}$$

where

$$Z = \int_0^\infty \frac{1}{\sigma^2(x)} \exp \left[ 2 \int_0^x \frac{F(y)}{\sigma^2(y)} dy \right], \ dx \tag{1.13}$$

is a normalising constant which ensures that $p_\infty(x)$ is a proper probability distribution.

### 1.2.4   The Chemical Langevin Equation

The Chemical Langevin Equation (CLE) is a stochastic differential equation that describes the time evolution of the state of a chemically reacting system. It is mathematically equivalent to the CFPE, and the probability density function of $\boldsymbol{x}(t)$ obeys Eq. (1.7) [79]. The CLE can be written in the form

$$\frac{dx}{dt} = F(x) + \sigma(x)\xi(t), \tag{1.14}$$

where $F(x)$ and $\sigma(x)$ are known functions, and $\xi(t)$ is a rapidly fluctuating stochastic term, known as white noise. A stochastic process $\xi(t)$ is called white noise if its time average value is zero i.e., $\langle \xi(t) \rangle = 0$, and for $t \neq t'$, $\xi(t)$ and $\xi(t')$ are statistically independent (i.e. the fluctuation function has no correlation at different times), i.e., $\langle \xi(t)\xi(t') \rangle = \delta(t - t')$.

The function $\sigma(x)$ in Eq. (1.14) denotes the form of the amplitude of the noise. If $\sigma(x)$ is a constant then the system is subject to additive noise, otherwise it is subject to multiplicative noise.

For a system with more than one species, the multi-dimensional CLE reads:

$$\frac{d\boldsymbol{x}}{dt} = \boldsymbol{F}(\boldsymbol{x}) + \boldsymbol{\sigma}(\boldsymbol{x})\xi(t). \tag{1.15}$$

We can derive the form of the the functions $\boldsymbol{F}(\boldsymbol{x})$ and $\boldsymbol{\sigma}(\boldsymbol{x})$ directly from the CME to obtain the CLE. A common derivation [79] begins with a formula that was originally proposed to speed up the CME simulation algorithm [88]. It involves approximately advancing the system from state $\boldsymbol{x}$ at time $t$ by a preselected time $\tau$, during which more than one reaction may occur. Placing two conditions on $\tau$ will allow us to transition from the CME to the CLE.

If we choose $\tau$ small enough such that all the propensity functions remain approximately constant during that time period, i.e., if

$$a_j(\boldsymbol{x}) = \text{const in } [t, t + \tau), \ \forall j, \tag{1.16}$$

then we can define the Poisson random variable with mean $a\tau$, $\mathcal{P}(a\tau)$, to be the number of reactions that will occur in a time $\tau$, given that the probability of a reaction occurring in the next infinitesimal time $dt$ is $adt$, where $a$ is a positive constant. Together with Eq. (1.1), this implies that the number of times reaction $j$ occurs in the next $\tau$ is $\mathcal{P}(a_j(\boldsymbol{x})\tau)$. Since reaction $j$ changes the state of the system by $\boldsymbol{\nu}_j$, the state of the system at time $t + \tau$ is

$$\boldsymbol{x}(t + \tau) = \boldsymbol{x} + \sum_{j=1}^{M} \mathcal{P}(a_j(\boldsymbol{x})\tau)\boldsymbol{\nu}_j. \tag{1.17}$$

Eq. (1.17) is called the *tau-leaping formula*, and it forms the basis of the *tau-leaping algorithm*. The errors introduced by this approximation will be small as long as the state vector updates are relatively small, and can be reduced by adaptively choosing the leap time $\tau$ based on the current state vector and propensity function values.

The second condition allows us to approximate the Poisson random variable by a normal random variable. Given that the expectation of $\mathcal{P}(a_j(\boldsymbol{x})\tau)$ is $a_j(\boldsymbol{x})\tau$, we require that each reaction will occur 'many times' during $\tau$. That is

$$a_j(\boldsymbol{x})\tau \gg 1, \forall j. \tag{1.18}$$

If this is the case we can use the result that

$$\mathcal{P}(\mu) \approx \mathcal{N}(\mu, \mu), \text{ for } \mu \gg 1,$$

which, together with the identity

$$\mathcal{N}(\mu, \sigma^2) = \mu + \sigma \mathcal{N}(0, 1),$$

allows us to further approximate Eq. (1.17) as follows:

$$
\begin{aligned}
\boldsymbol{x}(t + \tau) &\approx \boldsymbol{x} + \sum_{j=1}^{M} \mathcal{N}_j(a_j(\boldsymbol{x})\tau, a_j(\boldsymbol{x})\tau)\boldsymbol{\nu}_j, \\
&\approx \boldsymbol{x} + \sum_{j=1}^{M} \left[ a_j(\boldsymbol{x})\tau + \sqrt{a_j(\boldsymbol{x})\tau}\mathcal{N}_j(0, 1) \right] \boldsymbol{\nu}_j, \\
&\approx \boldsymbol{x} + \sum_{j=1}^{M} \boldsymbol{\nu}_j a_j(\boldsymbol{x})\tau + \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j(\boldsymbol{x})}\sqrt{\tau}\mathcal{N}(0, 1),
\end{aligned}
\tag{1.19}
$$

which is the Euler–Maruyama discretisation of the continuous time stochastic differential equation (SDE)

$$
\frac{d\boldsymbol{x}}{dt} = \sum_{j=1}^{M} \boldsymbol{\nu}_j a_j(\boldsymbol{x}) + \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j(\boldsymbol{x})}\xi_j(t),
\tag{1.20}
$$

where $\xi(t) \sim \mathcal{N}\left(0, \dfrac{1}{dt}\right) = \dfrac{1}{\sqrt{dt}}\mathcal{N}(0, 1)$, and the $\xi_j(t)$ are statistically independent Gaussian white-noise processes. In the limit $\tau \to 0$, the discrete time recurrence (1.19) converges to the continuous time process described by (1.20). Thus Eq. (1.19) allows us to compute approximate numerical solutions for Eq. (1.20).

Comparing this result with Eq. (1.15), we see that:

$$
\boldsymbol{F}(\boldsymbol{x}) = \sum_{j=1}^{M} \boldsymbol{\nu}_j a_j(\boldsymbol{x}); \quad \boldsymbol{\sigma}(\boldsymbol{x}) = \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j(\boldsymbol{x})}.
\tag{1.21}
$$

The discrete stochastic process $\boldsymbol{x}(t)$ has now been approximated by a continuous stochastic process, but this is only valid if the system satisfies both conditions given by Eqs. (1.16) and (1.18). It should also be noted that the CLE will not be accurate for describing rare events, because the approximation $\mathcal{P}(\mu) \approx \mathcal{N}(\mu, \mu)$, for $\mu \gg 1$ is a very poor approximation for values that are very far from their mean, $\mu$ (more than a few standard deviations $\sqrt{\mu}$ away) [89]. This means that the CLE nearly always underestimates the likelihood of rarely

occurring reactions, which in some circumstances may have a dramatic impact on the dynamics. Therefore the CLE is only valid over a limited period of time for systems that have rare events. Observed over a long period of time, a simulation of the original CME would generate the atypical behaviour resulting from rare events, but a simulation of the CLE will not generate the rare events, and therefore will only accurately describe the typical behaviour of the system.

### 1.2.5 Reaction Rate Equations

Reaction Rate Equations (RREs) consist of a set of $N$ coupled ordinary differential equations (ODEs), with one equation for each molecular component in a chemically reacting system. To obtain the deterministic RREs, we could simply ignore the stochastic part of the CLE, although we currently have no justification for doing so. The CLE is more commonly written in terms of the species concentrations

$$y_i(t) = \frac{x_i(t)}{\Omega}. \tag{1.22}$$

We can use what is known as *the thermodynamic limit* to derive the RREs from the CLE. The thermodynamic limit is a large-system limit in which the species populations, $x_i(t)$, and the containing volume, $\Omega$, all tend to infinity in such a way that the species concentrations, $y_i(t)$, remain constant (with respect to the limit). In this limit, the deterministic coefficients $(a_j)$, and therefore the deterministic drift terms, grow as the system size, but the stochastic diffusion terms grow as the square root of the system size. In this limit the CLE reduces to

$$\frac{d\boldsymbol{x}}{dt} = \sum_{j=1}^{M} \boldsymbol{\nu}_j a_j(\boldsymbol{x}), \tag{1.23}$$

and the continuous stochastic process $\boldsymbol{x}(t)$ has now been approximated by a continuous deterministic process.

In fact, being close to the thermodynamic limit is also a necessary condition for the validity of the CLE, which emphasises the fact that both the CLE and the RRE are only valid when there are large numbers of molecules of each species. Moreover, it has been proven [84] that both conditions (1.16) and (1.18) on $\tau$ are satisfied in the thermodynamic limit.

We note that the RRE provides us with an expression for the state of a chemically reacting system in thermodynamic limit, and not necessarily an expression for the dynamics of the mean abundance of each species. We can obtain the time evolution of the expected value of $\boldsymbol{x}$ from the CME (1.2) by multiplying both sides of Eq. (1.2) by $\boldsymbol{x}$, and then summing over all $\boldsymbol{x}$. Doing so, we get

$$\frac{\partial \langle \boldsymbol{x} \rangle}{\partial t} = \sum_{j=1}^{M} \boldsymbol{\nu}_j \langle a_j(\boldsymbol{x}) \rangle, \tag{1.24}$$

where $\langle \boldsymbol{x} \rangle$ is the expected state of the system, and we obtain the same expression by taking the expectation of both sides of the CLE (1.20). Now, if all chemical reactions in the system are unimolecular, the propensity functions are linear, and $\langle a_j(\boldsymbol{x}) \rangle = a_j(\langle \boldsymbol{x} \rangle)$. In this case, Eq. (1.24) is a set of closed ODEs for the means, $\langle \boldsymbol{x} \rangle$, and is identical to the RRE (1.23). However, if there are any bimolecular reactions, Eq. (1.24) contains at least one second order moment, and the evolution of the moments depend on higher order moments, which themselves depend on even higher order moments, and so on. Therefore, we would obtain an infinite set of open ended equations for all the moments, which is not the same as the RRE (1.23).

### 1.2.6 Summary of Mathematical Models of Regulatory Networks



Figure 1.4: Summary of the relationships between the four models of regulatory networks described in Section 1.2.2. The equations for the evolution of the probability density (and mass) functions are given in the rounded boxes in the left column, and those for the evolution of the state in the right column. The relationships between each equation are summarised on the corresponding arrows. The box in the bottom left hand corner (double line border) contains the definitions of three abbreviations used throughout the figure.

## 1.3    Noise in Gene Expression

Gene expression is a noisy process [90–92] since transcription and translation are inherently stochastic molecular processes that are also affected by environmental noise caused by fluctuations in the regulating TFs and other cellular components [93]. Consequently, the copy number of gene products in an individual cell fluctuates continuously [94–96]. As this variation can affect cell behaviour, cells have developed a range of mechanisms (e.g. negative feedback loops, see [97] and Section 2.2.1) to control this molecular noise [98]. However, molecular noise can have a positive role in cell fate determination [99, 100] because it enables an organism to adapt to changes in the environment without the need for genetic mutations [100]. For example, stochastic effects have been used to explain cell-to-cell variability in clonal populations [93, 96, 101, 102], and there is evidence to suggest that such variability in mammalian progenitor cell populations primes the cells for different lineages choices [103].

The models of variability in gene expression presented in this thesis are based on positive feedback and the intrinsic noise in transcription and translation. Although deterministic analyses can be useful for understanding the main properties of a dynamical system, they can fail to capture some behaviours. Deterministic models are a sufficient approximation for well-stirred, chemically reacting systems comprising a large number of molecules such that discreteness and stochasticity are not significant. However, if the copy number is small (and many gene products are present in small copy numbers [95, 104]) or the system is susceptible to noise amplification, for example, then stochasticity can have a major effect on the behaviour of the system.

We now demonstrate the importance of stochastic methods for describing the dynamics of regulatory networks by presenting two examples where deterministic models fail to capture the qualitative behaviour of the system, and illustrate why it is 'compulsory' to account for noise. Each example demonstrates that the presence of noise can lead to significant qualitative differences in the behaviour of a system, emphasising that the validity of macroscopic approaches to describe averages cannot be taken for granted.

**Intrinsic Noise in a Monostable System**

Since intrinsic noise in a system can enhance a signal, the deterministic description of the dynamics of a TRN does not necessarily correctly represent the evolution of the mean of the

inherently stochastic system [104, 105]. Consider two proteins $X$ and $Y$ that are synthesised at an equal rate $k$, and decay at rates $d_x$ and $d_y$, respectively. They can also irreversibly bind to form a heterodimer $Z$ with a constant association rate $a$. This chemical reaction system is described by the scheme

$$\phi \xrightarrow{k} X \xrightarrow{d_x} \phi, \quad \phi \xrightarrow{k} Y \xrightarrow{d_y} \phi, \quad X + Y \xrightarrow{a} Z, \quad Z \xrightarrow{d_z} \phi \tag{1.25}$$

In the deterministic setting, these reactions can be described by the RREs

$$\frac{dx}{dt} = k - d_x x - axy \tag{1.26}$$

$$\frac{dy}{dt} = k - d_y y - axy \tag{1.27}$$

$$\frac{dz}{dt} = axy - d_z z \tag{1.28}$$

where $x$, $y$ and $z$ are the concentrations of $X$, $Y$ and $Z$, respectively, and $x(0)$, $y(0)$, and $z(0)$ are the initial conditions. Eqs. (1.26) and (1.27) support one stable equilibrium point, whose basin of attraction is the entire positive quadrant. To illustrate the importance of stochasticity, sample trajectories of the system (Fig. 1.5, left) were simulated using Gillespie's stochastic simulation algorithm (SSA) (see Section 1.2.2) for two sets of parameter values:

$$k = 10^{-1}, \quad d_x = 10^{-6}, \quad d_y = 10^{-5}, \quad a = 10^{-5} \tag{1.29}$$

$$k = 1000, \quad d_x = 10^{-4}, \quad d_y = 10^{-3}, \quad a = 10^{-1} \tag{1.30}$$

The steady state values for $x$ and $y$ are given by $x_* \approx \sqrt{\frac{k}{a}} = 100$ and $y_* \approx \sqrt{\frac{k}{a}} = 1000$ for both sets of parameters, but there is a noticeable difference between the sample stochastic trajectories. For the first set of parameters the time series for the copy number of $Y$ (Fig. 1.5, left) remains close to the RRE steady state. However, the second set of parameters leads to an increase in the inherent noise and the copy number reaches up to six-fold the deterministic mean. In this case the noise has a severe effect on the system and demonstrates that a deterministic model masks important microscopic behaviour.

**Noise Induced Oscillations**

Many organisms have developed molecular mechanisms that generate oscillating levels of expression of proteins with a period roughly equal to 24 hours [105]. These oscillations, known as the circadian rhythm, allow the organism to synchronise their biological processes

Figure 1.5: Left: Stochastic simulation of the chemical system given by Eq. (1.25) using Gillespie's stochastic simulation algorithm (SSA). The first set of parameters (1.29) was used to generate the plot in black, and the second set (1.30) yields the plot in blue. Right: Time evolution of the repressor protein for the deterministic equation, whose solution was found using a nonstiff differential equation solver (ode45 in Matlab) (top) and stochastic version (bottom), simulated using the SSA using the same parameter values (1.32).

with environmental periodicity. An experimentally determined, minimal set of mechanisms required for the circadian system is described by [106] with a model of two genes: an activator $A$ and a repressor $R$. The activator protein $A$ binds to (unbinds from) its own promoter region and that of $R$ at rates $\gamma_A$ ($\phi_A$) and $\gamma_R$ ($\phi_R$) respectively, increasing the transcription rates of $A$ and $R$ mRNA from $\alpha_A$ and $\alpha_R$ to $\alpha'_A$ and $\alpha'_R$, respectively. The $A$ and $R$ mRNA are translated into their respective proteins at rates $\beta_A$ and $\beta_R$, and decay at rates $\delta_{MA}$ and $\delta_{MR}$, respectively. The proteins $A$ and $R$ decay at rates $\delta_A$ and $\delta_R$. Thus $A$ creates a positive feedback loop with itself. Meanwhile, the repressor protein can bind to the activator protein, creating dimer $C$, and preventing protein $A$ from binding to its promoter region, thus creating a negative feedback loop. The dimer $C$ decays at rate $\delta_A$, the product of which is protein $R$. This chemical reaction system is described by the scheme

$$D_A + A \xrightarrow{\gamma_A} D'_A \xrightarrow{\phi_A} D_A + A, \quad D_R + A \xrightarrow{\gamma_R} D'_R \xrightarrow{\phi_R} D_R + A,$$

$$D_A \xrightarrow{\alpha_A} D_A + M_A, \quad D'_A \xrightarrow{\alpha'_A} D'_A + M_A, \quad M_A \xrightarrow{\delta_{MA}} \phi,$$

$$D_R \xrightarrow{\alpha_R} D_R + M_R, \quad D'_R \xrightarrow{\alpha'_R} D'_R + M_R, \quad M_R \xrightarrow{\delta_{MR}} \phi,$$

$$M_A \xrightarrow{\beta_A} M_A + A, \quad A \xrightarrow{\delta_A} \phi, \quad M_R \xrightarrow{\beta_R} M_R + R, \quad R \xrightarrow{\delta_R} \phi,$$

$$A + R \xrightarrow{\gamma_c} C \xrightarrow{\delta_A} R. \tag{1.31}$$

In the deterministic setting, these reactions can be described by the RREs

$$\frac{dD_A}{dt} = \phi_A D'_A - \gamma_A D_A A$$

$$\frac{dD'_A}{dt} = \gamma_A D_A A - \phi_A D'_A$$

$$\frac{dD_R}{dt} = \phi_R D'_R - \gamma_R D_R A$$

$$\frac{dD'_R}{dt} = \gamma_R D_R A - \phi_R D'_R$$

$$\frac{dM_A}{dt} = \alpha'_A D'_A + \alpha_A D_A - \delta_{MA} M_A$$

$$\frac{dM_R}{dt} = \alpha'_R D'_R + \alpha_R D_R - \delta_{MR} M_R$$

$$\frac{dA}{dt} = \beta_A M_A + \phi_A D'_A + \phi_R D'_R - A\left(\gamma_A D_A + \gamma_R D_R + \gamma_C R + \delta_A\right)$$

$$\frac{dR}{dt} = \beta_R M_R - \gamma_C A R + \delta_A C - \delta_R R$$

$$\frac{dC}{dt} = \gamma_C A R - \delta_A C$$

where $D'_A$ and $D_A$ denote the concentrations of activator genes with and without protein $A$ bound to its promoter, respectively; similarly, $D'_R$ and $D_R$ refer to the repressor promoter; $M_A$ and $M_R$ denote the concentrations of $A$ and $R$ mRNA; $A$ and $R$ correspond to the activator and repressor proteins; and $C$ corresponds to the inactivated complex formed by $A$ and $R$; and $D_A(0)$, $D'_A(0)$, $D_R(0)$, $D'_R(0)$, $M_A(0)$, $M_R(0)$, $A(0)$, $R(0)$ and $C(0)$ are the initial conditions.

For the following set of parameter values

$$\alpha'_A = 500, \quad \alpha_A = 50, \quad \alpha'_R = 50, \quad \alpha_R = 0.01, \quad \beta_A = 50, \quad \beta_R = 5, \quad \delta_A = 1, \quad \delta_R = 0.2,$$

$$\delta_{MA} = 10, \quad \delta_{MR} = 0.5, \quad \gamma_C = 2, \quad \gamma_A = 1, \quad \gamma_R = 1, \quad \phi_A = 50, \quad \phi_R = 100 \tag{1.32}$$

the RRE model of the dynamics of the system exhibits oscillations in the repressor protein. When the degradation rate, $\delta_R$, of the repressor protein is reduced to 0.02, the increased presence of protein $R$ increases the number of opportunities for $A$ to be sequestered, thus reducing the number of activators in the system. Consequently, the number of repressor molecules falls to, and remains at a low steady state (Fig. 1.5, top right). The value of the degradation rate at which the system transitions from an oscillatory regime to a monostable regime is known as a Hopf bifurcation. However, when molecular noise is incorporated in the model, the oscillations reappear. This is illustrated by a sample trajectory of the system, simulated using the SSA (Fig. 1.5, bottom right). The frequency of the noise-induced

oscillations can be manipulated by changing the level of intrinsic noise in the system. This can be done, for example, by changing the speed of the molecular reactions.

## 1.4   Summary

We have introduced the molecular biology relevant to the studies that follow in the next two chapters. Central to this thesis are the TFs that interact with each other to regulate gene expression in stem and progenitor cells, and therefore control cell fate. The effects of these interactions can be represented in the form of a TRN, which can be converted into a set of equations that describe how expression levels change over time. Following a review of the historical understanding of cell fate determination, we described in detail three different types of differential equation used to model TRNs.

The CME is a time evolution equation for the probability, $p(\boldsymbol{x}, t)$, that the system is in state, $\boldsymbol{x} \in \mathbb{N}_0$, with regard to a continuous time variable $t$. We explained how it can be can be deduced from the laws of probability, and can be simulated using the SSA to obtain an evolving probability distribution, by computing sample trajectories of the state vector. A disadvantage of the SSA is that many trajectories are needed for the time-dependent solution of the CME, and only long time-integrations yield accurate results for its steady state solution.

Alternatively, the CME can be approximated by the CFPE – a time dependent partial differential equation – in which the state vector $\boldsymbol{x}$ becomes a real-valued vector, $\boldsymbol{x} \in \mathbb{R}^N$. Since the CFPE represents continuous processes, it is significantly more analytically tractable than the CME. As such, we were able to derive a simple expression for the one-dimensional steady-state probability density function, $p_\infty(x)$, from the one dimensional CFPE.

The CLE is a stochastic differential equation that describes the time evolution of the state vector, $\boldsymbol{x}(t)$. It is mathematically equivalent to the CFPE, in that the probability density function of $\boldsymbol{x}(t)$ obeys the CFPE. In the thermodynamic limit, the CLE reduces to the RREs, a set of coupled ODEs, with one equation for each molecular species.

All three of these differential equation models will be useful for helping us understand the mechanisms that control cell-to-cell variability in both adult and embryonic stem cell populations.

Since the processes involved in gene expression are noisy, it is often important to account for this stochasticity in models of regulatory networks. To illustrate this, we compared the CME model and the equivalent RREs of two chemically reacting systems. In the first example, sample trajectories of the system, simulated using the SSA, showed that an increase in the intrinsic noise can lead to a significant increase in the mean level of expression in the stochastic model, even when the deterministic mean remains unchanged. In the second example, the incorporation of molecular noise in the model resulted in oscillations that were not present in the solution of the RREs. In both cases, the deterministic model masked important microscopic behaviour.

# Chapter 2

# Mathematical Modelling of Pluripotency

## 2.1 Introduction

This chapter is comprised of two main parts, each of which aims to elucidate the underlying mechanisms that govern the experimentally observed variability in the expression levels of the protein Nanog – a TF in ESCs that is thought to be a key factor in maintaining pluripotency. In the first part we explore the combined role that feedback loops and the inherent stochasticity in transcription and translation have in controlling the variability of expression levels at the population level, and cell fate decisions at the single cell level. In the second part we consider the origins of the empirical bimodal distribution of Nanog expression. Although many studies report such heterogeneity, new evidence suggests that it might be a result of the measurement method interfering with the Nanog feedback mechanism – a mechanism that we explore in more detail.

### 2.1.1 The ESC TRN

The regulatory networks that underpin the ESC state contain many protein-protein and protein-DNA interactions forming a very complicated TRN with many positive feedback loops (Fig. 2.1a). Central to this network are three TFs: Oct4, Sox2 and Nanog [107]. A plethora of studies have revealed that these master regulators play a central role in the maintenance of stem cell identity [36, 108–114]. Genome-wide experiments have shown that Nanog, Oct4,

Figure 2.1: a) The *extended* embryonic stem cell transcriptional regulatory network, rich in regulatory loops. The arrows indicate interactions but not their type. b) The *core* embryonic stem cell regulatory network. The genes (blue rectangles) are expressed (dashed lines) and their products (proteins) bind to form dimers (pink ovals). The dimers can bind to the promoter region of the genes (solid lines) and positively affect their rate of expression.

and Sox2 bind cooperatively to promoter regions of several hundred genes involved in the regulation of pluripotency and differentiation [58, 115, 116], indicating that Oct4, Sox2 and Nanog directly affect their rate of transcription.

Experimental observations have shown that Oct4 and Sox2 are relatively homogeneously expressed, whereas Nanog is heterogeneous and exhibits a bistable pattern of expression [3, 107, 108, 117–119]. Studies on the function of Nanog reported that the level of Nanog expression influences developmental potential; the absence of Nanog leads to cell differentiation and the loss of pluripotency both in vivo and in vitro [3, 110, 114, 117], whereas high levels of Nanog maintain pluripotency despite the presence of differentiation inducing stimuli [3, 36]. Nanog shows significant temporal expression fluctuations at the single-cell level [117, 119–124], and it is these fluctuations that give rise to the heterogeneity within ESC populations, ensuring robustness of the population and the long-term regenerative potency of a single cell [124–126]. It has been shown that populations of ESCs that express Oct4 contain a subpopulation (10 - 20%) of cells that are Nanog-low, due to stochastic transitions between Nanog-high and Nanog-low states within individual cells [108]. These fluctuations transiently prime individual ESCs for differentiation without committing them to a particular state [117]. Such findings support the hypothesis that Nanog acts as molecular "gate-keeper"

for transient differentiation signals in fluctuating environments, while preparing the cell for differentiation when the appropriate and persistent stimuli do occur [3, 127, 128].

Oct4, Sox2 and Nanog form the *core* ESC TRN with many positive feedback loops (Fig. 2.1b). Each factor positively regulates the expression of itself and the other two either via an Oct4-Sox2 heterodimer, or a Nanog homodimer [111, 112]. In fact, it has been shown that Nanog dimerisation is necessary for the maintenance of ESC pluripotency and self-renewal [129]. Since positive feedback loops are self-sustaining (by perpetuating transient stimuli) the expression of the three master TFs can be maintained by this network structure. This was demonstrated by Ying et al. [130], who reported that extrinsic stimuli (growth factors or cytokines) were not required for the propagation and maintenance of pluripotency in culture of ESCs derived from mice, since the ESCs were able to self-renew when shielded from differentiation-inducing stimuli. As such, they proposed that the pluripotent state is an intrinsically self-maintaining state (termed "ground state") that does not require extrinsic input.

A study by MacArthur et al. [3] motivates the work that makes up the first part of this chapter, in which we aim to elucidate the molecular basis for the fluctuations in expression in the TRN. In order to reproduce the Nanog expression level fluctuations observed in wild-type ESCs and assess the dynamic response of gene expression upon loss of Nanog, they used a mouse ESC-line (NanogR) [114, 131] in which Nanog expression level can be precisely controlled with a chemical called doxycycline (dox). This cell-line is created by infecting a wild-type ESC with a type of retrovirus that delivers two genes to the nucleus. The expression of these genes is simultaneously driven by dox. The first gene is transcribed into a short hairpin RNA (shRNA) that mark the endogenous Nanog mRNA for degradation and ensures that it is never translated into Nanog proteins, and the second gene codes for a Nanog-GFP protein complex that is immune to shRNA knockdown. Thus exogenous Nanog is not subject to the same regulatory conditions as endogenous Nanog, since all Nanog feedback elements are removed, and dox can be used to turn Nanog expression on (in the presence of dox) and off (on the removal of dox) robustly and synchronously in the entire population, on demand. This Nanog switch was used by MacArthur et al. [3] to restore Nanog expression at several time points following Nanog depletion, and the ability of Nanog to restore the pluripotency network at each time point was assessed, based on expression changes in the extended TRN. When Nanog expression was switched off for 24 hours or less (mimicking a cell stochastically

transitioning from the Nanog-high to the Nanog-low state), the cells adopted a reversible primed state in which most elements of the extended TRN did not show significant changes in expression, and some differentiation genes were co-expressed. However, in the continued absence of Nanog (>24 hours), the pluripotency genes of the extended TRN were irrevocably downregulated and associated differentiation genes upregulated, thus committing them to their associated cell fates. In concordance with earlier studies, these results indicate that Nanog is a potent, nonspecific negative regulator of early lineage decisions, and suggest that early fate changes are reversible at the single cell level. Here, we explore the role of Nanog in the global feedback structure in the extended ESC TRN. Feedback loops (which can be positive, negative or mixed) commonly regulate phenotypic variability in diverse organisms and contexts by generating complex dynamics [132], such as multi-stability [133–138], excitability [108] and oscillations [139–141], and by modulating molecular noise [97, 142]. Accordingly, Nanog fluctuations may regulate early cell fate decisions and population variability by controlling feedback mechanisms in the extended ESC TRN. To investigate this possibility we explore how feedback in network structures relates to dynamics, and then analyse the feedback structure of the extended ESC TRN (Fig. 2.1a).

### 2.1.2   Faithfulness of Nanog Reporter Strategies

In the second part of this chapter, we question the faithfulness of a reporting method that has been used to measure protein expression, including that of Nanog, in a large number of experiments published in prominent journals. The basis of these studies are engineered cell-lines which express a green fluorescent protein (GFP) (either with or without Nanog) under the *Nanog* promoter. GFP, a protein originally isolated from the jellyfish *aequorea victoria*, exhibits bright green fluorescence when exposed to light in the blue to ultraviolet range [143]. There are two copies of *Nanog* on the autosomes [144], each of which has a slightly different DNA sequence [145], and is called an allele. *Heterozygous GFP knock-in* cell lines are created by replacing the coding region for *Nanog* with that for GFP on one allele [118]. Since the promoter region is unchanged, GFP is transcribed at the same rate as *Nanog*. GFP expression – as a proxy for Nanog – is then measured by passing the cells through a flow cytometer. This machine directs beams of light at the cells as they pass through a thin tube that is sufficiently narrow to only allow one cell through at a time, and detects the level of intensity of the resulting fluorescent signals returned by the cell [146]. In the

absence of other cellular material between the detector and the GFP molecules, the recorded fluorescence intensity (FI) is approximately proportional to GFP abundance, subject to a number of environmental variables that can affect the level of precision of the measurement [147]. However, since the GFP molecules are produced inside the nucleus of the cell, the fluorescent signal is attenuated, and the true relationship between FI and and copy number is unknown [148]. GFPs are most commonly used but different colour fluorescent proteins can be used simultaneously to report more than one gene at a time [144, 149].

The heterozygous knock-in strategy is one of an increasing number of reporter constructs that are being used to obtain information about gene products in live cells, both *in vivo* and *in vitro*, as well as over time. Such reporter strategies provide a vital tool for exploring spatial and temporal heterogeneity in development, homeostasis and disease. However, for a reporter to be of practical use and enable conclusions concerning gene expression at the single cell level to be drawn, the reporter signal must be representative of the level of expression of the protein of interest in that particular cell. For example, the faithfulness of reporter strategies at the single-cell- and population-level is highly relevant to experiments where cellular populations are sorted according to the strength of the reporter signal. A common sorting procedure involves dividing bimodal populations into low- and high-expressing subpopulations by flow cytometry. In this case, if a single cell reporter measurement does not reflect the actual expression level of the gene of interest, then functional studies could be confounded by mixed starting populations.

The majority of studies characterising heterogeneity in Nanog expression have used live-cell gene reporter strategies [149]. However, recent evidence suggests that these reporters do not give a faithful reflection of endogenous Nanog expression, because the strategies used apparently perturbs the underlying regulatory network. The extent of the perturbation depends on the reporter construct and the particular regulatory network. Here, we consider the faithfulness of the heterozygous GFP knock-in reporter for Nanog, since it is thought to be the most disruptive reporter strategy currently available [149]. It is thought that heterozygous knock-in reporters have the potential to significantly alter the population distribution of endogenous Nanog expression, because the creation of a *Nanog*-null allele disturbs normal *Nanog* transcriptional control [149] by interfering with its auto-regulation. This conclusion was drawn from a study by Faddah et al. [149], where protein expression from both alleles was monitored with a different colour fluorescent protein for each allele, using a reporter method that does

not alter Nanog function. In this strategy, the coding region for the fluorescent protein is inserted adjacent to the Nanog coding region, separated only by a short sequence that codes for a peptide - a very small protein that causes the two proteins cleave after translation. Contrary to previous reports, these results indicated that both *Nanog* alleles are expressed in the vast majority of cells. Since Nanog-low cells were very rare, a unimodal distribution of expression was observed, instead of a bimodal one. These results were obtained using two standard culture conditions, known as 0i and 2i. 0i contains fetal bovine serum and Leukemia inhibitory factor (LIF) - a cytokine that maintains self-renewal and prevents differentiation by activating the transcription factor STAT3 [110, 117, 144], and 2i is 0i conditions with the addition of mitogen-activated protein kinase and glycogen synthase kinase 3 inhibitors, which maintain 'ground state' pluripotency [130]. Therefore, although both conditions maintain the pluripotency state *in vitro*, the additional inhibitors in 2i culture conditions give the medium a stronger hold on pluripotency than 0i.

As part of a larger study that includes our work in this chapter on the faithfulness of Nanog reporter strategies, an experiment was carried out to investigate the relationship between reporter (GFP) and Nanog protein levels in NHET cells (a reporter cell-line that uses the heterozygous GFP knock-in strategy). Since the GFP reporter only monitors expression of the knocked-in allele, expression of the unaltered *Nanog* allele was measured using a method called immunostaining, where fluorescent proteins are attached to the target proteins via an antibody. Immunostaining kills the cells so can not be used to monitor expression in live cells. In this case a red fluorescent protein, called mCherry, was attached to Nanog proteins, to distinguish them from GFP. Flow cytometry was then used to observe mCherry and GFP fluorescence levels in each cell. The results show that although GFP expression is clearly bimodal, Nanog expression was unimodal under 2i conditions, and the bimodality under 0i conditions was of questionable significance (see Fig. 2.2). The difference between the modality and spread of the observed distributions within culture conditions could be due to, for example, differences in the translation and decay rates of Nanog and GFP, and the reporter methods used (live-cell knock-in or immunostaining).

To date, only two theoretical models have been proposed to explain the (presumed true) bimodality of Nanog expression [108, 139]. The first [108] involves only Nanog and Oct4 and models their dynamics as a noise-driven excitable system that gives rise to a small subpopulation with low Nanog expression through occasional random transient excursions

Figure 2.2: The histograms show the experimentally observed marginal distributions of Nanog (red bars) and GFP (green bars) expression in a GFP knock-in cell line, in 0i and 2i culture conditions. The joint distribution of the paired data is indicated by the scatterplot (blue dots). Data generated by Rosanna Smith, Centre for Human Development, Stem Cells and Regeneration, University of Southampton.

from the high to the low expression state. Such systems use a combination of positive and negative feedback loops to generate pulses of gene expression. In this case, the model is based on mutual Oct4 and Nanog activation, Oct4 and Nanog auto-regulation, and Nanog repression by high levels of Oct4, although the latter assumption has not been tested, and is required to achieve an excitable system. The resulting system has a single stable steady state, corresponding to high levels of both Nanog and Oct4. However, small noise-driven increases in Oct4 expression, amplified through auto-regulation, can perturb the system away from this state, and transiently lead to Nanog repression, enabling the cell to enter a low state of Nanog expression. This in turn decreases the promotion of Oct4, thus terminating the Oct4 pulse and allowing Nanog expression levels to return to the original, steady state.

The excitable system model was favoured over bistability by Kalmar et al. [108] because the observed proportion of Nanog-low in the parental, steady-state population is small (5%-20%), which led the authors to suggest that this state might be a relatively short-lived event. If this were the case, we would not expect to see the vast majority of the cells from the Nanog-low isolated subpopulation remaining in the Nanog-low state after 11 days [108, 117]. However, it appears that the cells do remain in the low state for a long time, suggesting the presence of two attractors [139]. For this reason, Glauche et al. [139] question the justification for an excitable system, stating that the time-scale of the transient escapes from the steady state might not be long enough to explain the time spent in the Nanog-low state.

Similarly, the second theoretical model proposed to explain the bimodality of Nanog expression [139] involves both positive regulation of Nanog by an Oct4-Sox2 dimer, and negative feedback on Nanog by an additional unknown factor X. Since positive feedback loops alone are sufficient to generate bistable distributions of expression, the model we present in Section 2.3 eschews the unnecessary addition of negative feedback loops, and is able to explain the experimental data without recourse to unknown mechanisms. Another notable source of variability we exclude from our model is cell-cycle phenomena. Although the expression level of some gene products can depend on the stage of the cell cycle [150–152], there is evidence to suggest that we can eliminate cell-cycle phenomena as an explanation for the observed variability in Nanog abundance. In one experiment, time-lapse microscopy of individual ESCs visually captured the stochastic nature of transitions between Nanog-low and Nanog-high states, revealing that there was no temporal (or spatial) pattern to the onset of Nanog expression [108]. This suggests that the fluctuations of Nanog levels are cell cycle independent. In addition, the lack of a spatial pattern supports the assumption that, with respect to Nanog expression, ESCs are independent of their neighbours. Two further experiments have demonstrated that the reconstitution of the steady-state distribution of Nanog expression from isolated subpopulations takes place on a timescale that is of an order of magnitude longer than the cell cycle time [108, 117]. Therefore, it seems that reconstitution cannot be entirely a cell cycle phenomena, because if it were, we would expect it to happen more quickly than actually it does.

Accordingly, we present a simple model of positive feedback for Nanog expression. We then use this model to explain how the Nanog regulatory system might be perturbed by the heterozygous knock-in strategy, and explore the possibility that the observed heterogeneity is a reporter artefact, as opposed to a biologically significant phenomenon.

## 2.2    Network Structure

The first of our studies on the ESC network explores the role of feedback structure. We begin by describing different types of small network structures and how they relate to dynamics, and then we investigate the role of Nanog in the global feedback structure in the extended ESC TRN.

Figure 2.3: Examples of network motifs consisting of three genes: a) an undirected cycle; b) a self-enhancing positive feedback loop; c) a negative feedback loop; d) a coherent feedforward loop; e) an incoherent feedforward loop. Activation of gene expression is represented by an arrow, and inhibition by a T-bar.

### 2.2.1 Network Motifs

Although complex in their overall structure [67, 153], molecular regulatory networks often contain certain types of small subnetworks at frequencies higher than expected by chance [59, 154], suggesting that these structural building blocks – or *motifs* as they are known – may perform specific regulatory functions [59]. Here, we shall consider the dynamics of three types of commonly occurring, and dynamically significant, motifs: positive feedback, negative feedback and feedforward loops. Examples of these motifs are shown in Fig. 2.3.

**Positive Feedback Loops**

In general, a necessary condition for bistability (or multistability -a dynamical system that supports more than one coexisting attractor is said to be *multistable*) is the presence of at least one *positive feedback loop* somewhere in the underlying regulatory network [155–157]. A feedback loop is a closed path in a network from a node back to itself in which each intermediary node is visited once. A positive feedback loop is one in which the net effect of the entire loop is positive: using the convention that activating links are denoted by $+1$ and inhibiting links by $-1$, a positive feedback loop is one in which the product of the link-signs is $+1$. A gene that activates its own transcription directly is a positive feedback loop of length 1. A positive feedback loop of length 3 is shown in Fig. 2.3b: in this case, all links are positive and the sign of the loop is therefore $(+1) \times (+1) \times (+1) = +1$. Since all the links are positive, this example feedback loop is *self-enhancing*. Once activated (perhaps in response to a transient activating signal), a self-enhancing feedback loop maintains the expression of all the genes in the loop. This kind of positive feedback often provides the molecular basis for irreversible switches which, by initiating all-or-none cell fate decisions, are important in cellular differentiation and development [132, 158, 159].

A common motif observed in genetic regulatory networks is a pair of mutually repressing genes, as depicted in Fig. 2.4(a) (this is a positive feedback loop of length 2 since its sign is $(-1) \times (-1) = +1$). This motif gives rise to a *toggle switch* [160], since it allows the cell to switch (or toggle) between two different states in response to a transient signal. The key property of this switch is that the mutual repression between the two genes does not allow them both to be co-expressed at high levels. Consider two attractors: $X$ (in which the first gene is active and, due to the repression exerted by gene 1, the second is inactive), $Y$ (in which the second gene is active and, due to the repression exerted by gene 2, the first is inactive), and a signal of magnitude $S$. Now suppose that if the system starts in the vicinity of $X$ it can be driven out of the basin of attraction of $X$ and into the basin of attraction of $Y$ in response to a signal that exceeds a critical value $S_{HIGH}$ in magnitude. Similarly, suppose that if the system starts in the vicinity of $Y$ it can switch to $X$ if the signal magnitude falls below $S_{LOW}$. For intermediate signal magnitudes ($S_{LOW} < S < S_{HIGH}$), the system is bistable since it admits two coexisting attractors. Thus, varying $S$ allows the system to toggle between the two alternative attractor states.

A particularly elegant example was designed and experimentally implemented in *Escherichia coli* by Gardner et al. [160], where the dynamics of the toggle switch are described by the following dimensionless ODE model:

$$\frac{dx}{dt} = \frac{\alpha_1}{1 + y^{\beta_1}} - x \tag{2.1}$$

$$\frac{dy}{dt} = \frac{\alpha_2}{1 + x^{\beta_2}} - y \tag{2.2}$$

where $x$ and $y$ are the concentrations of the repressors, $X$ and $Y$, $\alpha_1$ and $\alpha_2$ are the effective rates of synthesis, and $\beta_1$ and $\beta_2$ are levels of cooperativity of repression. Bistability is possible when $\beta_1, \beta_2 > 1$, and $\alpha_1, \alpha_2$ (the rates of synthesis) and $\beta_1, \beta_2$ (levels of cooperativity) are sufficiently close. In this case, there are three fixed points, two of which are stable and one of which is unstable. Since there are exactly two molecular species in this system, it can be studied using phase plane analysis. The solution curves in phase plane are shown in Fig. 2.4b. The first plot shows an example of where the conditions for bistability are met, so the system is drawn to one of two stable steady states, depending on the initial conditions. The separatrix (red line) divides the plane into two basins of attraction. In the second plot, the rate of synthesis of $Y$ is half that of $X$ so it is overwhelmed by its repressor and there is only one stable steady state. In the third plot the rates of synthesis are equal but the

Figure 2.4: a) A toggle switch consisting of two mutually repressing transcription factors. b) Phase plane (blue arrows) and solution curves (black lines) for the system described by Eqs. (2.1) and (2.2) for parameter sets $\alpha_1 = \alpha_2 = 4, \beta_1 = \beta_2 = 2$ (bistable), $\alpha_1 = 4, \alpha_2 = 3$ and $\beta_1 = \beta_2 = 2$ (monostable), and $\alpha_1 = \alpha_2 = 4, \beta_1 = \beta_2 = 1$ (monostable). The stable states are denoted by the filled, red circles and the red line in the first plot marks the separatrix of the two basins of attraction. c) Transient induction of expression of a clonal population initially in the low state causes it to switch to the high state over a period of 6 hours [160].

cooperativity of repression is equal to 1, so the nullclines ($dx/dt = 0$ and $dy/dt = 0$) do not have a sigmoidal shape and therefore only intersect once, producing a single stable steady state. Gardner et al. implemented the toggle switch by constructing a corresponding DNA sequence (called a plasmid) and inserting it in *Escherichia coli*. It was engineered so that production could be artificially induced, essentially increasing the effective rate of synthesis. Fig 2.4c shows how the distribution of expression of one of the two genes in a clonal population initially in the low state changed when the inducer was transiently applied - corresponding to a temporary change in the parameter values. The plasmid started in the low state and 3 hours after induction the expression began switching to the high state (first panel). A bimodal distribution then appeared at 4 hours after induction (2nd panel); by 5 hours the switching was nearly complete (3rd panels), and by 6 hours it was complete (last panel). In the theoretical setting, the parameters of the system were changed and the phase plane changed accordingly, so that the system moves from one state to another. The toggle switch requires only transient rather than sustained induction, as the rate of synthesis remains high when the switch to the high state is complete. The agreement between the theoretical model and experiment indicates that the theoretical design and implementation of genetic networks is an achievable goal.

**Negative Feedback Loops**

A negative feedback loop is a feedback loop in which the net effect of the entire loop is negative (one in which the product of the link-signs is $-1$). An example of a negative feedback loop is shown in Fig. 2.3c: in this case, the sign of the loop as a whole is $(+1) \times (+1) \times (-1) = -1$. Homeostasis (the maintenance of a constant internal state despite environmental variations) may be maintained by negative feedback, since negative feedback loops, in general, suppress fluctuations [161]. Similarly, negative feedback loops can, by introducing time-delays [162] and associated over- and under-compensation in gene expression, give rise to self-sustaining oscillatory behaviour [132]. An example of a synthetic three-gene negative feedback loop that generates self-sustained oscillations in protein levels in *Escherichia coli* was presented by Elowitz and Leibler [141], who showed that this system can function as a biological clock by inducing periodic bursts of protein synthesis.

**Feedforward Loops: Persistence Detectors & Pulse Generators**

Feedforward loops occur when a source gene regulates the expression of a target gene through two different paths. Figures 2.3d & e show examples of three-node feedforward loops. In both these cases, $A$ regulates $C$ both directly and indirectly via $B$. Feedforward loops are common in molecular regulatory networks, including the transcriptional regulatory networks of *Escherichia coli*, yeast [59, 63, 163] and other organisms [58, 164, 165]. Each of the regulatory interactions in the feedforward loop can be either be activating or repressing: if both paths in the feedforward loop have the same overall sign (both activating or both inhibiting) then the feedforward loop is said to be *coherent*, otherwise it is *incoherent*. The feedforward loop in Fig. 2.3d is coherent because the sign of the direct path from $A$ to $C$ has the same sign (positive in this case) as the indirect path from $A$ to $C$ via $B$. The feedforward loop in Fig. 2.3e is incoherent because the sign of the direct path from $A$ to $C$ is positive, while the sign of the indirect path from $A$ to $C$ via $B$ is negative.

Coherent and incoherent feedforward loops exhibit different dynamics. Coherent feedforward loops can filter out transient environmental fluctuations and act as persistence detectors [163, 166]. For instance, consider the 3 node feedforward loop in Fig. 2.3d. If expression of both $A$ *and* $B$ is needed to activate $C$ (if $C$ is regulated by an $AB$-dimer, for example) then an activating signal starting at $A$ must persist long enough for the concentration of $B$ to

reach the activation threshold before $C$ is activated. In contrast, when the activating signal is removed, and $A$ is down-regulated, the expression of $C$ also down-regulates without delay. If expression of $C$ only requires expression of $A$ *or* $B$ ($A$ and $B$ regulate $C$ independently), then the opposite effect is observed: there is no delay in activation of $C$ after activation of $A$, but there is a delay in down-regulation of $C$ when stimulation of $A$ stops [166].

Incoherent feedforward loops can act as pulse generators. For example, consider the 3 node feedforward loop in Fig. 2.3e. Node $A$ both directly activates $C$ and indirectly represses $C$ by activating the repressor $B$. Consequently, when a signal activates $A$, the production of $C$ is also rapidly activated. However, over time, levels of $B$ also accumulate until they reach the repression threshold. At this point production of $C$ decreases and its concentration drops, resulting in pulse-like expression of $C$.

**Does Structure Determine Function?**

Although certain network structures can be associated with defined dynamics, caution should be exercised when determining the relationship between structure and function in more complex regulatory networks. A study performed by Ingram et al. [167] showed that the function of even very simple motifs cannot always be determined by their structure. The authors investigated the behaviour of the bi-fan motif – in which the products of two source genes directly co-regulate the expression of two target genes – and found that there is no characteristic behaviour for this motif: the bi-fan can exhibit a range of possible responses. Given that the bi-fan is only slightly more complex than a feedforward loop, the authors conclude that " ... it is difficult to gain significant insights into biological function simply by considering the connection architecture of a gene network, or its decomposition into simple structural motifs". They add that additional information, such as the values of the kinetic parameters, or experimental time series data is required to make inferences about network dynamics.

### 2.2.2 The Combined Effect of Positive Feedback and Stochasticity

Two important mathematical theorems have helped us to identify the source of cell-to-cell variability in ESC populations. First, in 1949, the physicist Max Delbruck linked biological systems with the knowledge that dynamical systems can move between coexisting equilibrium

states under the influence of noise [73]. Second, in 1981 René Thomas presented a mathematical theorem now known as Thomas' Rule [155]. Informally, the rule asserts that for a complex dynamical system to have multiple coexisting stable equilibrium states, a positive feedback loop must be present in the interaction network. Formally, the interaction graph of a dynamical system is defined using the matrix of signs of the Jacobian matrix of the system. Consider the system $\frac{d\boldsymbol{x}}{dt} = f(\boldsymbol{x})$. In the interaction network of the system, there is a positive link from node $j$ to $i$ (e.g. promotion of transcription represented by an arrow) if there exists $\boldsymbol{x} \in \mathbb{R}^n$ such that $\frac{\partial f_i}{\partial x_j}(\boldsymbol{x})$ is positive, and a negative link from node $j$ to $i$ (e.g. repression of transcription represented by a T-bar) if $\frac{\partial f_i}{\partial x_j}(\boldsymbol{x})$ is negative, $(i, j = 1, \ldots, n)$ (the network can thus have both a positive and a negative link from one node to another) [168]. In such directed networks, a feedback loop is a closed path from a node back to itself in which each intermediary node is visited once. A positive (negative) feedback loop is one in which contains an even (odd) number of negative links (see Section 2.2.1). Thomas' rule states that if the system $\frac{d\boldsymbol{x}}{dt} = f(\boldsymbol{x})$ has several stable states, then the interaction network of the system has a positive feedback loop.

Combining the insights of Delbruck and Thomas, it can be seen that the combination of positive feedback and random fluctuations can give rise to phenotypic heterogeneity in an isogenic population [99, 169–171].

However, we should note that this rule applies to deterministic systems only; a stochastic system does not require the presence of a positive feedback loop a to yield a bimodal stationary state. Consider the following system which describes production of a species, $M$, from a gene that can switch randomly between active ($G_A$) and inactive ($G_I$) states, where $r_A$ ($r_I$) is the rate of switching out of the active (inactive) state. The terms inactive or active are simply used to label the two states since production can occur in either state. Specifically, when the gene is active (inactive), $M$ is produced at a rate $k_A$ ($k_I$), where $k_I << k_A$, and decays at rate $d$. In this case, slow switching between states – switching on a timescale that is slower than the timescales involved in production and decay – can result in a bimodal distribution of levels of $M$ at the population levels. This chemical reaction system can be described by the scheme

$$G_A \xrightarrow{r_A} G_I, \quad G_I \xrightarrow{r_I} G_A, \quad G_A \xrightarrow{k_A} G_A + M, \quad G_I \xrightarrow{k_I} G_I + M, \quad M \xrightarrow{d} \phi. \quad (2.3)$$

A simulation of the time series of the abundance of $M$, obtained using Gillespie's SSA, is

Figure 2.5: Dynamics of a chemical reaction system described by (2.3), and the corresponding stationary PMF, for parameter values $r_A = 2.5 \times 10^{-5}$, $r_I = 10^{-3}$ $k_A = 1.2$, $k_I = 0.2$, $d = 0.04$. The bistability arises from the slow switching between active and inactive states.

given in Fig. 2.5, together with the resulting bimodal equilibrium distribution. This simple two-state model demonstrates how a bimodal distribution of expression can arise without feedback if there is a strong separation of timescales.

### 2.2.3 The Importance of Nanog in the Global Feedback Structure

In this section, we investigate the role of Nanog in the global feedback structure in the extended ESC TRN. This work makes up part of the study by MacArthur et al. [3], in which we aim to elucidate the molecular basis for the observed fluctuations in expression in the constituent genes (see Section 2.1.1 for details). Since feedback loops commonly regulate cell-to-cell variability by generating complex dynamics, we analyse the feedback structure of the extended ESC TRN to investigate the possibility that Nanog fluctuations regulate population heterogeneity by controlling the feedback mechanisms. To do this, we compare the feedback architecture of the ESC TRN in NanogR cells (Fig. 2.7b) and in wild-type ESCs (Fig. 2.7a). The wild-type ESC TRN is self-perpetuating when shielded from differentiation-inducing stimuli [130]. However, in the NanogR cell line endogenous regulation of the Nanog gene does not contribute to Nanog protein levels. Consequently, all feedback loops that involve Nanog in the wild-type TRN are absent in the NanogR cells. In these cells the ESC TRN is therefore effectively held in a feedback-depleted state (Fig. 2.7b), and maintenance of pluripotency is dependent on continued exogenous expression with dox of Nanog rather than activation of self-perpetuating feedback loops.

As such, we enumerate the feedback loops that each transcription factor participates in, for both the wild-type ESC TRN and the NanogR TRN. In addition, we calculate a feedback index, which takes into account both the total number and the lengths of all closed paths present in the network, thus providing an adjusted measure of node involvement in the feedback structure in the extended ESC TRN. We now describe the methods used to carry out these analyses.

**Enumeration of the Feedback Loops**

A feedback loop (or cycle) of size $L$ is a closed path of $L$ links that starts and finishes at the same node and visits each intermediate node exactly once, and each link exactly once. First, we note that the longest feedback loop in the extended ESC TRN has a length of 5, because only 5 nodes in the network (Nanog, Oct4, Sox2, Dax1 and Rex1) have both incoming and outgoing edges. As the network is small, specific feedback loops may easily be found by exhaustive enumeration using a simple computational algorithm (see Appendix C for details). The total number obtained was verified with the adjacency matrix method of Harary and Manvel [172], which provides exact formulae for feedback loops of length 2 to 5. Both methods require as an input the adjacency matrix $\boldsymbol{A}$, which for an unsigned network $D$, with $p$ nodes, is defined as the $p \times p$ matrix with entry $a_{ij} = 1$ if there is a directed link from node $i$ to node $j$, and 0 otherwise. For example, for the wild-type ESC TRN, the adjacency matrix $\boldsymbol{A}$ is given by

$$\boldsymbol{A} = \begin{pmatrix} 1 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}. \tag{2.4}$$

A well-known result in graph theory is that the $i,j$ entry, $a_{ij}^{(L)}$, of $\boldsymbol{A}^L$ is the number of paths of length $L$ between two nodes $i$ and $j$ [173]. In particular, the diagonal entry, $a_{ii}^{(L)}$, is the number of paths of length $L$ that start and finish at node $i$ (closed paths). Therefore, the total number of closed paths of length $L$ is given by $\frac{1}{L}Tr(\boldsymbol{A}^L)$. The division by $L$ is necessary because each node that participates in the loop contributes once to the sum. However, this total includes closed paths that are not cycles. Thus, the number of cycles of length $L$ in a

Figure 2.6: a) A self loop is a closed walk of length 1. b) to k) All possible structures of closed walks of lengths 2 to 5 in a directed network without self-loops.

directed network, denoted $N_L$, is given by the total number of closed paths of length $L$ minus the number of closed paths of length $L$ that

1. visit at least one intermediate node more than once;

2. visit the start node more than twice; or

3. traverse at least one link more than once.

Now, the trace of the adjacency matrix is the number of self loops (Fig. 2.6a), i.e.

$$N_1 = Tr(\boldsymbol{A}),$$

and for $L \geq 2$ we can eliminate closed paths that are of type 1 and 2 due to the presence of self loops by setting the diagonals of the adjacency matrix equal to zero i.e. $a_{ii} = 0 \ \forall \ i$. Having eliminated self loops, the only possible closed paths of length 2 or 3 are cycles of the same length (Figs. 2.6b and c), and therefore

$$N_2 = \frac{1}{2}Tr(\boldsymbol{A}^2),$$

$$N_3 = \frac{1}{3}Tr(\boldsymbol{A}^3).$$

Figs. 2.6d to f show all possible closed paths of length 4, the first of which is the only cycle. Since $Tr(\boldsymbol{A^4})$ counts *all* closed walks in $D$, we must subtract the number of ways a closed path of length 4 can be achieved on the second and third structures. Now, structure (e) is two pairs of adjacent cycles of length 2, and structure (f) is one cycle of length 2, traversed twice. To count the number of times each of these structures occurs in a directed network that does not contain self loops, we first define $\bar{\boldsymbol{A}}$ to be the adjacency matrix of the undirected network, $U$, derived from the original directed network, $D$, by calling two nodes adjacent if and only if they form a cycle of length 2. Therefore, the entries of $\bar{\boldsymbol{A}}$ are defined by $\bar{a}_{ij} = a_{ij}a_{ji}$, where, $\bar{a}_{ij}^{(L)}$ denotes the $i, j$ entry of $\bar{\boldsymbol{A}}^L$. For example, for the wild-type ESC TRN, $\bar{\boldsymbol{A}}$ is given by

$$
\boldsymbol{A} = \begin{pmatrix}
1 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\
1 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0
\end{pmatrix} . \tag{2.5}
$$

For structure (e), the number of pairs of adjacent cycles of length 2 in $D$ is equal to the number of pairs of adjacent links in $U$. There are $\frac{1}{2}\sum\limits_{i\neq j} \bar{a}_{ij}^{(2)}$ pairs of adjacent links in $U$, where the $\frac{1}{2}$ cancels the duplication from each of its two endpoints - the number of pairs of adjacent links that start at node $i$ and end at node $j$, $a_{ij}$, accounts for the same pairs of adjacent links that start at node $j$ and end at node $i$, $a_{ji}$. However, each pair of adjacent links is a closed walk of length 4 in four different ways: once for starting at each end node (node A or C in Fig 2.6e), and twice for starting the middle node (node B, from which the walk can then visit either node A or C). Therefore, the number of closed walks of length 4 on structure (e) is $2\sum\limits_{i\neq j} \bar{a}_{ij}^{(2)}$.

Structure (f) in $D$ is a single link in $U$, traversed twice in each direction. There are $\frac{1}{2}\sum\limits_{i,j} \bar{a}_{ij}$ of these, where the $\frac{1}{2}$ cancels the duplication from each of its two endpoints. This structure is a closed walk of length 4 in two different ways - once for starting at each node - so the number of closed walks of length 4 on structure (f) is $\sum\limits_{i,j} \bar{a}_{ij}$.

After subtracting these from $Tr(\boldsymbol{A}^4)$ we are left with only closed walks of length 4 that are cycles, but each one is being counted at four different initial nodes, so we divide by four to

obtain the number of distinct cycles of length 4 in $D$ to obtain

$$N_4 = \frac{1}{4} \left[ Tr(\boldsymbol{A}^4) - \sum_{i,j} \bar{a}_{ij} - 2 \sum_{i \neq j} \bar{a}_{ij}^{(2)} \right].$$

Figs 2.6g to k show all possible closed paths of length 5, the first of which is the only cycle. Since $Tr(\boldsymbol{A^5})$ counts *all* closed walks in $D$, we must subtract the number of ways a closed path of length 5 can be achieved on the second and third structures. Both structures h and k are made up of a cycle of length 3 adjacent to a cycle of length 2. The number of cycles of length 3 starting from node $i$ is $a_{ii}^{(3)}$, and $\bar{a}_{ij}$ is the number of cycles of length 2 also starting at node $i$. Summing over all nodes in $D$ we obtain $\sum_i \left( a_{ii}^{(3)} \cdot \sum_j \bar{a}_{ij} \right)$. However, since structure k is symmetric, this term counts it twice (both nodes B and C in k are adjacent to a cycle of length 3 and length 2), so we must subtract the number of structures of type k. These can be counted by $\sum_{i \neq j} \left( a_{ij}^{(2)} \cdot \bar{a}_{ij} \right)$ because $a_{ij}^{(2)}$ counts the number of paths of length 2 starting at node $i$ and ending at node $j$ (CA, AB in structure k), and $\bar{a}_{ij}$ counts the number of cycles of length 2 that start at node $i$ and visit node $j$ (CB, BC in structure k). Since both structures h and k are a closed walk of length 5 in five different ways, so we obtain:

$$N_5 = \frac{1}{5} \left[ Tr(\boldsymbol{A}^5) - 5 \left\{ \sum_i \left( a_{ii}^{(3)} \cdot \sum_j \bar{a}_{ij} \right) - \sum_{i \neq j} \left( a_{ij}^{(2)} \cdot \bar{a}_{ij} \right) \right\} \right]$$

where the $\frac{1}{5}$ accounts for the fact that each cycle of length 5 is being counted at five different initial nodes.

**Feedback Centrality**

Feedback centrality given in [174, 175] is a measure of node involvement in both direct (non-intersecting) and indirect (self-intersecting) feedback. It takes into account both the total number and the lengths of all closed paths present in the network, such that shorter closed paths have more influence on the centrality of the node than longer closed paths. This rule is based on the following two observations:

- Consider a particle moving randomly along the paths of the network. It is more likely to return to the starting node and complete a closed walk if the path is of shorter length because it has fewer opportunities to exit the walk along the way. Thus shorter walks are typically more likely to occur than longer walks.

- Motifs in molecular regulatory networks (and other real-world networks) are typically small subgraphs (see discussion in Section 2.2.1).

The greater importance of shorter paths is taken into account by weighting the closed paths in decreasing order of their lengths. Thus the feedback centrality for the $i$th node, as given in refs. [174, 175], is

$$\sum_{L=0}^{\infty} \frac{1}{L!} a_{ii}^{(L)} - 1 = \exp(a_{ii}) - 1$$

where $\exp(a_{ij})$ is the matrix exponential of the network adjacency matrix. The $-1$ term is included for convenience to ensure that nodes that do not participate in any closed paths have an index of zero.

**Results**

We found that the extended ESC TRN is rich in feedback, containing a total of 28 distinct feedback loops. Furthermore, these feedback loops are not evenly distributed (Fig. 2.7c). Rather, the global feedback structure of this network is highly nested and is critically dependent on Nanog, Oct4 and Sox2, which participate in 68% (19/28), 68% (19/28) and 64% (18/28) of all feedback loops respectively.

The feedback centrality identified Nanog as the most central element in the global feedback structure (Fig. 2.7d). Removal of Nanog feedback destroys many of these feedback structures, leaving only 32% (9/28) of the feedback loops in tact (Fig. 2.7c, blue bars), and therefore severely compromises the global feedback structure. Consequently, fluctuations in Nanog expression levels can transiently activate the subnetwork shown in Fig. 2.7b in the ESC TRN, driving transitions between a (Nanog-expressing) feedback-rich, robust and self-perpetuating pluripotent state and a (Nanog-diminished), feedback-sparse and differentiation-sensitive state.

## 2.3 Stochastic Models of Transcriptional Regulation in a Perturbed Network

We now use mathematical models of the dynamics of the core ESC TRN to elucidate the effect of the heterozygous knock-in reporter strategy on Nanog protein expression, and demonstrate

Figure 2.7: a) The feedback-rich wild-type TRN. (b) The feedback-depleted NanogR TRN (+dox, Nanog active). (c) The total number of feedback loops that each transcription factor participates in for the wild-type ESC TRN (red) and NanogR TRN (blue). (d) Feedback centrality for the wild-type ESC TRN (red) and NanogR TRN (blue).

that the experimentally observed bistable distribution (see Fig. 2.2) might be a reporter artefact, as opposed to a biologically significant phenomenon.

We apply the theory introduced in Section 1.2 to develop continuous time Markov chain models. This provides us with a basis for understanding the characteristics that control the behaviour of the expression dynamics at a single cell level, and the resulting distribution at the population level. Furthermore, the models enable us to explain how the knock-in method perturbs the Nanog system, and predict the resulting change in the shape, location and modality of the distributions of expression.

Throughout the remainder of this chapter, the wild-type ESC-line, in which there are two *Nanog* alleles and no reporter, is referred to as 'Scenario 1', and 'Scenario 2' represents the

GFP knock-in cell-line in which there is one *Nanog* allele, and the coding region of the second *Nanog* allele is replaced with that for GFP.

### 2.3.1   Chemical Master Equation Model of the Core ESC Network

We begin by considering a CME model of the core ESC network (see Fig. 2.1b) in each scenario. In Scenario 1, there are 32 molecular species and 43 chemical reactions involved in the transcription and translation of Oct4, Sox2 and Nanog. In Scenario 2, there are two additional species and three additional reactions. The full specification of the models, including the stoichiometric matrix and propensity functions, can be found in Appendix A.1. The processes involved both models are illustrated in Fig. 2.8, and the differences between the two scenarios are labelled 'Scenario 1' and 'Scenario 2'. In Scenario 1, there are two *Nanog* alleles, labelled Nanog allele A and Nanog allele B. This represents the coding regions for Nanog mRNA. In Scenario 2, the DNA sequences on the promoter regions remain the same, but the coding region of *Nanog* allele B is replaced with that for GFP. This means that, given the same promoter binding state, GFP is transcribed at the same rate as Nanog. The GFP mRNA can either decay or be translated into GFP proteins, which in turn can decay. Since it is a GFP protein, it solely acts as a visual label and does not bind to any other molecular species, or influence any reactions they are involved in. We set the rates of GFP translation and decay equal to those of Nanog, since different rates would have an arbitrary effect on the distribution of GFP expression.

For simplicity, we let there be only one binding site for each of the two types of dimer on each promoter region. Each of the genes can be transcribed by RNAP, and the resulting mRNA can either decay or be translated into the corresponding protein. Nanog proteins (red-filled squares, labelled 'N') can form Nanog-Nanog homodimers (NN), and an Oct4 protein (yellow-filled circles labelled 'O') can bind to a Sox2 protein (blue-filled circles labelled 'S') to form an Oct4-Sox2 heterodimer (OS). The dimers can either disassociate, releasing their constituent proteins, or they can bind to a corresponding unbound binding site on the promoter region of any of the genes. Since a dimer is far more likely to disassociate than decay, we do not include dimer decay in the model. When a binding site is occupied, the rate of transcription of the corresponding gene is greater than when it is unoccupied, because the presence of the dimers promotes the recruitment of RNAP II. The dimers can unbind from the promoter site, which results in the reduction of the rate of transcription.

Figure 2.8: An illustration of the transcription and translation processes involved in the core embryonic stem cell network, which consists of the three master genes *Nanog* (red-filled boxes labelled 'N'), *Sox2* (blue-filled circles labelled 'S') and *Oct4* (yellow-filled circles labelled 'O'). There are two copies of each gene (two alleles), one on each chromosome. In Scenario 1 there are two Nanog alleles, and in Scenario 2 there is one *Nanog* allele and the second is replaced with the *GFP* coding region (the promoter region remains unchanged). Reactions (arrows) are labelled with rounded rectangles, together with the propensity functions.

Figure 2.9: Time series and distributions for Nanog, Sox2, Oct4 and GFP proteins generated by simulations of the master equations specified in Appendix A.1, and illustrated in Fig. 2.8 for the parameter values given by (2.6). Scenario 1 refers to the wild-type cell-line and Scenario 2 to the knock-in cell-line.

The corresponding CME cannot be solved analytically. To obtain the marginal probability mass function (PMF) for the copy number of Nanog, Sox2, Oct4 and GFP proteins we simulate the CME using Gillespie's algorithm. The following set of parameter values gives rise to bistability in the knock-in cell line, and homogeneity in the wild-type cell:

$$g_M^N = 0.01, \ g_M^S = 0.5, \ g_M^O = 0.5, \ k_2^{N/NN} = 10, \ k_2^{N/OS} = 3, \ k_2^{S/NN} = 2, \ k_2^{S/OS} = 15,$$

$$k_2^{O/NN} = 2, \ k_2^{O/OS}, \ d_M^N = 2, \ d_M^S = 5, \ d_M^O = 5, \ g_P^N = 5, \ g_P^S = 20, \ g_P^O = 20, \ d_P^N = 3$$

$$d_P^S = 5, \ d_P^O = 5, \ u^{NN} = 10, \ u^{OS} = 10, \ a^{NN} = 1, \ a^{OS} = 1, \ k_3^{N/NN} = 1, \ k_3^{N/OS} = 5,$$

$$k_3^{S/NN} = 1, \ k_3^{S/OS} = 10, \ k_3^{O/NN} = 1, \ k_3^{O/OS} = 10, \ k_1^{N/NN} = 1, \ k_1^{N/OS} = 0.3,$$

$$k_1^{S/NN} = 0.1, \ k_1^{S/OS} = 5, \ k_1^{O/NN} = 0.1, \ k_1^{O/OS} = 5. \tag{2.6}$$

These parameter values were chosen to demonstrate that it is possible to obtain bimodality in Nanog in the knock-in cell-line, whilst Nanog is unimodal in the wild-type cell, under the same parameter values. Very few of these parameter values have been experimentally determined, and they are dependent on the culture conditions, which themselves vary between batches. Therefore, guided by the available literature, we ensured that the orders of magnitude between the parameter values were realistic. For example, the background rate of mRNA transcription is typically at least an order of magnitude smaller than the enhanced rate of transcription due to a bound promoter region, and dimer association is at least an order smaller than disassociation.

The results of the simulations are shown in Fig. 2.9. A representative sample of the time series for each protein is plotted together with the corresponding equilibrium distributions. The simulations were initially run until all molecular species had reached their equilibrium behaviour (determined by visually inspecting the time series plots), from which point data for the Nanog, Oct4 and Sox2 proteins (and GFP for Scenario 2) was collected until the probability distributions reached their equilibrium state (determined by visual inspection). The time series are plotted for the first $10^3$ units of time, although the simulations were run for $10^5$ units of time to achieve the corresponding equilibrium probability distributions.

In Scenario 1, the distributions for all three proteins are unimodal. In Scenario 2, the distributions exhibit features similar to those observed in 2i culture conditions: Nanog and GFP are bimodal, where the low peaks are much smaller than the high peaks, and Oct4 and Sox2 are homogeneously expressed. At the single cell level, a cell can be in a Nanog-/Oct4-/Sox2-high

state and can switch into a Nanog-low state whilst Oct4 and Sox2 remain relatively constant. Crucially, the heterogeneity in the Nanog expression arises due to strong positive feedback via the Nanog homodimer, plus the intrinsic noise in the molecular reactions. Underlying these dynamics is a bistable system and the intrinsic noise drives it between the two distinct states.

We note that the probability distribution for GFP has slightly greater spread than that for Nanog, the reason for which is as follows: consider a cell that occupies the state of highest probability with respect to both GFP and Nanog. When the GFP level makes an excursion away from this state, it does not affect Nanog expression, but when the Nanog level makes such an excursion it does have an effect on GFP because it interacts with both the Nanog and GFP promoter region. Therefore, GFP has more freedom to vary without affecting the dynamics of Nanog, resulting in a probability distribution with greater spread for GFP than for Nanog.

Although the model agrees with the experimental data quantitatively, it is too complicated to enable us to identify the components that control stability. Since we are particularly interested in elucidating the underlying mechanisms that control the heterogeneity in Nanog expression, and how the substitution of one Nanog coding region with that of GFP might affect it, we now focus on the Nanog autoregulatory loop alone.

### 2.3.2   Models of the Nanog Autoregulatory Loop

We now explore a subset of the processes of the core ESC network that includes only the *Nanog* gene and its products. A simple model of the Nanog autoregulatory loop will capture the bimodal expression in Nanog, and by replacing the coding region of one of the two alleles with that of GFP we can see how the kinetics of the system can be disturbed by the knock-in method.

It is known that the heterogeneity in Nanog expression arises due to strong positive feedback via the Nanog homodimer, but is only weakly regulated by the Oct4-Sox2 dimer [129]. This is reflected in the simulation results from the full model, as illustrated in Fig. 2.9, which indicate that the dynamics of the Nanog protein are decoupled from those of Oct4 and Sox2. We therefore assume that the extrinsic noise in Nanog expression is due to fluctuations in Oct4 and Sox2 expression. In the reduced model for Scenario 1 there are 7 molecular

Figure 2.10: The processes involved in the Nanog auto-regulatory loop. Reactions (arrows) are labelled with rounded rectangles, together with the propensity functions. The variables $M, P$, and $C$ denote the number of Nanog mRNA, protein and homodimers; $D_A, D'_A$ and $D_B, D'_B$ the number of bound and unbound promoter sites on alleles A and B, respectively; and the variables $M_G$ and $P_G$ denote the number of GFP mRNA and proteins. The parameters $g_M, g_P, a, u, k_1, k_2, k_3, d_M, d_P$ denote the reaction rate constants.

species and 11 chemical reactions. We retain the Nanog mRNA, proteins, homodimer, and the (un)bound Nanog dimer binding sites on the promoter region of each of the two alleles. In Scenario 2 there are GFP mRNA and proteins in addition to the molecular species in Scenario 1, bringing the total number of molecular species to 9 and chemical reactions to 14. The processes involved in the Nanog autoregulatory loop are illustrated in Fig. 2.10. They are a subset of those in the core ESC network, as illustrated in Fig. 2.8.

Full details of CME model of the Nanog autoregulatory loop, including the propensity functions and stoichiometric matrix are given in Appendix A.2. Time series simulations of the model for both scenarios are shown in Fig. 2.11, together with the equilibrium probability distributions. As for the more complicated model, the simulations were initially run until all molecular species had reached their equilibrium behaviour (determined by visually inspecting the time series plots), from which point data for the Nanog protein (and GFP for Scenario 2) were collected until the probability distributions reached equilibrium (determined by visual

Figure 2.11: Time series simulations (top row) of the master equation for the Nanog feedback loop using Gillespie's algorithm, and the equilibrium probability distributions (bottom row). The first two columns refer to Scenario 1 using parameter values (2.7) (left column), and again except $k_3 \to k_3/2$ and $u \to u/2$ (second column). The former achieves bistability and the latter homogeneity. The second parameter set is also used for Scenario 2, and the translation and decay rates for GFP are identical to that for Nanog (last column).

inspection). The time series are plotted for the first $10^3$ units of time, although the simulations were run for $10^5$ units of time to achieve the corresponding equilibrium probability distributions.

Simulation results are shown for two parameter sets under Scenario 1, the first of which is:

$$g_M = 10, \ k_2 = 120, \ d_M = 3.8, \ d_P = 8.1, \ u = 105, \ k_3 = 105, \ g_P = 35, \ a = 1, \ k_1 = 1 \quad (2.7)$$

and demonstrates that the model can achieve bistability (first column). In the second simulation (second column) the same parameters were used except the values of both $k_3$ and $u$ were halved, resulting in a unimodal distribution. The simulation of Scenario 2 (last column) uses exactly the same parameter values as for the second simulation of Scenario 1, with the addition of translation and decay rates for GFP, equal to those for Nanog. We see again - as was the case for the model of the core ESC TRN under Scenario 2 - that the probability distribution for GFP has greater spread than that for Nanog, and occurs for the same reasons as previously discussed.

Thus, we have demonstrated theoretically that the modelled GFP knockout strategy results in a bimodal distribution of both GFP and Nanog expression, when the distribution was unimodal in the wild-type system with the two Nanog alleles (middle column). However, although these simulations suggest that the introduction of a GFP reporter can affect Nanog dynamics, the mechanism by which this occurs is unclear. We can better understand how this behaviour change occurs by studying the corresponding RREs for the system.

**The Reaction Rate Equations**

From Eq. (1.23), the corresponding RREs that approximate this stochastic system can be obtained by multiplying the stoichiometric matrix, $\boldsymbol{S}$, Eq. (A.1), by the propensity vector $\boldsymbol{a}(\boldsymbol{x})$, Eq. (A.2). Doing so we obtain the following system of RREs for Scenario 1:

$$\frac{dM}{dt} = 2g_M + k_2(D_A + D_B) - d_M M \tag{2.8}$$

$$\frac{dP}{dt} = g_P M - d_P P - aP^2 + 2uC \tag{2.9}$$

$$\frac{dC}{dt} = aP^2/2 - uC - k_1 C(D'_A + D'_B) + k_3(D_A + D_B) \tag{2.10}$$

$$\frac{dD'_A}{dt} = -k_1 C D'_A + k_3 D_A \tag{2.11}$$

$$\frac{dD'_B}{dt} = -k_1 C D'_B + k_3 D_B \tag{2.12}$$

$$\frac{dD_A}{dt} = k_1 C D'_A - k_3 D_A \tag{2.13}$$

$$\frac{dD_B}{dt} = k_1 C D'_B - k_3 D_B \tag{2.14}$$

where $D'_A + D_A = 1$, $D'_B + D_B = 1$, and $M(0)$, $P(0)$, $C(0)$, $D'_A(0)$, $D'_B(0)$, $D_A(0)$, and $D_B(0)$ are the initial conditions.

In Scenario 2, Eq. (2.8) in Scenario 1 is replaced by

$$\frac{dM}{dt} = g_M + k_2 D_A - d_M M \tag{2.15}$$

The remaining RREs in Scenario 1 are unchanged and we gain the following two equations for GFP mRNA and proteins:

$$\frac{dM_G}{dt} = g_M + k_2 D_B - d_M M_G \tag{2.16}$$

$$\frac{dP_G}{dt} = g_P M_G - d_P P_G \tag{2.17}$$

with initial conditions $M_G(0)$ and $P_G(0)$.

We can reduce the number of dimensions of both systems of ODEs to reveal the difference in the dynamics of Nanog expression between Scenarios 1 and 2. By assuming that dimer (dis)association, and DNA (un)binding occurs at a much faster time scale than transcription and translation, and using a quasi-equilibrium approximation, we obtain the following RRE for the Nanog protein level under Scenario 1 (see Appendix A.3 for the full derivation):

$$\frac{dP}{dt} = 2g_M \frac{g_P}{d_P} + \frac{2k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P \tag{2.18}$$

and for Nanog and GFP protein levels in Scenario 2:

$$\frac{dP}{dt} = g_M \frac{g_P}{d_P} + \frac{k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P \tag{2.19}$$

$$\frac{dP_G}{dt} = g_M \frac{g_P}{d_P} + \frac{k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P_G \tag{2.20}$$

The first term on the right-hand side of Eqs. (2.18) to (2.20) accounts for protein production at a constant baseline rate. This term for Nanog in Scenario 2 is half that in Scenario 1 (as highlighted by the 2 in red font in Eq. (2.18)) because there are two *Nanog* alleles in Scenario 1 and only one in Scenario 2, thus halving the output of Nanog mRNA. The second term on the right-hand side of the RREs is a Hill function that accounts for increased production due to a nonlinear positive feedback loop. Its value ranges from 0 to $2k_2 \frac{g_P}{d_P}$ in Scenario 1, and to $k_2 \frac{g_P}{d_P}$ in Scenario 2, and increases in a continuous manner with expression level. Thus, $2k_2 \frac{g_P}{d_P}$ is the maximum rate of production of Nanog in Scenario 1, which again is halved in Scenario 2, since the opportunities for mRNA production are shared equally with that for GFP. The aggregate parameter $\frac{2uk_3}{ak_1}$ is the protein concentration at which the half-maximum rate of production occurs. The third term on the RHS of the equations says that both $P$ and $P_G$ decay linearly (with identical constant half-life). In conclusion, it has become clear that, on the level of the RREs, replacing one of the Nanog genes with GFP has the effect of halving the rate of production. Since Eq. (2.19) decouples from the rest of the system, we can now carry out a stability analysis on Eqs. (2.18) and (2.19) to determine how the reduced rate of production influences the behaviour of Nanog expression at the population level.

There are nine free parameters in each of the scenarios, so we begin the stability analysis by rescaling the systems in order to aggregate these parameters and simplify calculations. Eqs. (2.18) and (2.19) may then be expressed in nondimensional form (see Appendix A.4 for a full

derivation):

$$\frac{dX}{dt} = \alpha + \frac{X^2}{\gamma^2 + X^2} - X \tag{2.21}$$

where the dimensionless parameters are defined as

$$\alpha = \frac{g_M}{k_2}, \quad \gamma = \begin{cases} \dfrac{d_M d_P}{2k_2 g_P}\sqrt{\dfrac{2uk_3}{ak_1}} = \gamma_1 & \text{(Scenario 1)} \\[4mm] \dfrac{d_M d_P}{k_2 g_P}\sqrt{\dfrac{2uk_3}{ak_1}} = \gamma_2 = 2\gamma_1 & \text{(Scenario 2)} \end{cases}$$

Thus $\alpha$ is the same for both scenarios, but the GFP knockout has the effect of doubling $\gamma$.

By writing the system in a dimensionless form it becomes clear that the dynamics are governed by two dimensionless parameters, $\alpha \geq 0$ and $\gamma \geq 0$. The latter may be thought of, phenomenologically, as governing the strength of the auto-activatory feedback loop: at $\gamma = 0$ the feedback loop is fully active and production of Nanog occurs at its maximal rate, while as $\gamma \to \infty$ feedback is inhibited and production tends to zero. The size of $\gamma$, and thus the strength of the feedback loop, may vary depending on internal or external signals and the dynamics of the system change as $\gamma$ is varied.

In the absence of time-delays and noise, sustained oscillations (or more exotic stable dynamic behaviour) can only occur when there is more than one species present. Since this model is one-dimensional (the RRE for GFP is decoupled) we do not need to look for these here, so we look for fixed-point equilibrium solutions. Fixed points are solutions to the cubic

$$X^3 - (1 + \alpha)X^2 + \gamma^2 X - \alpha\gamma^2 = 0 \tag{2.22}$$

We can ascertain the stability of the fixed points by plotting the curve $\dfrac{X^2}{\gamma^2 + X^2}$ and the straight line $X - \alpha$ on the same plot, for varying values of $\gamma$, as shown in Fig. 2.12. We summarise the direction of flow from the sign of $f(x)$, as indicated by the arrows marked on the $X$-axes (this is the phase portrait). Fig. 2.12 shows that the system may be monostable or bistable depending on the values of $\alpha$ and $\gamma$.

To summarise the behaviour of the system (2.21) we divide the $\alpha - \gamma$ parameter plane into distinct behavioural regimes, by calculating the values of $\alpha$ and $\gamma$ at which the bifurcations take place. When $\gamma(\alpha)$ is equal to either of the critical values, $\gamma_- = \gamma_-(\alpha)$, and $\gamma_+ = \gamma_+(\alpha)$, there are three real fixed points, two of which are repeated roots. We can find these critical values by using the fact that the discriminant, $D$, of a cubic is equal to zero if and only if

Figure 2.12: The Hill function $\frac{X^2}{\gamma^2+X^2}$ (solid line), and the function $X - \alpha$ for $\alpha = 0.05$ (dashed line), for $\gamma = 0.3162$, $\gamma_-$, $0.5196$, $0.5385$, $\gamma_+$, $0.7071$, where $0.3162 < \gamma_- < 0.5196$ and $0.5385 < \gamma_+ < 0.7071$. The filled circles mark the fixed-points. The arrows on the horizontal axes indicate the direction of flow between the fixed-points. When $\alpha = 0.05$, this system exhibits bistability for $\gamma_- < \gamma < \gamma_+$.

there are two repeated real roots and one other distinct real root. Given a cubic of the form $ax^3 + bx^2 + cx + d = 0$, the corresponding discriminant is

$$D = b^2c^2 - 4b^3d - 4ac^3 + 18abcd - 27a^2d^2 \tag{2.23}$$

For Equation (2.22), this gives:

$$\gamma^2 + (2\alpha^2 - 5\alpha - \frac{1}{4})\gamma + \alpha(1 + \alpha)^3 = 0 \tag{2.24}$$

which is a quadratic in $\gamma$ with solutions

$$\gamma_{\pm}(\alpha) = -\left(\alpha^2 - \frac{5}{2}\alpha - \frac{1}{8}\right) \pm \left(-2\left(\alpha - \frac{1}{8}\right)\right)^{\frac{3}{2}}. \tag{2.25}$$

This function may be used to divide the $\alpha$-$\gamma$ parameter plane into three regions, for $0 < \alpha < \frac{1}{8}$, $0 < \gamma < \frac{27}{64}$, as shown in Fig. 2.13, left. The values of the scaled fixed points over the parameter plane (computed numerically) are plotted as three surfaces in Fig. 2.13 (right).

Figure 2.13: Left: The $\alpha$-$\gamma$ parameter plane divided into three regions as given by Eq. (2.25) (solid curves). The monostable Nanog high region can be divided into a further three regions as shown by the different patterned areas (see main text for details). The white arrows indicate the region to which a wild-type cell (Scenario 1) originally in the ones of these areas of the Nanog high region transitions as a result of the insertion of the GFP knock-in (Scenario 2). Right: Positions of the scaled fixed points points on the $\alpha$-$\gamma$ parameter plane.

The bifurcation point, $(x_\pm, \alpha_\pm, \gamma_\pm)$ is a saddle node bifurcation for $f(x)$, since for $f(x_\pm, \alpha_\pm, \gamma_\pm) = 0$ the following conditions are satisfied [176]:

$$\frac{\partial f}{\partial x}(x_\pm, \alpha_\pm, \gamma_\pm) = 0, \quad \frac{\partial^2 f}{\partial x^2}(x_\pm, \alpha_\pm, \gamma_\pm) \neq 0, \quad \frac{\partial f}{\partial \alpha}(x_\pm, \alpha_\pm, \gamma_\pm) \neq 0, \quad \frac{\partial f}{\partial \gamma}(x_\pm, \alpha_\pm, \gamma_\pm) \neq 0$$
$$(2.26)$$

The monostable Nanog high region can be further divided into three regions: The values of $\gamma$ for which the insertion of the GFP knock-in (Scenario 2) in the wild-type cell (Scenario 1) will result in a transition of the system from the Nanog high region (yellow area) to:

1. the bistable region; these are the values of $\gamma$ that satisfy $\frac{\gamma_-(\alpha)}{2} < \gamma(\alpha) < \frac{\gamma_+(\alpha)}{2}$, as indicated by the unpatterned area in the monostable-high regime (plain yellow) in Fig. 2.13, left;

2. within the Nanog-high region; these are the values of $\gamma$ that satisfy $\gamma(\alpha) < \frac{\gamma_-(\alpha)}{2}$, as indicated by the vertical hatched area;

3. the monostable Nanog-low region; these are the values of $\gamma$ that satisfy $\gamma(\alpha) > \frac{\gamma_+(\alpha)}{2}$, as indicated by the dotted region.

Wild-type cells that reside in the unpatterned area of the Nanog high region are at risk of a qualitative change in population behaviour, and those in either of the patterned areas are at risk of quantitative change. In contrast, wild-type cells that reside anywhere in the bistable region are at risk of a qualitative change in population behaviour, as doubling the value of $\gamma$ would result in either a shift into the monostable low region, or a significant change in the proportion of cells in each state.

These analyses demonstrate that the introduction of the GFP knock-in reporter can induce a qualitative change in Nanog concentration that alters the behaviour and therefore the function of the cell. In conclusion, we have proposed a model that, on the RRE level, reveals a basis for which the experimentally observed bimodal distributions of Nanog expression can arise and how the introduction of heterozygous knock-in reporters can affect Nanog dynamics.

## A Stochastic Differential Equation Model

In Section 2.1.2, we provided an example of an experimentally observed bimodal distribution of Nanog protein expression in 2i culture conditions, as reported by the GFP knock-in cell line (Fig. 2.2). We now fit a model to this empirical distribution that will allow us to infer the shape of the distribution of Nanog protein expression in the wild-type cell, given the assumption that the translation and decay rates of GFP mRNA and proteins are the same as those of Nanog. Instead of deriving a complicated bivariate SDE from the CME, we fit the simplest possible model that takes into account the presence of positive feedback loop(s) in the underlying regulatory network. Since we have demonstrated that positive feedback can give rise to a Hill function, we use the form of Eq. (2.20) to develop an SDE to describe the dynamics of Nanog protein expression. Since experimental data suggests that protein expression fluctuations often scale linearly with expression level [177], a natural choice for the noise term is $\sqrt{2\sigma X^2}$, where $\sigma$ is a constant. Therefore, we use the following SDE to describe the dynamics of Nanog protein expression, $x(t)$:

$$\frac{dx}{dt} = \alpha_0 + \frac{\alpha_1 x^2}{K^2 + x^2} - \beta x + \sqrt{2\sigma x^2}\xi(t) \qquad (2.27)$$

Thus, from Eq. (1.12), the stationary distribution, $p_\infty(x)$, is

$$p_\infty(x) = A \exp\left[-\frac{\alpha_0^*}{x} + \frac{\alpha_1^*}{K}\arctan\left(\frac{x}{K}\right) - (2 + \beta^*)\ln(x)\right] \qquad (2.28)$$

where A is a normalising constant which ensures that $p_\infty(x)$ is a proper probability distribution, and $\alpha_{0,1}^* = \dfrac{\alpha_{0,1}}{\sigma}$, $\beta^* = \dfrac{\beta}{\sigma}$.

Had we used an additive noise term, the fixed points of the deterministic part of the SDE would correspond to the positions of the maxima and minima of the stationary distribution. This means that we could estimate the location of the empirical distribution on the $\alpha$-$\gamma$ plane in Fig. 2.13 using its fitted values of $\alpha$ and $\gamma$. However, the maxima and minima of the stationary distribution corresponding to the SDE obtained by the addition of a multiplicative noise term do not correspond to the fixed points of the deterministic part of the SDE. Thus, we now determine how the parameters of our chosen SDE are related to the behavioural regimes by finding the maxima and minima of Eq. (2.28) as a function of its parameters. Differentiating Eq. (2.28) with respect to $x$ and setting the result equal to zero, we obtain

$$x^3 - \frac{\alpha_0^* + \alpha_1^*}{2 + \beta^*}x^2 + K^2 x - \frac{\alpha_0^* K^2}{2 + \beta^*} = 0.$$

Applying the scaling $x = \dfrac{\alpha_1^*}{2 + \beta^*}X$, we obtain the familiar cubic, Eq. (2.22),:

$$x^3 - (1 + \alpha)x^2 + \gamma^2 x - \alpha\gamma^2 = 0$$

where $\alpha = \frac{\alpha_0^*}{\alpha_1^*}$, $\gamma = \frac{K(\beta^*+2)}{\alpha_1^*}$. Since this is the cubic that was used to divide the $\alpha\gamma$ plane in Fig. 2.13, we can use this plane to estimate the location of the empirical distribution, and predict that of Nanog in the wild-type cell.

We now consider how to fit the model to the empirical distribution. Since noise scales with abundance, expression levels are heavily skewed on a linear scale, and the low population is dominated by the high population. This is why experimentally observed histograms of expression levels are conventionally viewed on a $\log_{10}$-scale. We perform the model fitting procedure on the $\log_{10}$-scale to enable a good fit to the small Nanog low population. Thus we apply the change-of-variable technique to obtain the following expression for the probability density for the log Nanog protein expression, $y(t)$:

$$p_\infty(y) = A' \exp\left[-\alpha_0^* 10^{-y} + \frac{\alpha_1^*}{K}\arctan\left(\frac{10^y}{K}\right) - \ln 10(1 + \beta^*)y\right] \tag{2.29}$$

Maximum likelihood estimation was used to fit this model to the experimentally observed distribution under 2i culture conditions. The observed and fitted distributions are shown in Fig. 2.14, and the parameter estimates are $\hat{\alpha}_0^* = 2412$, $\hat{\alpha}_1^* = 279370$, $\hat{K} = 13889$, $\hat{\beta}^* = 5.085$.

Figure 2.14: Emprical distribution of Nanog expression in the knock-in cell line (thin black, solid line), fitted stationary distribution given by (2.29) (thick green line), and the inferred distribution of Nanog expression in the wild-type cell (red dotted line), obtained by multiplying the fitted production parameters by 2 as per Eq. (2.18).

The fitted model places the empirical distribution at $(0.0086, 0.3522)$ on the $\alpha$-$\gamma$ plane, and the predicted value of $\gamma$ in the wild-type cell, in 2i conditions, is $\hat{\gamma}/2 = 0.1761$, placing it in the monostable high region. Multiplying the production parameter estimates by 2 as per Eq. (2.18), i.e. $\hat{\alpha}^*_{0,1} \to 2\hat{\alpha}^*_{0,1}$, we can obtain an estimate of the distribution of Nanog expression in the wild-type cell, under the assumption that Nanog and GFP translation, mRNA and protein decay rates are the same (Fig. 2.14, red dotted line). This means that, theoretically, the empirical bimodal distribution of Nanog expression as reported by GFP in the heterozygous knock-in cell-line does not reflect the distribution of Nanog in the wild-type cell. Rather, the model predicts that the distribution of endogenous Nanog expression is homogeneous, and expression levels are a half order of magnitude greater than the reported high state.

## 2.4   Conclusions

In the first part of this chapter, our analysis of the extended ESC TRN revealed that the global feedback-rich structure is highly nested and critically dependent on Nanog, Oct4 and Sox2, all of which participate in two thirds of all feedback loops. The feedback centrality identified Nanog as the most central element in the global feedback structure, and the removal of Nanog leaves only a third of the feedback loops in tact. Taken together with other findings in [3], our results indicate that Nanog fluctuations regulate population heterogeneity by

transiently activating different subnetworks in the extended ESC TRN, driving transitions between a Nanog-expressing, feedback-rich, robust and self-perpetuating pluripotent state and a Nanog-diminished, feedback-sparse and differentiation-sensitive state.

We note that although the feedback structure of the extended TRN is severely compromised on removal of Nanog, it is not entirely destroyed: a small number of key feedback loops still remain, most notably those involving Oct4, Sox2, Dax1 and Rex1. MacArthur et al. [3] suggest that this may explain why, although they are prone to differentiate, ESCs can be maintained in a self-renewing state in the absence of Nanog [117]. In this case self-renewing ESCs may adapt to depend on a depleted feedback structure, which highlights the remarkable robustness of the pluripotency TRN [3].

The results of our analyses complement those of the experiments carried out by MacArthur et al. [3], which showed that pluripotency decayed over time following the removal of Nanog, and the efficiency of the rescue of Nanog levels progressively diminished as the core network disintegrated. They found that when Nanog was reintroduced after three days, the system had crossed a critical point and it was no longer possible to return to the pluripotent state. The authors proposed a simple RRE model of Nanog regulation of pluripotency, the analysis of which suggests that the observed irreversibility is due to the Nanog-dependent positive feedback loops in the TRN, which give rise to a one-way switch.

In the second part of this chapter, we presented mathematical models of positive feedback that explain how the experimentally observed bimodal distribution of Nanog protein expression could be a reporter artefact, as opposed to a significant biological phenomenon. We began by considering a CME model of the core ESC network for both the wild-type cell and the heterozygous knock-in cell-line. We demonstrated that while the network was able to give rise to a bimodal distribution of expression for Nanog and GFP in the model knock-in cell-line, these distributions exhibited homogeneity in the model wild-type cell under the same parameter values.

Although the model of the core ESC network agreed quantitatively with the experimental data, it was too complicated to enable the identification of the components that control stability. Therefore we simplified the model to the Nanog autoregulatory loop alone, and by observing from the time series simulations that the dynamics of the Nanog protein are decoupled from those of Oct4 and Sox2, we reasoned that the extrinsic noise in Nanog expression is due to fluctuations in Oct4 and Sox2 expression. We obtained the same results for

this reduced model as for the core ESC network model, in that the modelled GFP knock-in cell-line can achieve a bimodal distribution of both GFP and Nanog expression, while under the same parameter values, the Nanog distribution is unimodal in the wild-type system.

In order to determine the mechanism by which the creation of a Nanog null allele affects Nanog dynamics, we studied the corresponding RREs for the system. By reducing the dimensions of the systems of ODEs we found that the baseline and feedback production rates in the knock-in cell line are half that of those in the wild-type cell, and the stability analysis revealed this reduction weakens the strength of the auto-activatory feedback loop.

The dynamics of Nanog expression are governed by two parameters $\alpha$ and $\gamma$, whose values determine if the system is monostable or bistable. By dividing the $\alpha$-$\gamma$ plane into behavioural regimes, we were able to see the regions to which a wild-type cell in the Nanog high region would transition as a result of the insertion of the GFP knock-in; depending on the values of $\alpha$ and $\gamma$ in the wild-type cell, the introduction of the reporter construct will result in either a bistable, or a lower monostable, distribution of Nanog and GFP expression. Wild-type cells that reside anywhere in the bistable region are at risk of a shift into the monostable low region, or a significant change in the proportion of cells in each state. In summary, our analyses demonstrated that the introduction of the GFP knock-in reporter can induce a qualitative change in Nanog concentration that alters the behaviour and therefore the function of the cell.

In the last section of this chapter we fitted a model to the experimentally observed bimodal distribution of Nanog protein expression in 2i culture conditions, as reported by the GFP knock-in cell line, and this enabled us to infer the distribution of Nanog expression in the wild-type cell. We found that, in theory, the observed bimodal distribution does not reflect that of Nanog in the wild-type cell. Instead, the model predicts that the distribution of endogenous Nanog expression is homogeneous, and expression levels are a half order of magnitude greater than the reported high state. However, the inference was based on the assumption that the translation and decay rates of GFP mRNA and proteins are the same as those of Nanog. Since this is unlikely to be true, the observed GFP distribution may not even reflect the distribution of Nanog expression in the reporter cell-line, and the differences in reaction rates must be taken into account to obtain a more accurate estimate of the distribution of Nanog expression in the wild-type cell.

# Chapter 3

# Mathematical Modelling of Multipotency

## 3.1 Introduction

In this chapter, we consider the role of cell-to-cell variability within multipotent (as opposed pluripotent) stem cell populations. In regard to our findings in Chapter 2 on reporter strategies, we note that the studies discussed hereafter did not use a heterozygous knock-in reporter cell line, and therefore normal regulation of the gene of interest was not disturbed by the removal one of the alleles. Moreover, the experimentally observed distributions of expression to which we fit a mathematical model were obtained by immunostaining. This method is thought to have no direct effect on the transcription, translation or decay processes, because GFP proteins are attached to the target proteins (via an antibody) after these processes have occurred, and therefore the distributions of interest have already been established.

### 3.1.1 Functional Diversity in Cellular Populations

Clonal populations of unicellular organisms often exhibit phenotypic diversity – in which qualitatively different subpopulations of cells coexist – which confers selective advantage under adverse environmental conditions. Well-known examples include antibiotic bacterial persistence, the lysis-lysogeny switch of $\lambda$ phage, competence development and sporulation of *B. subtilis*, and lactose uptake by *E. coli* [178, 179]. The ubiquity of this phenomenon indicates that it is a generic, evolvable mechanism that facilitates collective cellular dynamics

by enabling robust, rapid responses to diverse environmental changes. Recently, stochastic fluctuations in the expression of important marker proteins have been seen to generate functional diversity within multipotent mammalian stem cell populations, suggesting a similar role for cell-to-cell variability in higher organisms [108, 180–182]. These observations have motivated speculation that functional multipotency (the ability to differentiate along a number of distinct cellular lineages) is a collective property of stem and progenitor cell populations, reflective of fitness constraints imposed at the population – rather than the individual cell – level [183–185]. This perspective is appealing since such regulated cell–to–cell variability, in principle, allows a cellular population to remain primed to respond quickly to a range of different differentiation cues while remaining robust to cell loss. However, convincing demonstrations of the potency of individual stem cells appear to argue strongly against such a collective view. For example, single long-term repopulating hematopoietic stem cells are able to fully reconstitute the blood system of lethally irradiated adult mice, and small numbers of pluripotent stem cells are able to rescue the development of genetically compromised embryos [186, 187]. Thus, it is still unclear how population-level and cell-intrinsic regulatory programs interact to control mammalian stem and progenitor cell dynamics. In this chapter, we use tools from statistical mechanics and information theory to propose a theoretical framework for the functional role of this variability. In order to illustrate our perspective, we model the concentration of a single protein in blood forming stem cells, *in vitro*, as a continuous Markov process, which we describe by a stochastic differential equation. Our theoretical framework is based on three important mathematical concepts that we now define.

### 3.1.2 Ergodicity and Entropy

A continuous time Markov chain $\{X(t) : t \geq 0\}$ with state space $S$ is ergodic if for any real valued function, $f : S \to \mathbb{R}$, the time average of the values of $f(X(t))$ converge to the spatial average over the entire space, with respect to the stationary probability density function, $p_\infty$. That is

$$\lim_{T \to \infty} \frac{1}{T} \int_0^T f\left(X(t)\right) \ dt = \int f(x) \ dp_\infty.$$

The notion of ergodicity is strongly related to the thermodynamic concept of entropy. In this thesis we will not be concerned with entropy in its thermodynamic context but rather as a general measure of "information". The first major application of entropy in fields outside of thermodynamics was introduced by Claude Shannon in his pioneering paper of 1948

[188], in which he simultaneously created the field of information theory and solved most of its fundamental problems. Information theory is the mathematical study of the accurate communication of information or data over a communication channel - the medium used to transmit the signal from transmitter to receiver, which could be a band of radio frequencies, or a beam of light, for example. In this context, "information" is thought of as a set of messages, where the goal is to send these messages over a noisy channel, and then to have the receiver reconstruct the message with low probability of error, in spite of the channel noise. Information theory is concerned with the theoretical limitations and potential performance of systems that can turn noisy channels into reliable communication channels using computational encoding and decoding methods.

Shannon began his paper [188] by defining a measure of how much information is produced by an ergodic discrete Markov process. Given a set of possible events whose probabilities of occurrence are $p_1, p_2, \ldots, p_n$, Shannon stated that the measure, $S$, of how much "choice" is involved in the selection of the event or of how uncertain we are of the outcome must satisfy the following three properties:

1. $S$ should be continuous with the $p_i$.

2. If all the $p_i$ are equal then $S$ should be a monotonic increasing function of the number of states. This is because there is more uncertainty, or more choice when there are more possible events.

3. Any succession of choices is the weighted sum of the individual values of $S$.

He deduced that the only $S$ that satisfies these conditions is of the form $S = -K \sum\limits_{i=1}^{n} p_i \log_b p_i$, where $K$ is a positive constant. In information theory $K = 1$, and the logarithm is usually to the base 2 (in which case, entropy is measured in bits), while in thermodynamics $K$ is the Boltzmann constant and the logarithm is natural. Informally, entropy is a measure of how flat a probability distribution is. For a probability distribution that is not subject to any constraints, its entropy is at its maximum value when all probabilities are equal, in which case it is equal to the log of the total number of states.

The second law of thermodynamics states that during any spontaneous process, the total entropy change for an isolated system is positive. This means that the entropy of an isolated system out of equilibrium increases over time toward the maximum entropy uniform

distribution (in which all states are equally likely). However, for systems that interact with their environment and are likely to be subject to constraints, the second law is not directly applicable. In this case, the equilibrium distribution is not expected to be uniform, and the entropy may not always increase. Rather, subject to certain reasonable assumptions, a related quantity known as the relative entropy or Kullback-Leibler divergence (with respect to the equilibrium distribution) decreases with time [189, 190]. More formally, if $p_\infty$ is the unique stationary distribution (i.e. the process is ergodic), $D(p \mid\mid p_\infty) = \int p \log \frac{p}{p_\infty} \, dx$ decreases with time.

## The Principle of Maximum Entropy

*The principle of maximum entropy* as a method for estimating a probability distribution was first proposed by Edwin Jaynes in 1957 [191]. It states that, given a set of known constraints on the target distribution, then, among all distributions satisfying these constraints, we should choose the one that is "maximally non-committal with regard to missing information" [191], i.e., the one with the largest Shannon entropy.

More precisely, consider a random variable $X$ that takes known values $x_1, \ldots, x_n$, with unknown probabilities $p_1, \ldots, p_n$, and with $m$ constraint functions $f_k(x)$ with $1 \leq k \leq m < n$, where

$$\langle f_k(X) \rangle = F_k,$$

and the $F_k$ are fixed. Then the maximum entropy principle assigns probabilities in such a way that maximises the information entropy, $S(p_1, \ldots, p_n)$, of $X$ under the above constraints, along with the constraint that the probabilities must sum to one, i.e.,

$$\sum_{i=1}^{n} p_i = 1.$$

The usual approach is to solve this optimisation problem using Lagrange Multipliers. Thus, we introduce $m + 1$ Lagrange multipliers $\lambda_k, \ k = 1, \ldots, m, \mu$ (one for each constraint), and the function

$$L(p_1, \ldots, p_n; \lambda_1, \ldots, \lambda_m, \mu) = -\sum_{i=1}^{n} p_i \log p_i - \sum_{k=1}^{m} \lambda_k \left( \sum_{i=1}^{n} f_k(x_i) p_i - F_k \right) - \mu \left( \sum_{i=1}^{n} p_i - 1 \right)$$

which we would like to maximise with respect to $p_1, p_2 \ldots, p_n; \lambda_1, \lambda_2, \ldots, \lambda_m; \mu$. Thus we solve

$$\nabla_{p_1,\ldots,p_n;\lambda_1,\ldots,\lambda_m,\mu} L = 0.$$

Differentiating with respect to $p_i$, we obtain

$$-\log p_i - 1 - \sum_{k=1}^{m} \lambda_k f_k(x_i) - (\lambda_0 - 1) = 0,$$

using the notation $\lambda_0 = \mu + 1$. Hence

$$p_i = \exp\left(-\lambda_0 - \sum_{k=1}^{m} \lambda_k f_k(x_i)\right). \tag{3.1}$$

The Lagrange multipliers, $(\lambda_0, \lambda_1, \ldots, \lambda_m)$, are then found from the relevant constraints. The constraint on the sum of probabilities gives

$$1 = \sum_{i=1}^{n} p_i = e^{-\lambda_0} Z,$$

where

$$Z = Z(\lambda_1, \ldots, \lambda_m) = \sum_{i=1}^{n} \exp\left(-\sum_{k=1}^{m} \lambda_k f_k(x_i)\right)$$

is the *partition function*, and therefore

$$e^{-\lambda_0} = \frac{1}{Z}, \quad \lambda_0 = \log Z. \tag{3.2}$$

The remaining constraints give

$$F_k = \sum_{i=1}^{n} f_k(x_i) \, p_i = e^{-\lambda_0} \sum_{i=1}^{n} f_k(x_i) \exp\left(-\sum_{r=1}^{m} \lambda_r f_r(x_i)\right) = -\frac{1}{Z} \frac{\partial Z}{\partial \lambda_k},$$
$$= -\frac{\partial \log Z}{\partial \lambda_k}, \tag{3.3}$$

which are $m$ simultaneous, implicit equations, sufficient to determine the $m$ unknowns $\lambda_k$, and are usually solved by numerical methods. Using Eq. (3.2) the probabilities (3.1) are fully determined:

$$p_i = \frac{1}{Z} \exp\left(-\sum_{k=1}^{m} \lambda_k f_k(x_i)\right).$$

and the maximum-entropy distribution has the generic form

$$p^*(\boldsymbol{x}) = \frac{1}{Z} \exp\left(\sum_k \lambda_k f_k(\boldsymbol{x})\right).$$

It can be proved [192, 193] that the maximum entropy distribution and the maximum likelihood distribution are equal when the family of distributions is the exponential family, i.e., when the family is defined by

$$p(\boldsymbol{x}|\boldsymbol{\lambda}) = \frac{1}{Z(\boldsymbol{\lambda})} \exp\left(\sum_k \lambda_k f_k(\boldsymbol{x})\right),$$

where $Z$ is the normalising constant, the functions $f_k(\boldsymbol{x})$ are given, and the parameters $\boldsymbol{\lambda} = \{\lambda_k\}$ are not known. By differentiating the log-likelihood, it can be shown [192, 193] that the maximum likelihood parameters $\{\boldsymbol{\lambda}_{ML}\}$ satisfy

$$\sum_{\boldsymbol{x}} p(\boldsymbol{x}|\boldsymbol{\lambda}_{ML}) f_k(\boldsymbol{x}) = \frac{1}{N} \sum_{n=1}^{N} f_k(\boldsymbol{x}^{(n)})$$

where the left-hand sum is over the entire state space $\boldsymbol{x}$, and the right-hand sum is over the set of $N$ data points, $\{\boldsymbol{x}^{(n)}\}$. However, the maximum-likelihood setting in classical statistics (see, for example, [194], Chapter 6) differs from the principle of maximum entropy setting. In maximum likelihood, the true distribution is assumed to be from the same family as the distributions over which the likelihood is maximised, whereas the principle of maximum entropy poses no parametric assumptions on the truth.

## 3.2  A Theoretical Framework for Multipotency

In application to dynamic variation in gene expression, ergodic theory says that if there exists a unique equilibrium distribution $p_\infty(x)$, toward which every initial condition (subpopulation of cells) converges, then the underlying stochastic processes that regulate expression are ergodic, and each cell in the population will independently explore the same state space in accordance with $p_\infty(x)$. In this case the population is intrinsically robust to targeted removal of any particular subpopulation because the remaining cells in the population will eventually "recolonize" (in state space) the removed subpopulation and reconstitute the equilibrium distribution $p_\infty(x)$. Similar reconstitution has been experimentally observed with respect to a number of proteins that are associated with stem and progenitor cells, including Nanog [108,

117], Rex1 [121], and Stella [122], as well as between distinct phenotypic states within cancer cell populations [195], and with respect to the stem cell surface marker Sca1 in hematopoietic progenitor cells [196]. These studies strongly suggest that the underlying stochastic processes that regulate expression variation are "ergodic-like".

With this in mind we propose a theoretical framework that views cellular multipotency as an instance of maximum entropy statistical inference. In this view, individual cells satisfy any minimal regulatory constraints imposed upon them (such as basic metabolic requirements, etc.), yet, in the absence of defined instructions, are maximally noncommittal with respect to their remaining molecular identity, thereby generating a diverse population that is able to respond optimally to a range of unforeseen future environmental changes. Thus, rather than viewing the multipotent cell state as an attractor of the underlying molecular regulatory dynamics (i.e., associating cellular identities with well-defined, stable patterns of gene expression – a common modelling assumption that has received some experimental validation for differentiated cell types [77]) individual multipotent cells are characterised by fundamental *uncertainty* in their molecular state, and their populations exhibit variability in accordance with this intrinsic uncertainty. However, since this model exchanges the attractor hypothesis at the single cell level for an ergodicity assumption for the underlying stochastic processes, each individual cell has the latent potential to assume every identity within the population, and it thereby retains the regenerative capacity of the entire population. As this view is fundamentally stochastic, its corollary is that regulation of multipotency occurs at the level of probabilities (i.e., at the population level), rather than at the individual cell level. In order to illustrate this perspective we will consider here the expression dynamics of Sca1, a key protein expressed in mouse hematopoietic progenitor cells, and known to have an important role in regulating cell fate [196–199].

### 3.2.1 Hematopoietic Stem Cells

Formation and maintenance of the blood system depends on hematopoietic stem cells (HSCs), which reside in small numbers in the bone marrow of adult mammals [200]. Since mature blood cells are predominantly short lived, HSCs are critically important for replenishing them and therefore sustaining the life of the individual [201]. In a human adult, $10^{11}$–$10^{12}$ blood cells are produced daily to maintain steady state levels [202]. As with all other stem cells, HSCs are capable of self-renewal, i.e., they possess the ability to divide into identical copies

of themselves without forming any newly differentiated features, enabling them to maintain their population size. HSCs can also differentiate to all blood cell lineages [203]. They sit at the top of a hierarchy of progenitors that become progressively restricted to several or single lineages (see Fig. 3.1) [204]. These progenitors yield blood precursors committed to differentiate down a single lineage and produce mature blood cells. While progenitor cells do not have the same self renewal capacity as the parent stem cells, they retain some capacity for further divisions.

HSCs differentiate in the bone marrow to lymphoid or myeloid stem cells [205]. Lymphoid stem cells give rise to B-cells and T-cells, which are the major cellular components of the adaptive immune system [206, 207]. Myeloid stem cells give rise to a second level of lineage-specific cells that go on to produce the following mature cells [208]:

- Granulocytes (basophils, neutrophils, and eosinophils) are types of white blood cell that help the body fight bacterial infections [209].

- Megakaryocytes and platelets help wounds heal and stop bleeding by forming blood clots [210].

- Macrophages reside in every tissue of the body and engulf dead cells and pathogens [211].

- Mast cells release histamine and are part of the immune system [212].

- Erythrocytes and reticulocytes are the oxygen carrying red blood cells [213].

- Dendritic cells act as messengers between the innate and the adaptive immune systems by processing antigen material and presenting it to the T-cells of the immune system [214].

To maintain hematopoietic homeostasis, HSC numbers need to be precisely regulated [216]. The cell fate decisions (life and death, and self renewal and differentiation) of HSCs are important processes that regulate the size and lifespan of the HSC pool. Defects in these processes can contribute to hematopoietic insufficiencies and the development of hematopoietic malignancies [216]. While human HSCs are frequently transplanted (obtained from the same or a different individual), to treat a range of diseases and cancers of the blood malignancies and congenital immunological defects [217], there is chronic shortage of suitably matched

Figure 3.1: The lineages of the hematopoietic system. Adapted from [215].

HSCs for transplantation, which has generated a great deal of interest in proliferating HSC colonies *in vitro* [218]. One major challenge in expanding HSC colonies is preventing them from differentiating and eventually dying [219]. Recent protocols use cultures with defined media (in which all of the chemical components are known) and specially engineered growth factors because of their simplicity and efficacy [220].

### 3.2.2 The Experimental Data

To investigate the origins of the cell-to-cell variability within clonal populations of HSCs, Chang et al. [196] analysed the dispersion of expression levels of stem cell surface marker Sca1 in populations of a multipotent erythroid, myeloid, and lymphocytic (EML) mouse hematopoietic progenitor cell line. Immunostaining and flow cytometry (see Section 2.1.2 for details) was used to generate experimental data that revealed a heterogeneous, bimodal distribution of Sca1 expression (see Fig. 3.3, bottom panel and [196]). Sca1 is functionally important since it has a role in HSC and progenitor lineage fate. Specifically, low Sca1 levels

predispose HSC and progenitor cells to myeloid and T-cell differentiation [196, 221, 222], whereas high Sca1 levels are required for erythroid, megakaryocyte, platelet and B-cell differentiation [196, 221–223]. In order to characterise the dynamics by which the population heterogeneity arises, Chang et al. [196] observed the evolution of the distributions of selected subpopulations. Cells from the lowest, highest and middle 15% of the distribution of Sca1 expression were isolated, and the evolution of the distributions of these three subpopulations was observed as they equilibrated in culture over a period of 18 days. All three initial conditions eventually reconstituted the parental distribution (see Fig. 3.3), and this convergence took approximately an order of magnitude longer than the cell cycle time.

With the aim of ascertaining the phenomenon that drives the regeneration of the parental distribution from the three sorted population fractions, [196] took a phenomenological approach to determine which general class of models of stochastic processes best describes the observed behaviour. The first model described Sca1 expression as a simple mean reverting (Ornstein-Uhlenbeck) process that includes both noise-driven diffusion (capturing the generation of cell-to-cell variability) and a drift towards the deterministic equilibrium (representing relaxation to the parental distribution mean), but this process described the data poorly, because it failed to recapitulate the regeneration of the low state (subpopulation) from the Sca1 high fraction.

Thus, the simple Ornstein-Uhlenbeck model was extended to include transitions between distinct states using a Gaussian mixture model (GMM) as a first approximation to a multimodal system. To corroborate the notion of multiple subpopulations with respect to Sca1 steady-state expression, it was shown that the observed histogram evolution can be better described by a two-component GMM than a single Gaussian distribution. The cells in every measured histogram (time point) were then partitioned into two subpopulations according to the fitted GMMs and the expression values of the individual cells, thus providing the time evolution of the relative abundance (weight) for each subpopulation. For the Sca1-mid and Sca1-high fractions, the weights exhibited a steep change after 96 hours, before eventually reaching a plateau, suggesting that the stochastic transitions between the subpopulations had a dominant role in the eventual relaxation to equilibrium.

A further experiment was carried out to determine whether clonal heterogeneity in Sca1 protein expression correlates with heterogeneity of the differentiation potential of these cells. On the application of a stimulus (a protein signaling molecule called erythropoietin), cells

with low Sca1 expression showed an increased propensity to differentiate to the erythroid lineage, and those with high Sca-1 expression showed an increased propensity to differentiate to the myeloid lineage.

Taken together, these results indicate that the universal reconstitution of the parental distribution is a result of noise-driven transitions between co-existing Sca1 high and low states, which transiently prime individual cells for erythroid and myeloid differentiation, respectively, and generate a characteristically bimodal Sca1 expression distribution within the population.

Since the two-component GMM does not explain the observed sigmoidal evolution of the relative weights of the two subpopulations, Chang et al. [196] proposed a non-linear ODE model for the size of each subpopulation. The model is based on cell-to-cell communication mediated by signalling molecules secreted by the cells, the abundance of which alters the rate of switching between states. Although the model captures the observed sigmoidal kinetics, different parameter estimates were required for each of the three initial fractions, thus rendering the model not useful for predicting weight dynamics for new initial conditions. We have chosen to model the observed Sca1 expression dynamics, not only to illustrate the theoretical framework proposed in Section 3.2, but also to demonstrate that cell-to-cell communication is not required to obtain sigmoidal weight relaxation dynamics. Furthermore, our simple SDE model can predict the evolution of the distribution of Sca1 expression for any initial condition, not just that of the weights.

### 3.2.3   A Stochastic Differential Equation Model

Since the underlying mechanisms by which the stochastic fluctuations of Sca1 levels are regulated are not known, we assume here that the intracellular dynamics of the Sca1 expression level $z(t)$ are described by a generic stochastic differential equation:

$$\frac{dz}{dt} = a(z) + \sqrt{2d(z)}\xi(t), \qquad (3.4)$$

where $\xi(t)$ is a standard one-dimensional white noise process, $a(z)$ describes the deterministic dynamics, and $d(z)$ accounts for fluctuations in Sca1 levels due to both intrinsic sources (i.e., noise in the molecular processes involved in Sca1 production or decay, such as transcription, translation, translocation, and degradation, etc.) and extrinsic sources (i.e., fluctuations in upstream regulators and uncontrolled environmental noise). Rather than model Sca1 levels

directly, it is convenient to introduce a reaction coordinate $x(z)$ such that the Fokker-Planck Equation (FPE) for the probability density $p(x, t)$ has the form

$$\frac{\partial p}{\partial t} = L(p), \quad L(p) = \frac{\partial}{\partial x}\left(\frac{\mathrm{d}\psi}{\mathrm{d}x}\, p\right) + \sigma \frac{\partial^2 p}{\partial x^2}, \tag{3.5}$$

with scalar potential $\psi(x)$ and diffusion coefficient $\sigma$. Such a transformation, which maps the original dynamics to those of a Brownian particle in a one-dimensional potential field, may be achieved by application of Itō's lemma, which reads:

$$\frac{dx}{dt} = a(z)\frac{dx}{dz} + d(z)\frac{d^2x}{dz^2} + \sqrt{2d(z)}\frac{dx}{dz}\xi(t).$$

The reaction coordinate $x(z)$ can be chosen such that the noise term in this equation is constant, say $\sqrt{2\sigma}$, which gives the transformation

$$x = \int \sqrt{\frac{\sigma}{d(z)}}dz. \tag{3.6}$$

Since the dynamics are one-dimensional we may also introduce a potential $\psi(x)$ such that

$$-\frac{d\psi}{dx} = a(z)\frac{dx}{dz} + d(z)\frac{d^2x}{dz^2},$$

to obtain

$$\frac{dx}{dt} = -\frac{d\psi}{dx} + \sqrt{2\sigma}\xi(t),$$

which is the stochastic differential equation corresponding to the FPE given by Eq. (3.5). Experimental data suggest that protein expression fluctuations often scale linearly with expression level [177]. Thus, a natural choice for the noise term is $d(z) = \sigma z^2$. Substituting this into Eq. (3.6) gives $x = \log(z)$. This approach is similar to that taken in [224].

The stationary solution of Eq. (3.5) is the Boltzmann-Gibbs distribution

$$p_\infty(x) = Z^{-1}\exp(-\psi/\sigma), \quad Z = \int \exp(-\psi/\sigma)\,\mathrm{d}x. \tag{3.7}$$

Figure 3.2: Smoothed probability density estimate of the empirical equilibrium Sca1 distribution, calculated from the aggregate of all three observations at $(t = 432)$ (left), and the corresponding potential function $\psi(x)$, (right). Experimental data, obtained from Chang et al. [196], are shown in dark red and the fitted model is shown in black.

This solution exists so long as $\psi(x)$ grows sufficiently rapidly as $|x| \to \infty$ that the partition function $Z$ remains finite. In this case, the dynamics are ergodic and the free energy is

$$
\begin{aligned}
F(p) &= \int \psi p \, \mathrm{d}x + \sigma \int p \log p \, \mathrm{d}x, \qquad (3.8) \\
&= E(p) - \sigma S(p),
\end{aligned}
$$

where $E(p)$ and $S(p)$ are the energy and entropy functionals, respectively. We use free energy later in Section 3.2.5 to assess the convergence of the probability density to the equilibrium state.

### 3.2.4 Model Fitting Based on a Data-Driven Potential

In the absence of detailed information on how Sca1 fluctuations are regulated, the potential $\psi(x)$ may be estimated numerically from the empirical Sca1 distribution by inverting Eq. (3.7). The model then has a single free parameter, the diffusion coefficient $\sigma$, which sets the time scale for the dynamics.

Estimates of $\sigma$ and $\psi(x)$ were obtained by model fitting using maximum likelihood estimation to the experimentally observed evolving Sca1 expression distributions starting from the three preselected populations. This was done numerically by fitting Eq. (3.5) using the partial differential equation solver 'pdepe', in Matlab. The equilibrium probability density and fitted potential is shown in Fig. 3.2.

Despite the simplicity of this model, excellent agreement with the experimental time-series data was observed from all three initial conditions, using the same numerically estimated potential and the same estimate of $\sigma$ (Figs. 3.3–3.5).

Figure 3.3: Model fit to experimental data. Model simulations using the same estimates of $\psi(x)$ and $\sigma$ are shown against the three independent experimental time-series; simulations differ only in the experimentally prescribed initial conditions. Data, obtained from Chang et al. [196], are in dark red and the fitted model is in black. The potential $\psi(x)$ was estimated numerically via Eq. (3.7) using aggregated data from the final time point.

### 3.2.5 Quantification of Convergence to Equilibrium

In their original publication, Chang et al. [196] argued, based upon analysis of changing proportions of cells in the Sca1 high and low states, that the observed dynamics are characterised by slow "sigmoidal" relaxation towards the stationary state. Since a constant probability flux across a barrier naturally leads to exponential relaxation, they suggested that these dynamics

Figure 3.4: Convergence to equilibrium with respect to the free energy. Exponential convergence was observed from all three initial conditions for large time, in accordance with Eq. (3.5).

indicate deviation from expected first-order kinetics, possibly due to regulation of Sca1 fluctuations by cell-to-cell communication or autocrine signalling, and proposed a simple model of cell-to-cell communication to explain this data. However, it is apparent that such recourse is not needed since in all three cases the experimental system is initially far from equilibrium, and therefore far from the regime in which first-order kinetics apply. Rather, in accordance with standard reaction-rate theory, the dynamics are characterised by an initial transient period during which local equilibrium is first established within each potential well, before transitions between wells occur [225].

The natural way to examine convergence to equilibrium for Eq. (3.5) is via the free energy, which is a Lyapunov functional for the dynamics [226, 227]. Examination of the free energy shows that this separation of time scales generates the observed convergence dynamics without the need to include additional regulatory mechanisms in the model (see Fig. 3.4). Taken together, these results indicate that the observed Sca1 expression dynamics are well described by a simple ergodic process in which individual cells behave independently with respect to Sca1 fluctuations.

### 3.2.6 First Passage Time Distributions

The ergodicity is useful since it allows inference of the behaviour of individual cells from the population dynamics. While stochastic excursions into the Sca1 high and low states have been seen to transiently confer different lineage biases to individual progenitor cells in culture, the time scales upon which these excursions occur at the single cell level are not known. Thus, the distribution of first passage times (FPTs) out of the Sca1 low and high states are of particular interest. Defining the ranges of Sca1 low and high expression as $L = (-\infty, x_0)$ and $H = (x_0, \infty)$, respectively, where $x_0$ is the intermediate maxima in $\psi(x)$, the FPT $T(x)$ out

Figure 3.5: First passage time (FPT) distributions in the Sca1 low (left panel) and high (right panel) states. The FPT distributions $F_X(x_X, t)$ starting at the local minima of the potential $\psi(x)$ are shown in black; the expected FPT distributions $\langle F_X \rangle(t)$ averaging over all initial conditions in $X \in \{L, H\}$ are shown in blue.

of $X$ for a cell initially at $x \in X$ (where $X \in \{L, H\}$) may be obtained from the backward FPE associated with Eq. (3.5). Denoting $G(x, t) = P(T(x) \geq t)$, we solve

$$\frac{\partial G}{\partial t} = -\frac{d\psi}{dx}\frac{\partial G}{\partial x} + \sigma\frac{\partial^2 G}{\partial x^2}, \tag{3.9}$$

with initial conditions $G(x, 0) = 1$ for $x \in X$ and boundary conditions $G(x_0, t) = \partial G/\partial x(\pm\infty, t) = 0$, from which the FPT distributions $F_X(x, t) = P(T(x) = t) = -\partial G/\partial t$, for $X \in \{L, H\}$, may be obtained. Conventionally, the FPT distribution $F_X(x, t)$ is evaluated from the local minima $x_X$ of $\psi(x)$ in $X$ since this is the state of highest probability. Alternatively, we can weight each initial position within $X$ according to the probability that the cell is at this position at equilibrium. We thus define the expected FPT distribution with respect to the Gibbs measure,

$$\langle F_X \rangle(t) = \int_{x \in X} \frac{p_\infty(x)}{w_X} F_X(x, t)\, dx,$$

where $w_X = \int_{x \in X} p_\infty(x) dx \in [0, 1]$ is the weight of the population in $X$. Numerical approximations to $F_X(x_X, t)$ and $\langle F_X \rangle(t)$ are shown in Fig. 3.5. These distributions yield mean FPTs of 60 (56) hours for the low state and 1573 (1487) hours for the high state using $F_X(x_X, t)$ ($\langle F_X \rangle(t)$). These time scales are substantially longer than the EML cell population doubling time (approximately 18–20 h [198]), and they therefore suggest that Sca1 fluctuations are not simply a consequence of the cell cycle. Rather, by setting the expected length of time that a pair of cells initially at the same position (e.g., daughter cells from the same cell division) will forget their common origin – and therefore the expected length of time that their identities will be coupled – Sca1 switching appears to encode an elementary form of epigenetic memory that endows individual cells with a transient functional identity. Since the rate of

switching is slower than the rate of cell division, this allows the formation of communities of cells that maintain the same characteristics though divisions, and are therefore able to adopt a temporarily stable functional phenotype. Yet, by allowing mixing between the communities on a feasible time scale, Sca1 fluctuations also safeguard long-term cell-to-cell variability and ensure that a robustly heterogeneous population, able to rapidly respond to a range of environmental challenges and resilient to the removal of cellular subtypes, is maintained.

### 3.2.7 Parameterised Model of the Potential

These results indicate that regulated fluctuations in Sca1 levels may be an intrinsic feature of EML cells in culture since they provide a mechanism by which the population hedges against unforeseen future environmental challenges and thereby retains the capacity to differentiate along both erythroid and/or myeloid lineages as required. If this is the case, then it is natural to ask if the experimentally observed stationary Sca1 distribution is optimal for this purpose; that is, if it is *maximally* variable in some appropriately defined way. To investigate this, it is convenient to introduce a parameterisation of the potential $\psi(x)$, in order to compare distributions. A parsimonious model, which allows for observed bimodality without introducing large numbers of parameters, is:

$$\frac{d\psi}{dx} = \beta x - \alpha_0 - \frac{\alpha_1 x^n}{K^n + x^n}, \tag{3.10}$$

where $n$ is a positive even integer. This restriction ensures that $\psi(x)$ is continuous and real for all $x \in \mathbb{R}$. Although in principle $x$ may be negative, Sca1 levels are sufficiently high that we did not observe negative values in practice. Intuitively, this is a simple model of a positive-feedback based bistable switch of the kind that commonly regulates cell fate changes [134, 135, 156]. We saw such a model for Nanog in ESCs in Chapter 2, where we demonstrated that positive feedback can give rise to a Hill function.

Eq. (3.10) may be integrated to obtain the following expression for $\psi(x)$ (see Appendix B.1 for the full derivation):

$$\psi(x) = \tfrac{1}{2}\beta x^2 - (\alpha_0 + \alpha_1)x + \frac{\alpha_1 K}{n} \sum_{k=1}^{n/2} \left[ 2S_k \arctan\left( \frac{\frac{x}{K} - C_k}{S_k} \right) - C_k \log\left( \left(\frac{x}{K}\right)^2 - 2C_k \frac{x}{K} + 1 \right) \right]$$

Figure 3.6: Fit of the parameterised model for the stationary distribution to the experimental data (left) and the corresponding potential function $\psi(x)$. In both panels, experimental data are shown in dark red and the fitted model is shown in black. The maximum likelihood estimates of the parameters are $\hat{\alpha}_0 = 0.118$, $\hat{\alpha}_1 = 0.0868$, $\hat{K} = 3.5308$, $\hat{\beta} = 0.0436$, $\hat{\sigma} = 0.0047$, and $\hat{n} = 12$.

where $S_k = \sin\left((2k-1)\pi/n\right), C_k = \cos\left((2k-1)\pi/n\right)$. Substituting this equation into Eq. (3.7) it can be seen that the stationary distribution $p_\infty(x)$ is characterised by five non-negative parameters: $n, \alpha_0/\sigma, \alpha_1/\sigma, \beta/\sigma, K$. Estimates of these parameters were calculated by model fitting to the aggregated data from the final time point using maximum likelihood estimation. The fitted stationary distribution is shown in Fig. 3.6 (left). An estimate of $\sigma$ was obtained by model fitting to the observed evolving Sca1 expression distributions, using maximum likelihood estimation. This was done numerically by fitting the Fokker-Planck equation, (3.5), using the partial differential equation solver 'pdepe', in Matlab. As expected, the estimate is similar to that for the numerical potential (0.004746 vs 0.004546). Since the plot of the model fit to the experimental time series data is indistinguishable from Fig. 3.3, it is not shown here. Again, all three experimentally observed time series are well described by the one dimensional FPE with the same parameter values, indicating that it is a suitable model of the dynamics of Sca1 expression.

**Stability Analysis**

To characterise the different aspects of the dynamics of the observed Sca1 system relative to the family of distributions given by Eq. (3.10), we begin by carrying out a stability analysis. Since the form of the SDE for Sca1 (as implied by Eq. (3.10)) is identical to that for Nanog

expression (2.19) in Chapter 2, except for the Hill coefficient, $n$ (where $n = 2$ for Nanog), the stability analysis here is similar to that for Nanog. Therefore, we avoid repetition of the previous stability analysis by only including the details that differ for general $n$.

Rescaling the system to aggregate the model parameters (see Appendix B.2 for the derivation), we obtain the nondimensional SDE

$$\frac{dx}{dt} = \alpha + \frac{x^n}{\gamma^n + x^n} - x + \sqrt{2\sigma_d}\,\xi(t), \tag{3.11}$$

where $\alpha = \dfrac{\alpha_0}{\alpha_1}$, $\gamma = \dfrac{K\beta}{\alpha_1}$, and $\sigma_d = \dfrac{\beta}{\alpha_1^2}\sigma$. Thus the stationary distribution $p_\infty(x)$ is characterised by four nonnegative dimensionless parameters: $\boldsymbol{\theta} = [n, \alpha, \gamma, \sigma_d]$.

To find the parameter values for which the system is bistable, we look for the fixed points, and therefore the solutions to the polynomial

$$x^{n+1} - (1 + \alpha)x^n + \gamma^n x - \alpha\gamma^n = 0. \tag{3.12}$$

In general there are $n + 1$ solutions. Using Descartes' rule of signs, when $n$ is positive and even, either one or three solutions are positive and real, and none are negative and real, i.e., the remaining solutions are complex.

We can ascertain the stability of the fixed points by plotting the curve $\dfrac{x^n}{\gamma^n + x^n}$ and the straight line $x - \alpha$ on the same plot, for varying values of $\gamma$, as shown in Fig. 3.7. We surmise the direction of flow from the sign of $F(x) = \alpha + \dfrac{x^n}{\gamma^n + x^n} - x$, as indicated by the arrows marked on the $x-$axes. Since $n$ reflects the steepness of the sigmoid, for values of $n \geq 1$ the line will intersect the curve in the same way as was described for Nanog in Section 2.3.2, as $\gamma$ varies from small to large. To find the values of $\alpha$ and $\gamma$ at which the bifurcations take place, we use the discriminant of Eq. (3.12). By expressing the discriminant of a polynomial in terms of its roots we can see that the discriminant is equal to zero if and only if there is a repeated real root

$$(-1)^{m(m-1)/2} a_m^{2m-2} \prod_{i \neq j}(r_i - r_j),$$

where $a_m$ is the leading coefficient and $r_1, \ldots, r_m$ are the roots of the polynomial. The discriminant $D$ of a polynomial, $p$ is given by the formula

$$D(p) = (-1)^{m(m-1)/2} \frac{1}{a_m} R(p, p').$$

Figure 3.7: (*Top row*) The Hill function $\frac{x^n}{\gamma^n + x^n}$ for $n = 12$ (solid line), and the function $x - \alpha$ (dashed line) for $\alpha = 0.2$ (dashed line), with $\gamma = 0.05$, $\gamma_-$, 0.45, 0.75, $\gamma_+$, 1.2, where $0.05 < \gamma_- < 0.45$ and $0.75 < \gamma_+ < 1.2$. The filled circles mark the fixed-points, and the arrows on the horizontal axes indicate the direction of flow between the fixed-points. When $\alpha = 0.2$, this system exhibits bistability for $\gamma_- < \gamma < \gamma_+$.

where $p'$ is the derivative of $p$, and $R(p, p')$, is the resultant of $p$ and $p'$, which is equal to the determinant of the Sylvester matrix, a $(2m - 1) \times (2m - 1)$ matrix, whose $m - 1$ first rows contain the coefficients of the polynomial and the $m$ last rows contain the coefficients of its derivative. The resultant, $R(p, p')$, of Eq. (3.12), is given in Appendix B.3. Substituting $m = n + 1$, we obtain

$$(-1)^{n(n+1)/2} R = 0$$

which we can now solve for $\gamma$ to obtain an expression for the critical values as function of $\alpha$. In this case, we found in data fitting that $n = 12$, and the positive, real solutions for $\gamma_\pm(\alpha)$ consist of a pair of complicated polynomials of degree 12, as given by Eq. (B.5) in Appendix B.3.

These two solutions divide the $\alpha\gamma$ parameter plane into three distinct behavioural regimes as illustrated in Fig. 3.8 (left). Since the conditions (2.26) for each bifurcation point, $(x_\pm, \alpha_\pm, \gamma_\pm)$, are met, they are saddle node bifurcations for $F(x) = \alpha + \frac{x^n}{\gamma^n + x^n} - x$.

Figure 3.8: The $\alpha\gamma$ parameter plane is divided into three regime regions as determined by the discriminant of the polynomial given by Eq. (3.12). Right: Hysteresis plot for $\alpha = 0.7$.

It is clear from Fig. 3.8 (left panel) that the Sca1 system exhibits hysteresis. The hysteresis plot, Fig. 3.8 (right panel), illustrates the following example. Consider a cell that occupies the monostable high region at the point $\alpha = 0.7$, $\gamma = 0.1$. As $\gamma$ increases the cell moves through the bistable region, and Sca1 expression remains high until it crosses the upper bifurcation curve and moves into the monostable low region. In order for the cell to return to a state of high expression, the value of $\gamma$ must now decrease such that the cell moves back through the bistable region and crosses both bifurcation curves. Thus the value of $\gamma$ at which the cell switches back to the original state of high expression is lower than the value at which it switches to the low expression state, resulting in a delay as it traverses the bistable region.

### 3.2.8 Maximum Entropy and Self-Organised Criticality

For fixed $\boldsymbol{\theta}$, the conditional probability $p_\infty(x \mid \boldsymbol{\theta})$ is the minimiser of the free energy $F(p)$, and may therefore be viewed as the most noncommittal way to assign probabilities subject to the particular constraints imposed upon the dynamics by $\psi(x; \boldsymbol{\theta})$ (i.e., an instance of maximum entropy statistical inference) [191]. As each set of model parameters defines a different potential, which places different constraints upon the dynamics, we may therefore determine the extent to which Sca1 fluctuations optimize population diversity by assessing the proximity of the empirical stationary Sca1 distribution to the maximum entropy distribution

$p_\infty^{\mathrm{max}}(x) = p_\infty(x \,|\, \boldsymbol{\theta}^*)$, where $S(p_\infty(x \,|\, \boldsymbol{\theta}^*)) = \max_{\boldsymbol{\theta}} S(p_\infty(x \,|\, \boldsymbol{\theta}))$. The relative entropy,

$$D(p_\infty \,||\, p_\infty^{\mathrm{max}}) = \int p_\infty \log \left( \frac{p_\infty}{p_\infty^{\mathrm{max}}} \right) \, dx,$$

is a natural way to measure this proximity. Since the Hill coefficient $n$ is, informally, a measure of the sensitivity of the underlying switch to the input stimulus, it primarily affects the curvature of the potential around the local minima $x_0$ (where present) and does not have a strong effect on the entropy. However, by governing a bifurcation that determines whether the underlying switch is in a monostable or bistable state, $\alpha$ and $\gamma$ can affect the entropy of the stationary distribution considerably. Fig. 3.9 shows how the relative entropy of $p_\infty(x)$ varies over the biologically relevant bistable region of the $\alpha\gamma$ plane. Note that $p_\infty(x)$ also depends upon $\sigma_d$, the relative strength of stochastic fluctuations. However, since this parameter includes the effects of unregulated extrinsic noise, we assume that it is not within the cell's capacity to regulate and fix it at the experimentally determined value. It can be seen that the point estimate for the experimentally observed Sca1 distribution is remarkably close to the maximum entropy distribution $p_\infty^{\mathrm{max}}(x)$. However, while the maximum entropy distribution is in the centre of the bistable regime, the empirical distribution is close to one of the critical lines that separate the bistable and monostable regimes (shown in blue in Fig. 3.9, right panel). *Self-organized criticality* (SOC) is the spontaneous evolution of a dynamical system to a critical state. The concept was put forward by Per Bak, Chao Tang and Kurt Wiesenfeld in 1987 [228], and is considered to be one of the mechanisms by which complexity arises in biological systems [229]. A dynamical system that is close to criticality will typically remain in one state under the vast majority of small, transient perturbations, but some perturbations can help generate large fluctuations that lead to a transition to a different state. Such systems are especially well suited for adaptation and information processing in the sense that adaptability is associated with the possibility of finding adequate new states in changing environments at very fast rates. Here, proximity to criticality specifically regulates the rate of mixing between the Sca1 high and low subpopulations, and therefore the response time of the population to environmental changes. To illustrate this, we consider how the mean first passage time (MFPTs) in each state varies in the vicinity of the maximum entropy state in the $\alpha\gamma$ plane. The MFPTs in the low and high states, denoted $\tau_-$ and $\tau_+$ respectively, can be found by solving the ordinary differential equation [190]

$$F(x)\frac{d\tau_\pm}{\partial x} + \sigma \frac{d\tau_\pm^2}{dx^2} = -1 \tag{3.13}$$

Figure 3.9: (Left panel) Entropy of the stationary distribution relative to the maximum entropy distribution over the $\alpha\gamma$ plane. The empirical distribution is marked with a magenta cross, and the maximum entropy distribution $p_\infty^{\max}(x)$ is marked with a green circle. Colour shows the percentiles. (Right panel) Minimum MFPT $\tau$ in the vicinity of the maximum entropy distribution (the dashed box in the left panel). The critical lines separating the bistable and monostable regimes are shown in blue. The empirical distribution lies in the small region of the $\alpha\gamma$ plane that is both close to critical and of high entropy. Colour shows dimensionless time.

with the same boundary conditions as for Eq. (3.9), i.e., $\tau_\pm(x_0) = d\tau_\pm/dx(\pm\infty) = 0$, where $x_0$ is the intermediate maxima in $\psi(x)$. Integrating once we obtain

$$\frac{d\tau_\pm}{dx} = \frac{-\frac{1}{\sigma}\int e^{-\Phi(x)/\sigma}dx}{e^{-\Phi(x)/\sigma}} \tag{3.14}$$

Applying the boundary conditions for the low state and integrating again we find:

$$\tau_- = \frac{1}{\sigma}\int_x^{x_0}\int_{-\infty}^u e^{\frac{\Phi(u)-\Phi(s)}{\sigma}}ds\,du \tag{3.15}$$

and similarly for the high state:

$$\tau_+ = \frac{1}{\sigma}\int_{x_0}^x\int_u^\infty e^{\frac{\Phi(u)-\Phi(s)}{\sigma}}ds\,du \tag{3.16}$$

The MFPT is conventionally evaluated from the local minima $x_X$ of $\psi(x)$ in $X$, since this is the state of highest probability. Fig. 3.9 (right panel) shows how $\tau = \min[\tau_-, \tau_+]$, varies in the vicinity of the maximum entropy state in the $\alpha\gamma$ plane. It can be seen that the minimum MFPT in the maximum entropy state is approximately an order of magnitude greater than that of the empirical distribution. Thus, while a population distributed according to the

maximum entropy distribution would ultimately be able to adapt better to environmental changes than the empirical population, it could not do so as rapidly. In this regard, close proximity to criticality is vital since it ensures that a diverse population is produced, yet mixing between subpopulations occurs on a physically relevant time scale. These results suggest that Sca1 levels are regulated by fitness constraints that involve a trade-off between maximising cell-to-cell variability and maintaining the ability to respond rapidly to environmental changes.

## 3.3   Conclusions

In summary, we have proposed an information-theoretic interpretation of stem cell dynamics that views cellular multipotency as an instance of maximum entropy statistical inference. We illustrated this view by analysing expression fluctuations of Sca1 using a simple stochastic model. The Langevin equation is a powerful tool for understanding the process of gene expression, since it allows us to obtain statistical information about the expression dynamics at the single cell level, and the resulting distribution at the population level. The choice of appropriate reaction coordinate provided a constant diffusion coefficient that allowed us to use a simple Langevin equation. Both the diffusion coefficient and the shape of a potential were inferred without knowing the underlying regulatory network. Moreover, the model accurately predicts the complex kinetics with which the sorted fractions relaxed back to the steady-state distribution of Sca1 expression.

These results indicate that the observed Sca1 expression dynamics are well described by a simple ergodic process in which individual cells behave independently with respect to Sca1 fluctuations. This ergodicity is useful since it allows inference of the behaviour of individual cells from the population dynamics. The predictive value of our modelling approach may have important practical implications for applications in stem cell biology. For example, the model implies that switching from a low to high state is a diffusive process, and therefore cells with concentrations near the state barrier are closer to switching than cells away from the barrier. In contrast, the two-state model described in Section 2.2.2, and given by Eq. (2.3), implies the rate of switching is uniform over the range of expression within each state. The two-state model was rejected in favour of the Langevin approach, since the gradual spreading of the shape of the distributions that are observed during relaxation suggests a process that

is driven by small step fluctuations (i.e., diffusion), and therefore an SDE model more accurately describes the dynamics of Sca1 expression [196]. The non-uniformity of the switching probabilities over the state space is reflected by the dependence of the FPT distributions, $F_X(x,t)$, predicted from our diffusive model, on the initial position, $x$. These FPT distributions quantify how the cells explore the state space differently depending on where they reside in the potential. A practical implication of this finding is that, for cell sorting experiments, the details of the window size and location used to select subpopulations can have a measurable influence on the apparent stability of the selected subpopulations.

In order to assess the proximity of the empirical stationary Sca1 distribution to the maximum entropy distribution, we had to introduce extra model parameters. However, the simple model of a positive-feedback based bistable switch described all three experimentally observed time series well, with the same parameter values, indicating that it was a suitable model of the dynamics of Sca1 expression. This model allowed us to see that the point estimate for the experimentally observed Sca1 distribution is both remarkably close to the maximum entropy distribution, and to one of the critical lines that separate the bistable and monostable regimes. We showed that the minimum MFPT in the maximum entropy state is approximately an order of magnitude greater than that of the empirical distribution, thus demonstrating how proximity to criticality specifically regulates the rate of mixing between the Sca1 high and low subpopulations. Therefore, our Langevin equation model has enabled us to use the tools of statistical mechanics to provide evidence to suggest that Sca1 levels are regulated by fitness constraints that involve a trade-off between maximising cell-to-cell variability and maintaining the ability to respond rapidly to environmental changes.

Since the biological mechanisms of gene transcription are extraordinarily complex and not completely understood, our model relies on a number of simplifying assumptions, both biological and physical in nature. One of the most explicit is the assumption that the reaction rates are constant. In reality, the rates are affected by many other processes including chromatin remodelling, translational regulation, and protein folding [142, 230]. In addition, the parameterised SDE describes the auto-activating influence at a high-level of abstraction, where the rate of Sca1 expression depends exclusively on itself. This means that the influence of other molecules and cellular processes is not taken into account directly. Although our models are phenomenological, these trade-offs are necessary since we do not know the structure of the

underlying regulatory network. Even with these limitations, SDE models of this kind can be useful in deciphering basic aspects of gene expression dynamics.

It has been shown [231–233] that culture conditions can have a significant effect on the width and shape of a distribution of gene expression in clonal populations of stem cells (see also Fig. 2.2 in Section 2.1.2). For example, a recent study by [233] demonstrated that the expression distribution of key regulators suggested a more-homogenous transcriptional network in mouse ESCs cultured with 2i than with two commonly used serum-based cultures – a result that is consistent with what we see in Fig. 2.2. Rather than viewing such disparity in the resulting distribution of expression as a problem that must be overcome, [233] conclude their findings show that, with proper targeting of the molecules that regulate the gene of interest, variability among ESCs is largely controllable without hampering pluripotency and self-renewal. Thus, since we are yet to see how the complex and evolving conditions *in vivo* affect the distribution of Sca1 expression, we cannot generalise our conclusions from the analyses in this chapter to the *in vivo* situation.

Although we have focused on Sca1 dynamics, comparable expression fluctuations are known to generate functional diversity in other mammalian stem cell systems [108, 117, 120–124, 234]. Thus, the generation of ergodic expression fluctuations may be a generic way in which cell populations maintain robust multilineage differentiation potential under environmental uncertainty. If so, then molecular noise processing could be particularly important in regulating stem cell function in a range of contexts. A better understanding of the relationship between molecular noise and stem cell identity should help to distinguish variability due to interchangeable subpopulations of cells from that due to the presence of distinct, noninterconvertible, cell types (i.e., to determine which underlying stochastic processes are ergodic) [235, 236]. We anticipate that advances in single cell profiling techniques will help to address these issues in the near future.

# Chapter 4

# Discussion

In this thesis, we have proposed and applied mathematical models of the regulation of gene expression to investigate the source of the experimentally observed cell-to-cell variability in both adult and embryonic stem cell populations.

In Chapter 2, we investigated the possibility that Nanog fluctuations regulate population variability by controlling feedback mechanisms in the extended ESC TRN, first by exploring how feedback in network structures relates to dynamics, and then by investigating the role of Nanog in the global feedback structure in the extended ESC TRN. To do this, we examined the effect of removing the Nanog feedback elements, by enumerating the feedback loops that each transcription factor participates in, for both the wild-type ESC TRN and the NanogR TRN. We found that this network is rich in feedback, and that the global feedback structure of this network is critically dependent on Nanog, Oct4 and Sox2, each of which participate in over two thirds of all feedback loops in the network. Thus, removal of Nanog destroys many of the feedback structures of the ESC TRN, leaving only a third of the feedback loops intact, and therefore severely compromises the global feedback structure. Furthermore, the feedback centrality, an adjusted measure of node involvement in the feedback structure in the network, identified Nanog as the most central element in the global feedback structure. These analyses indicate that Nanog fluctuations regulate population heterogeneity by transiently activating different subnetworks in the extended ESC TRN, driving transitions between a Nanog-expressing, feedback-rich, robust and self-perpetuating pluripotent state and a Nanog-diminished, feedback-sparse and differentiation-sensitive state.

We note that the single-gene perturbation we have studied here does not reflect the full complexity of transcriptional regulation in ESCs, and it remains to be determined whether feedback-controlled population heterogeneity has a role *in vivo*. A better understanding of the role of feedback in controlling ESCs will facilitate the maintenance of more defined pluripotent populations and the development of more robust differentiation protocols.

In the second study in Chapter 2, we presented mathematical models of positive feedback that explain why heterozygous knock-in reporters might not give a faithful reflection of the endogenous expression of autoregulating genes such as Nanog. We applied a combination of CME, CLE and RREs models to demonstrate that the creation of a Nanog-null allele disturbs normal Nanog transcriptional control by interfering with the Nanog feedback mechanism; our analyses indicated that the dynamics of *Nanog* are perturbed by the heterozygous knock-in strategy because the GFP proteins that are produced instead of Nanog cannot feedback on the *Nanog* allele, or that of *GFP*, thus ultimately reducing the rates of production of Nanog.

In our simulations we used the same parameter values for GFP translation and mRNA and protein decay as for those for Nanog in order to illustrate the effect of the knock-in reporter strategy alone. However, the rates of GFP reactions are almost certainly not identical, which would result in further differences between the true distribution of Nanog protein expression and that reported by GFP. Once the reaction rates for both genes have been experimentally determined, we could potentially then use our model to infer the true distribution of Nanog protein expression from the observed GFP distribution. Since GFP is currently seen as a direct proxy for Nanog, we could use our model to more accurately assess Nanog from GFP.

The heterozygous knock-in reporter method has commonly been used to measure gene expression, including that of Nanog, in a large number of experiments published in prominent journals [118, 144, 149]. Since the conclusions drawn are based on the assumption that the shape of the reporter distribution represents that of the true distribution, and that there is a linear relationship between the reporter and the protein of interest, and further studies are in turn based on these conclusions, our work demonstrates why it is important to use mathematical models to provide evidence that a reporter gives a faithful reflection of the endogenous expression of the target gene, or at least a method for inferring it.

Although many Nanog reporter constructs involve single-allele, knock-in strategies, there is an increasing number of more sophisticated methods aimed at reducing perturbation to the wild-type system. Illustrations of a range of Nanog reporter constructs are given in Fig. 4.1.

The first reporter construct, NHET, is the heterozygous knock-in reporter strategy analysed in Chapter 2, and TNGA is another single-allele construct used by [117]. The improved constructs NGR and NGNC used by [144, 149], are examples of double fluorescence transcription reporters, which have one or more copies of the coding region for a fluorescent protein inserted after that for Nanog (the coding region for Nanog remains). In these constructs, both Nanog alleles undergo this treatment – one allele reports GFP and the other an alternative fluorescent protein, such as mCherry. The coding region for the fluorescent protein(s) is separated from that for Nanog by a short sequence that codes for a peptide that causes the two proteins cleave after translation - and therefore the function of the Nanog protein is in principle not disturbed. The construct NVNK – a type of fluorophore-fused reporter – used by [237], is similar to NGR and NGNC, except the fluorescent protein remains fused to the Nanog protein, and therefore may interfere with its normal function. Nd [238] is a type of bacterial artificial chromosome (BAC) reporter that leaves the endogenous Nanog alleles intact, and carries the reporter in an additional DNA sequence with the same promoter region as that of Nanog.

Every reporter strategy has its own technological limitations, but further issues may occur depending on the underlying network of the protein on which it is reporting. We are hopeful that we could apply our model to these improved reporter strategies to determine how the dynamics of Nanog and other proteins could be perturbed by their implementation. For example, we could account for the increased coding length of the Nanog alleles in the double fluorescence transcription reporters by reducing the baseline and feedback production rates according to the number of inserted copies of the fluorescent protein. Since the fluorophore-fused reporter might alter the function of the Nanog protein, we could model the impaired function by reducing the effective concentration for feedback production.

In summary, researchers should be aware of the limitations of live-cell reporters, and understand how the system under observation can impact the faithfulness of the reporter.

In Chapter 3, we proposed an information-theoretic interpretation of stem cell dynamics that views cellular multipotency as an instance of maximum entropy statistical inference. Motivated by the goal of understanding of the functional role of the considerable cell-cell variability commonly exhibited by adult stem cells, we used a simple stochastic model to analyse the dynamics of the Sca1 protein in murine hematopoietic stem cells *in vitro*.

Figure 4.1: Illustrations of the Nanog gene in wild-type cells, and six types of Nanog reporter construct, used by [117, 118, 144, 149, 237, 238], respectively.

In future work, we would like to determine the effect of experimental measurement error on the estimates of the model parameters and the potential. The flow cytometer calibration histograms (Fig. 1b [196]) show the flow cytometer records a range of fluorescence intensities for any known numbers of fluorescent proteins. This is due to machine measurement error,

which can stem from uncontrollable factors such as acidity, temperature, and variation in the intensity of the beam of light directed at the cells. The spread of the signal obtained from the calibration histograms can provide an approximation of the error in the measurement of fluorescence intensity. Initial investigations suggest that taking account of the measurement error will make the model fitting very difficult. This is because it creates an inverse problem, whose numerical solution involves approximations that ultimately introduce more error in the parameter estimates than the procedure is designed to eliminate. It is therefore understandable that researchers chose to exclude this technical variability in their modelling efforts, although Chang et al. [196] did so because the heterogeneity in the distribution of FI was significantly larger than the measurement error [196].

Our results showed that observed dynamics naturally self-organise close to a critical state with near-optimal information-processing capacity. Although we have focused on Sca1 dynamics, comparable expression fluctuations are known to generate functional diversity in other mammalian stem cell systems [108, 117, 120–124, 234]. Thus, the generation of ergodic expression fluctuations may be a generic way in which cell populations maintain robust multilineage differentiation potential under environmental uncertainty. If so, then molecular noise processing could be particularly important in regulating stem cell function in a range of contexts. A better understanding of the relationship between molecular noise and stem cell identity should help to distinguish variability due to interchangeable subpopulations of cells from that due to the presence of distinct, noninterconvertible, cell types (i.e., to determine which underlying stochastic processes are ergodic) [235, 236]. We anticipate that advances in single cell profiling techniques will help to address these issues in the near future.

# Appendix A

# Details of the Models of the Core ES TRN

## A.1 The CME Model of the Core ES TRN

For Scenario 1 (2), the CME model of the core ES Network consists of 32 (34) molecular species and 43 (46) chemical reactions. The molecular species and the corresponding labels for their copy number are given in Table A.1. The state vector, $\boldsymbol{x}$, whose $i$th element is the copy number of molecular species $i$, is:

$$
\begin{aligned}
\boldsymbol{x} = (&M_N, M_S, M_O, P_N, P_S, P_O, C_N N, C_O S, D_{NA}^{NN'}, D_{NA}^{OS'}, D_{NB}^{NN'}, D_{NB}^{OS'}, D_O^{NN'}, D_O^{OS'}, \\
&D_S^{NN'}, D_S^{OS'}, D_{NA}^{NN}, D_{NA}^{OS}, D_{NB}^{NN}, D_{NB}^{OS}, D_O^{NN}, D_O^{OS}, D_S^{NN}, D_S^{OS}, {\color{red}M_G}, {\color{red}P_G})'
\end{aligned}
$$

where the elements in black are those for Scenario 1, and the additional two species for Scenario 2 are in red.

Similarly, the chemical reactions in Scenario 1 are in black and cyan (reaction 2 does not occur in Scenario 2), and the additional reactions in Scenario 2 are in red. The chemical reactions are described in Table A.2. The corresponding propensity functions, $a_j(\boldsymbol{x})$, and the non-zero values of the stoichiometric vectors, $\boldsymbol{\nu}_j$, are given in Table A.3.

Table A.1: The molecular species in the CME model of the core ES network, and the notation for their corresponding copy numbers. Those in Scenario 1 are in black and the additional species in Scenario 2 are in red.

| Species | Copy number notation |
|---|---|
| Nanog mRNA | $M_N$ |
| Sox2 mRNA | $M_S$ |
| Oct4 mRNA | $M_O$ |
| Nanog proteins | $P_N$ |
| Sox2 proteins | $P_S$ |
| Oct4 proteins | $P_O$ |
| Nanog-Nanog dimers | $C_{NN}$ |
| Oct4-Sox2 dimers | $C_{OS}$ |
| unbound NN dimer binding sites on N allele A (0 or 1) | $D_{NA}^{NN'}$ |
| unbound OS dimer binding sites on N allele A (0 or 1) | $D_{NA}^{OS'}$ |
| unbound NN dimer binding sites on N allele B (0 or 1) | $D_{NB}^{NN'}$ |
| unbound OS dimer binding sites on N allele B (0 or 1) | $D_{NB}^{OS'}$ |
| unbound NN dimer binding sites on S allele A (0 or 1) | $D_{SA}^{NN'}$ |
| unbound OS dimer binding sites on S allele A (0 or 1) | $D_{SA}^{OS'}$ |
| unbound NN dimer binding sites on S allele B (0 or 1) | $D_{SB}^{NN'}$ |
| unbound OS dimer binding sites on S allele B (0 or 1) | $D_{SB}^{OS'}$ |
| unbound NN dimer binding sites on O allele A (0 or 1) | $D_{OA}^{NN'}$ |
| unbound OS dimer binding sites on O allele A (0 or 1) | $D_{OA}^{OS'}$ |
| unbound NN dimer binding sites on O allele B (0 or 1) | $D_{OB}^{NN'}$ |
| unbound OS dimer binding sites on O allele B (0 or 1) | $D_{OB}^{OS'}$ |
| bound NN dimer binding sites on N allele A (0 or 1) | $D_{NA}^{NN}$ |
| bound OS dimer binding sites on N allele A (0 or 1) | $D_{NA}^{OS}$ |
| bound NN dimer binding sites on N allele B (0 or 1) | $D_{NB}^{NN}$ |
| bound OS dimer binding sites on N allele B (0 or 1) | $D_{NB}^{OS}$ |
| bound NN dimer binding sites on S allele A (0 or 1) | $D_{SA}^{NN}$ |
| bound OS dimer binding sites on S allele A (0 or 1) | $D_{SA}^{OS}$ |
| bound NN dimer binding sites on S allele B (0 or 1) | $D_{SB}^{NN}$ |
| bound OS dimer binding sites on S allele B (0 or 1) | $D_{SB}^{OS}$ |
| bound NN dimer binding sites on O allele A (0 or 1) | $D_{OA}^{NN}$ |

Table A.1: The molecular species in the CME model of the core ES network, and the notation for their corresponding copy numbers. Those in Scenario 1 are in black and the additional species in Scenario 2 are in red.

| Species | Copy number notation |
|---|---|
| bound OS dimer binding sites on O allele A (0 or 1) | $D_{OA}^{OS}$ |
| bound NN dimer binding sites on O allele B (0 or 1) | $D_{OB}^{NN}$ |
| bound OS dimer binding sites on O allele B (0 or 1) | $D_{OB}^{OS}$ |
| GFP mRNA | $M_G$ |
| GFP proteins | $P_G$ |

Table A.2: The chemical reactions of the CME model of the core ES network. Those in Scenario 1 are in black and cyan, and those in Scenario 2 are in black and red.

| Reaction | Description |
|---|---|
| 1 | Nanog mRNA transcription from Nanog allele A |
| 2 | Nanog mRNA transcription from Nanog allele B |
| 3 | Sox2 mRNA transcription from Sox2 allele A |
| 4 | Sox2 mRNA transcription from Sox2 allele B |
| 5 | Oct4 mRNA transcription from Oct4 allele A |
| 6 | Oct4 mRNA transcription from Oct4 allele B |
| 7 | Nanog mRNA decay |
| 8 | Sox2 mRNA decay |
| 9 | Oct4 mRNA decay |
| 10 | Nanog protein translation |
| 11 | Sox2 protein translation |
| 12 | Oct4 protein translation |
| 13 | Nanog protein decay |
| 14 | Sox2 protein decay |
| 15 | Oct4 protein decay |
| 16 | Nanog proteins dimerise |
| 17 | Sox2 and Oct4 proteins dimerise |
| 18 | Nanog dimer disassociates |
| 19 | Sox2-Oct4 dimer disassociates |
| 20 | Nanog-Nanog dimer binds to unbound binding site on Nanog allele A |
| 21 | Oct4-Sox2 dimer binds to unbound binding site on Nanog allele A |

Table A.2: The chemical reactions of the CME model of the core ES network. Those in Scenario 1 are in black and cyan, and those in Scenario 2 are in black and red.

| Reaction | Description |
|----------|-------------|
| 22 | Nanog-Nanog dimer binds to unbound binding site on Nanog allele B |
| 23 | Oct4-Sox2 dimer binds to unbound binding site on Nanog allele B |
| 24 | Nanog-Nanog dimer binds to unbound binding site on Sox2 allele A |
| 25 | Oct4-Sox2 dimer binds to unbound binding site on Sox2 allele A |
| 26 | Nanog-Nanog dimer binds to unbound binding site on Sox2 allele B |
| 27 | Oct4-Sox2 dimer binds to unbound binding site on Sox2 allele B |
| 28 | Nanog-Nanog dimer binds to unbound binding site on Oct4 allele A |
| 29 | Oct4-Sox2 dimer binds to unbound binding site on Oct4 allele A |
| 30 | Nanog-Nanog dimer binds to unbound binding site on Oct4 allele B |
| 31 | Oct4-Sox2 dimer binds to unbound binding site on Oct4 allele B |
| 32 | Nanog-Nanog dimer unbinds from binding site on Nanog allele A |
| 33 | Oct4-Sox2 dimer unbinds from binding site on Nanog allele A |
| 34 | Nanog-Nanog dimer unbinds from binding site on Nanog allele B |
| 35 | Oct4-Sox2 dimer unbinds from binding site on Nanog allele B |
| 36 | Nanog-Nanog dimer unbinds from binding site on Sox2 allele A |
| 37 | Oct4-Sox2 dimer unbinds from binding site on Sox2 allele A |
| 38 | Nanog-Nanog dimer unbinds from binding site on Sox2 allele B |
| 39 | Oct4-Sox2 dimer unbinds from binding site on Sox2 allele B |
| 40 | Nanog-Nanog dimer unbinds from binding site on Oct4 allele A |
| 41 | Oct4-Sox2 dimer unbinds from binding site on Oct4 allele A |
| 42 | Nanog-Nanog dimer unbinds from binding site on Oct4 allele B |
| 43 | Oct4-Sox2 dimer unbinds from binding site on Oct4 allele B |
| 44 | GFP mRNA transcription from Nanog allele B |
| 45 | GFP mRNA decay |
| 46 | GFP protein translation |
| 47 | GFP protein decay |

Table A.3: The propensity functions and the non-zero entries of the 43 (46) Stoichiometric vectors of the CME model of the core ES network for Scenario 1 (2). Those that apply to Scenario 1 (2) are in black and cyan (black and red).

| Reaction $j$ | Propensity Function $a_j(x)$ | stoichiometric vector $\nu_j$ |
|---|---|---|
| 1 | $g_M^N + k_2^{N/NN} D_{NA}^{NN} + k_2^{N/OS} D_{NA}^{OS}$ | $\nu_{1,1} = 1$ |
| 2 | $g_M^N + k_2^{N/NN} D_{NB}^{NN} + k_2^{N/OS} D_{NB}^{OS}$ | $\nu_{2,1} = 1$ |
| 3 | $g_M^S + k_2^{S/NN} D_{SA}^{NN} + k_2^{S/OS} D_{SA}^{OS}$ | $\nu_{3,2} = 1$ |
| 4 | $g_M^S + k_2^{S/NN} D_{SB}^{NN} + k_2^{S/OS} D_{SB}^{OS}$ | $\nu_{4,2} = 1$ |
| 5 | $g_M^O + k_2^{O/NN} D_{OA}^{NN} + k_2^{O/OS} D_{OA}^{OS}$ | $\nu_{5,3} = 1$ |
| 6 | $g_M^O + k_2^{O/NN} D_{OB}^{NN} + k_2^{O/OS} D_{OB}^{OS}$ | $\nu_{6,3} = 1$ |
| 7 | $d_M^N M_N$ | $\nu_{7,1} = -1$ |
| 8 | $d_M^S M_S$ | $\nu_{8,2} = -1$ |
| 9 | $d_M^O M_O$ | $\nu_{9,3} = -1$ |
| 10 | $g_P^N M_N$ | $\nu_{10,14} = 1$ |
| 11 | $g_P^S M_S$ | $\nu_{11,5} = 1$ |
| 12 | $g_P^O M_O$ | $\nu_{12,6} = 1$ |
| 13 | $d_P^N P_N$ | $\nu_{13,4} = -1$ |
| 14 | $d_P^S P_S$ | $\nu_{14,5} = -1$ |
| 15 | $d_P^O P_O$ | $\nu_{15,6} = -1$ |
| 16 | $\frac{a^{NN}}{2} P_N(P_N - 1)$ | $\nu_{16,4} = -2, \nu_{16,7} = 1$ |
| 17 | $a^{OS} P_S P_O$ | $\nu_{17,5} = -1, \nu_{17,6} = -1, \nu_{17,8} = 1$ |
| 18 | $u^{NN} C_{NN}$ | $\nu_{18,4} = 2, \nu_{18,7} = -1$ |
| 19 | $u^{OS} C_{OS}$ | $\nu_{19,5} = 1, \nu_{19,6} = 1, \nu_{19,8} = -1$ |
| 20 | $k_1^{N/NN} C_{NN} D_{NA}^{NN'}$ | $\nu_{20,7} = -1, \nu_{20,9} = -1, \nu_{20,21} = 1$ |
| 21 | $k_1^{N/OS} C_{OS} D_{NA}^{OS'}$ | $\nu_{21,8} = -1, \nu_{21,10} = -1, \nu_{21,22} = 1$ |
| 22 | $k_1^{N/NN} C_{NN} D_{NB}^{NN'}$ | $\nu_{22,7} = -1, \nu_{22,11} = -1, \nu_{22,23} = 1$ |
| 23 | $k_1^{N/OS} C_{OS} D_{NB}^{OS'}$ | $\nu_{23,8} = -1, \nu_{23,12} = -1, \nu_{23,24} = 1$ |
| 24 | $k_1^{S/NN} C_{NN} D_{SA}^{NN'}$ | $\nu_{24,7} = -1, \nu_{24,13} = -1, \nu_{24,25} = 1$ |
| 25 | $k_1^{S/OS} C_{OS} D_{SA}^{OS'}$ | $\nu_{25,8} = -1, \nu_{25,14} = -1, \nu_{25,26} = 1$ |
| 26 | $k_1^{S/NN} C_{NN} D_{SB}^{NN'}$ | $\nu_{26,7} = -1, \nu_{26,15} = -1, \nu_{26,27} = 1$ |
| 27 | $k_1^{S/OS} C_{OS} D_{SB}^{OS'}$ | $\nu_{27,8} = -1, \nu_{27,16} = -1, \nu_{27,28} = 1$ |
| 28 | $k_1^{O/NN} C_{NN} D_{OA}^{NN'}$ | $\nu_{28,7} = -1, \nu_{28,17} = -1, \nu_{28,29} = 1$ |
| 29 | $k_1^{O/OS} C_{OS} D_{OA}^{OS'}$ | $\nu_{29,8} = -1, \nu_{29,18} = -1, \nu_{29,30} = 1$ |

Table A.3: The propensity functions and the non-zero entries of the 43 (46) Stoichiometric vectors of the CME model of the core ES network for Scenario 1 (2). Those that apply to Scenario 1 (2) are in black and cyan (black and red).

| Reaction $j$ | Propensity Function $a_j(x)$ | stoichiometric vector $\nu_j$ |
|---|---|---|
| 30 | $k_1^{O/NN}C_{NN}D_{OB}^{NN'}$ | $\nu_{30,7} = -1,\ \nu_{30,19} = -1,\ \nu_{30,31} = 1$ |
| 31 | $k_1^{O/OS}C_{OS}D_{OB}^{OS'}$ | $\nu_{31,8} = -1,\ \nu_{31,20} = -1,\ \nu_{31,32} = 1$ |
| 32 | $k_3^{N/NN}D_{NA}^{NN}$ | $\nu_{32,7} = 1,\ \nu_{32,9} = 1,\ \nu_{32,21} = -1$ |
| 33 | $k_3^{N/OS}D_{NA}^{OS}$ | $\nu_{33,8} = 1,\ \nu_{33,10} = 1,\ \nu_{33,22} = -1$ |
| 34 | $k_3^{N/NN}D_{NB}^{NN}$ | $\nu_{34,7} = 1,\ \nu_{34,11} = 1,\ \nu_{34,23} = -1$ |
| 35 | $k_3^{N/OS}D_{NB}^{OS}$ | $\nu_{35,8} = 1,\ \nu_{35,12} = 1,\ \nu_{35,24} = -1$ |
| 36 | $k_3^{S/NN}D_{SA}^{NN}$ | $\nu_{36,7} = 1,\ \nu_{36,13} = 1,\ \nu_{36,25} = -1$ |
| 37 | $k_3^{S/OS}D_{SA}^{OS}$ | $\nu_{37,8} = 1,\ \nu_{37,14} = 1,\ \nu_{37,26} = -1$ |
| 38 | $k_3^{S/NN}D_{SB}^{NN}$ | $\nu_{38,7} = 1,\ \nu_{38,15} = 1,\ \nu_{38,27} = -1$ |
| 39 | $k_3^{S/OS}D_{SB}^{OS}$ | $\nu_{39,8} = 1,\ \nu_{39,16} = 1,\ \nu_{39,28} = -1$ |
| 40 | $k_3^{O/NN}D_{OA}^{NN}$ | $\nu_{40,7} = 1,\ \nu_{40,17} = 1,\ \nu_{40,29} = -1$ |
| 41 | $k_3^{O/OS}D_{OA}^{OS}$ | $\nu_{41,8} = 1,\ \nu_{41,18} = 1,\ \nu_{41,30} = -1$ |
| 42 | $k_3^{O/NN}D_{OB}^{NN}$ | $\nu_{42,7} = 1,\ \nu_{42,19} = 1,\ \nu_{42,31} = -1$ |
| 43 | $k_3^{O/OS}D_{OB}^{OS}$ | $\nu_{43,8} = 1,\ \nu_{43,20} = 1,\ \nu_{43,32} = -1$ |
| 44 | $g_M^N + k_2^{N/NN}D_{NB}^{NN} + k_2^{N/OS}D_{NB}^{OS}$ | $\nu_{44,33} = 1$ |
| 45 | $d_M^N M_G$ | $\nu_{45,33} = -1$ |
| 46 | $g_P^N M_G$ | $\nu_{46,34} = 1$ |
| 47 | $d_M^N M_G$ | $\nu_{47,34} = -1$ |

## A.2    The CME Model of the Nanog Autoregulatory Loop

For Scenario 1 (2), the CME model of the Nanog autoregulatory loop consists of 7 (9) molecular species and 11 (14) chemical reactions. The state vector, $x$, whose $i$th element is the copy number of molecular species $i$, is:

$$x = (M, P, C, D_A', D_B', D_A, D_B, M_G, P_G),$$

where the elements in black are those for Scenario 1, and the additional two species for Scenario 2 are in red. The chemical reactions are described in Table A.4, where the chemical reactions in Scenario 1 are in black and cyan (reaction 2 does not occur in Scenario 2), and the additional reactions in Scenario 2 are in red.

Table A.4: Descriptions of the chemical reactions involved in the Nanog autoregulatory loop. Reactions in black occur in both Scenarios 1 and 2, that in cyan occurs in Scenario 1 only, and those in red occur in Scenario 2 only. Reaction 12 replaces 2 in Scenario 2.

| Reaction | Description |
|---|---|
| 1 | N mRNA transcription from N allele A |
| 2 | N mRNA transcription from N allele B |
| 3 | N mRNA decay |
| 4 | N protein translation |
| 5 | N protein decay |
| 6 | N proteins dimerise |
| 7 | N dimer disassociates |
| 8 | NN dimer binds to unbound binding site on N allele A |
| 9 | NN dimer binds to unbound binding site on N allele B |
| 10 | NN dimer unbinds from binding site on N allele A |
| 11 | NN dimer unbinds from binding site on N allele B |
| 12 | GFP mRNA transcription from allele B |
| 13 | GFP mRNA decay |
| 14 | GFP protein translation |
| 15 | GFP protein decay |

The state-change vectors $\boldsymbol{\nu}_j$ and propensity functions $a_j(\boldsymbol{x})$ for all the reactions can be written in the form of a stoichiometric matrix $\boldsymbol{S}$, whose $j$th column is the state-change vector $\boldsymbol{\nu}_j$, and a propensity vector $\boldsymbol{a}(\boldsymbol{x})$, whose $j$th entry is the propensity function $a_j(\boldsymbol{x})$. For this model we have:

$$\boldsymbol{S} = \begin{pmatrix} 1 & 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 & -2 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & -1 & -1 & -1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{pmatrix} \tag{A.1}$$

$$\boldsymbol{a}(\boldsymbol{x}) = (g_M + k_2 D_A,\ g_M + k_2 D_B,\ d_M M,\ g_P M,\ d_P P,\ aP(P-1)/2,\ uC, \tag{A.2}$$

$$k_1 C D'_A,\ k_1 C D'_B,\ k_3 D_A,\ k_3 D_B\ g_M + k_2 D_B,\ d_M M_G,\ g_P M_G,\ d_P P_G)'$$

where the elements in black and cyan are for Scenario 1, and those in red and black are those for Scenario 2 (those in cyan are zero for Scenario 2).

## A.3    RRE Dimension Reduction

To reduce the number of dimensions of the RREs in Section 2.3.2, we begin by assuming that dimer (dis)association, and DNA (un)binding occurs at a much faster time scale than transcription and translation. Since our goal is to describe the effective dynamics of the system over long times, we use the quasi-equilibrium approximation. Setting Eqs. (2.13) and (2.14) equal to zero, we obtain

$$D_A = D_B = \frac{C}{k_3/k_1 + C}, \tag{A.3}$$

$$C = \frac{a}{2u}P^2, \tag{A.4}$$

and substituting Eq. (A.4) in (A.3), we find

$$D_A = D_B = \frac{P^2}{\frac{2uk_3}{ak_1} + P^2}. \tag{A.5}$$

Now substituting Eq. (A.5) in Eq. (2.8), and Eq. (A.4) in (2.9), we obtain the following effective dynamics for $M$ and $P$ in Scenario 1:

$$\frac{dM}{dt} = 2g_M + \frac{2k_2 P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M M \tag{A.6}$$

$$\frac{dP}{dt} = g_P M - d_P P, \tag{A.7}$$

and similarly for Scenario 2, the effective dynamics for $M$, $P$, $M_G$ and $P_G$ become

$$\frac{dM}{dt} = g_M + \frac{k_2 P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M M \tag{A.8}$$

$$\frac{dP}{dt} = g_P M - d_P P \tag{A.9}$$

$$\frac{dM_G}{dt} = g_M + \frac{k_2 P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M M_G \tag{A.10}$$

$$\frac{dP_G}{dt} = g_P M_G - d_P P_G. \tag{A.11}$$

Since we are interested in protein concentration, we eliminate the RREs for $M$ and $M_G$ using Eqs. (A.7) and (A.11) to see that the nullclines of $P$ and $P_G$ satisfy

$$P = \frac{g_P}{d_P} M, \quad P_G = \frac{g_P}{d_P} M_G,$$

together with the chain rule

$$\frac{dM}{dt} = \frac{dM}{dP}\frac{dP}{dt} = \frac{d_P}{g_P}\frac{dP}{dt}$$

to obtain, for Scenario 1:

$$\frac{dP}{dt} = 2g_M \frac{g_P}{d_P} + \frac{2k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P, \tag{A.12}$$

and for Scenario 2:

$$\frac{dP}{dt} = g_M \frac{g_P}{d_P} + \frac{k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P \tag{A.13}$$

$$\frac{dP_G}{dt} = g_M \frac{g_P}{d_P} + \frac{k_2 \frac{g_P}{d_P} P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P_G. \tag{A.14}$$

## A.4    Non-dimensionalisation of the RREs

The following RRE for Nanog concentration was derived in Section 2.3.2:

$$\frac{dP}{dt} = 2g_M\frac{g_P}{d_P} + \frac{2k_2\frac{g_P}{d_P}P^2}{\frac{2uk_3}{ak_1} + P^2} - d_M P. \tag{A.15}$$

To nondimensionalise the system we introduce the scalings $P = aX$ and $t = b\tau$, where $X$ and $\tau$ are dimensionless variables. Thus

$$\frac{dP}{dt} = \frac{dP}{dX}\frac{dX}{d\tau}\frac{d\tau}{dt} = \frac{a}{b}\frac{dX}{d\tau}.$$

So Equation (A.15) becomes:

$$\frac{a}{b}\frac{dX}{d\tau} = 2g_M\frac{g_P}{d_P} + \frac{a^2 2k_2\frac{g_P}{d_P}X^2}{\frac{2uk_3}{ak_1} + a^2X^2} - ad_M X$$

$$\Rightarrow \quad \frac{dX}{d\tau} = \frac{b}{a}2g_M\frac{g_P}{d_P} + \frac{ab2k_2\frac{g_P}{d_P}X^2}{\frac{2uk_3}{ak_1} + a^2X^2} - bd_M X$$

$$\Rightarrow \quad \frac{dX}{d\tau} = \frac{b}{a}2g_M\frac{g_P}{d_P} + \frac{X^2}{\frac{1}{ab}\frac{2uk_3d_P}{2k_2ak_1g_P} + \frac{a}{b}\frac{d_P}{2k_2g_P}X^2} - bd_M X.$$

To eliminate the coefficients of $X$ and $X^2$, we set

$$b = \frac{1}{d_M}, \quad a = \frac{2k_2g_P}{d_P d_M}.$$

Therefore the ODE becomes:

$$\frac{dX}{dt} = \alpha + \frac{X^2}{\gamma^2 + X^2} - X, \tag{A.16}$$

where

$$\alpha = \frac{g_M}{k_2}, \quad \gamma^2 = \frac{2uk_3}{ak_1}\left(\frac{d_M d_P}{2k_2g_P}\right)^2$$

are dimensionless.

# Appendix B

# Derivations

## B.1 Derivation of the Parameterised Potential Function

The ordinary differential equation for the deterministic dynamics of Sca-1 expression is:

$$\frac{\mathrm{d}\psi}{\mathrm{d}x} = \beta x - \alpha_0 - \frac{\alpha_1 x^n}{K^n + x^n}.$$

Integrating both sides we obtain:

$$\psi(x) = \int \left( -\alpha_0 - \frac{\alpha_1 x^n + \alpha_1 K^n - \alpha_1 K^n}{K^n + x^n} + \beta x \right) \, dx = \int \left( -(\alpha_0 + \alpha_1) + \frac{\alpha_1 K^n}{k^n + x^n} + \beta x \right) \, dx$$

$$= \frac{1}{2}\beta x^2 - (\alpha_0 + \alpha_1)x + \alpha_1 K^n I_n,$$

where $I_n = \int \frac{1}{K^n + x^n} \, dx$. Now we derive an expression for $I_n$ for $n$ positive and even.

Let $F(x, n, K) = K^n + x^n$, so that $I_n = \int 1/F \, dx$. The aim is to express $1/F$ as a sum of easily integrable terms. We begin by expressing $F$ as a product of linear factors; specifically,

$$F(x, n, K) = (x - r_1)(x - r_2) \dots (x - r_n),$$

where $r_1, r_2, \dots, r_n$ are the roots of $F$. Now the roots of $K^n + x^n$ are equal to the $n$th roots of $-K^n$, which can be found using Euler's identity $\exp(i\pi) = -1$ together with Euler's formula, which gives $\exp((2k - 1)i\pi) = -1$, for $k = 1, \dots n$, and then multiplying both sides by $K^n$

we obtain:

$$- K^n = K^n \exp((2k - 1)i\pi) \Rightarrow (-K^n)^{\frac{1}{n}} = K \exp((2k - 1)i\pi/n).$$

Thus the roots of $F$ are:

$$r_k = K \exp((2k - 1)i\pi/n), \tag{B.1}$$

for $k = 1, \ldots n$. Now we use partial fraction decomposition to express $1/F$ as the sum of terms of the form $\frac{A_k}{(x - r_k)}$, for suitable coefficients $A_k$. Therefore, we wish to solve

$$\frac{1}{F} = \sum_{k=1}^{n} \frac{A_k}{(x - r_k)}$$

for $A_k$. Cross-multiplication leads to

$$\sum_{k=1}^{n} A_k (x - r_1)(x - r_2) \ldots (x - r_{k-1})(x - r_{k+1}) \ldots (x - r_n) = 1,$$

and substituting $x = r_k$ gives

$$A_k = (r_k - r_1)(r_k - r_2) \ldots (r_k - r_{k-1})(r_k - r_{k+1}) \ldots (r_k - r_n).$$

This product can be simplified as follows: Since

$$K^n + x^n = (x - r_1)(x - r_2) \ldots (x - r_n),$$

dividing by both sides by $(x - r_k)$ gives:

$$\frac{K^n + x^n}{(x - r_k)} = (x - r_1)(x - r_2) \ldots (x - r_{k-1})(x - r_{k+1}) \ldots (x - r_n).$$

Now using polynomial division on the left-hand-side, we obtain:

$$x^{n-1} + r_k x^{n-2} + r_k^3 x^{n-3} + \ldots r_k^{n-1} = (x - r_1)(x - r_2) \ldots (x - r_{k-1})(x - r_{k+1}) \ldots (x - r_n),$$

and substituting $x = r_k$ gives

$$(r_k - r_1)(r_k - r_2) \ldots (r_k - r_{k-1})(r_k - r_{k+1}) \ldots (r_k - r_n) = r_k^n \, r_k^{-1} = -n \, K^n \, r_k^{-1},$$

and therefore

$$A_k = -r_k/nK^n,$$

and

$$1/F = \sum_{k=1}^{n} \frac{-r_k/nK^n}{x - r_k}.$$

Since the complex roots of $F$ come in conjugate pairs, say $r_k$ and $\bar{r}_k$,

$$\begin{aligned} 2/F &= \sum_{k=1}^{n} \left( \frac{-r_k/nK^n}{x - r_k} + \frac{-\bar{r}_k/nK^n}{x - \bar{r}_k} \right) \\ &= -\frac{1}{nK^n} \sum_{k=1}^{n} \frac{r_k(x - \bar{r}_k) + \bar{r}_k(x - r_k)}{(x - r_k)(x - \bar{r}_k)}. \end{aligned} \tag{B.2}$$

Now using Euler's formula $r_k$ and $\bar{r}_k$ can be written as

$$r_k = K\cos((2k-1)\pi/n) + iK\sin((2k-1)\pi/n)$$

$$\bar{r}_k = K\cos((2k-1)\pi/n) - iK\sin((2k-1)\pi/n),$$

so

$$r_k + \bar{r}_k = 2K\cos((2k-1)\pi/n); \quad r_k\bar{r}_k = K^2.$$

Substituting this result into Eq. (B.2) we obtain:

$$1/F = \frac{1}{nK^n} \sum_{k=1}^{n} \frac{1 - C_k(x/K)}{(x/K)^2 - 2C_k(x/K) + 1}, \tag{B.3}$$

where $S_k = \sin((2k-1)\pi/n)$ and $C_k = \cos((2k-1)\pi/n)$. Next, when we substitute the identity $S_k^2 + C_k^2 = 1$ in the numerator, the $k$th summand becomes:

$$\frac{S_k^2}{(x/K)^2 - 2C_k(x/K) + 1} - \frac{C_k\left((x/K) - C_k\right)}{(x/K)^2 - 2C_k(x/K) + 1}.$$

Substituting the same identity in the denominator of the first term and then completing the square, we obtain:

$$\frac{S_k^2}{((x/K) - C_k)^2 + S_k^2} - \frac{C_k\left((x/K) - C_k\right)}{(x/K)^2 - 2C_k(x/K) + 1},$$

which easily integrates to

$$S_k K \arctan\left(\frac{(x/K) - C_k}{S_k}\right) - \frac{C_k K}{2} \log\left((x/K)^2 - 2C_k(x/K) + 1\right).$$

For $n$ even $C_1 = C_n,\ C_2 = C_{n-1}, \ldots, C_{n/2} = C_{n+2/2}$, and $S_1 = -S_n,\ S_2 = -S_{n-1}, \ldots, C_{n/2} = C_{n+2/2}$, and thus the integral, $I_n$ becomes:

$$I_n = \frac{1}{nK^{n-1}} \sum_{k=1}^{n/2} \left( 2S_k \arctan\left( \frac{\frac{x}{K} - C_k}{S_k} \right) - C_k \log\left( \left(\frac{x}{K}\right)^2 - 2C_k\frac{x}{K} + 1 \right) \right),$$

and

$$\psi(x) = \frac{1}{2}\beta x^2 - (\alpha_0 + \alpha_1)x +$$

$$\frac{\alpha_1 K}{n} \sum_{k=1}^{n/2} \left[ 2S_k \arctan\left( \frac{\frac{x}{K} - C_k}{S_k} \right) - C_k \log\left( \left(\frac{x}{K}\right)^2 - 2C_k\frac{x}{K} + 1 \right) \right],$$

where $S_k = \sin\left((2k-1)\pi/n\right), C_k = \cos\left((2k-1)\pi/n\right)$.

So for example for $n = 2$, $S_1 = \sin\left(\pi/2\right) = 1,\ C_1 = \cos\left(\pi/2\right) = 0$, therefore

$$\psi(x) = \frac{1}{2}\beta x^2 - (\alpha_0 + \alpha_1)x + \alpha_1 K \arctan\left(\frac{x}{K}\right).$$

## B.2 Non-Dimensionalisation of the SDE Model for Sca1 Expression

In order to rescale the parameterised stochastic differential equation for Sca1

$$\frac{\mathrm{d}x}{\mathrm{d}t} = \beta x - \alpha_0 - \frac{\alpha_1 x^n}{K^n + x^n} + \sqrt{2\sigma}\xi(t), \tag{B.4}$$

we begin by introducing the scalings $x = aX$ and $t = b\tau$, where $X$ and $\tau$ are dimensionless variables:

$$\frac{a\,dx}{b\,dt} = \alpha_0 + \frac{\alpha_1(ax)^n}{K^n + (ax)^n} - \beta ax + \frac{\sqrt{2\sigma}}{\sqrt{b}}\xi(t)$$

$$\Rightarrow \frac{dx}{dt} = \frac{b}{a}\alpha_0 + \frac{b}{a}\frac{\alpha_1(ax)^n}{K^n + (ax)^n} - b\beta x + \frac{\sqrt{b}}{a}\sqrt{2\sigma}\xi(t).$$

First, let $b = 1/\beta$:

$$\frac{dx}{dt} = \frac{\alpha_0}{a\beta} + \frac{1}{a\beta}\frac{\alpha_1(ax)^n}{K^n + (ax)^n} - x + \frac{1}{a}\sqrt{\frac{2\sigma}{\beta}}\xi(t),$$

since $\frac{dW_t}{dt} \sim \frac{1}{\sqrt{dt}}\mathcal{N}(0,1)$. Now let $a = \frac{\alpha_1}{\beta}$:

$$\frac{dx}{dt} = \frac{\alpha_0}{\alpha_1} + \frac{x^n}{\left(\frac{K\beta}{\alpha_1}\right)^n + x^n} - x + \frac{\beta}{\alpha_1}\sqrt{\frac{2\sigma}{\beta}}\xi(t)$$

$$= \frac{\alpha_0}{\alpha_1} + \frac{x^n}{\left(\frac{K\beta}{\alpha_1}\right)^n + x^n} - x + \sqrt{2\frac{\beta}{\alpha_1^2}\sigma}\xi(t)$$

$$= \alpha + \frac{x^n}{\gamma^n + x^n} - x + \sqrt{2\sigma_d}\xi(t),$$

where $\alpha = \frac{\alpha_0}{\alpha_1}$, $\gamma = \frac{K\beta}{\alpha_1}$, and $\sigma_d = \frac{\beta}{\alpha_1^2}\sigma$. Therefore Eq. (B.4) may then be expressed in nondimensional form:

$$\frac{dX}{dt} = \alpha + \frac{x^n}{\gamma^n + x^n} - x + \sqrt{2\sigma_d}\xi(t).$$

## B.3   The Resultant Matrix

The resultant, $R$, of the polynomial

$$X^{n+1} - (1+\alpha)X^n + \gamma^n X - \alpha\gamma^n = 0,$$

is $R =$

$$
\begin{vmatrix}
1 & -(1+a) & 0 & \ldots & 0 & \gamma^n & -\alpha\gamma^n & 0 & \ldots & 0 & 0 \\
0 & 1 & -(1+a) & 0 & \ldots & 0 & \gamma^n & -\alpha\gamma^n & 0 & \ldots & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & \ldots & 0 & 1 & -(1+a) & 0 & \ldots & 0 & \gamma^n & -\alpha\gamma^n \\
n+1 & -(1+a)n & 0 & \ldots & 0 & \gamma^n & 0 & 0 & \ldots & 0 & 0 \\
0 & n+1 & -(1+a)n & 0 & \ldots & 0 & \gamma^n & 0 & 0 & \ldots & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & \ldots & 0 & n+1 & -(1+a)n & 0 & \ldots & 0 & \gamma^n & 0 \\
0 & 0 & 0 & \ldots & 0 & n+1 & -(1+a)n & 0 & \ldots & 0 & \gamma^n
\end{vmatrix}
$$

Solving

$$(-1)^{n(n+1)/2} R = 0$$

for $\gamma$ to we can obtain an expression for the critical values as function of $\alpha$. In this case, $n = 12$, and the positive, real solutions are: $\gamma_\pm(\alpha) =$

$$
\begin{aligned}
&\Bigg[ \frac{11^{11}}{2^{25}3^{12}} + \frac{2\times 11^{10}}{2^{22}3^{10}}\alpha + \frac{2\times 11^8\, 61}{2^{20}3^9}\alpha^2 + \frac{11^6\, 6833}{2^{17}3^9}\alpha^3 + \frac{11^4\, 19\times 2039}{2^{16}3^6}\alpha^4 + \frac{2\, 11^2\, 267667}{2^{12}3^6}\alpha^5 \\
&+ \frac{2\times 13\times 613\times 1381}{2^{10}3^6}\alpha^6 + \frac{11\times 13\times 179\times 7^2}{2^9 3^4}\alpha^7 + \frac{11\times 13\times 757}{2^9 3^2}\alpha^8 + \frac{2\times 11\times 13\times 83}{2^6 3^5}\alpha^9 \\
&+ \frac{2\times 11\times 13}{2^4\, 3}\alpha^{10} + \frac{5\times 7}{2}\alpha^{11} - \alpha^{12} \\
&\pm (121 - 48\alpha)^{3/2}\Bigg( \frac{11^8}{2^{25}3^{12}} + \frac{23\times 11^6}{2^{22}3^{10}}\alpha + \frac{3\times 11^4\times 17\times 29}{2^{20}3^9}\alpha^2 + \frac{2\times 3\times 11^2\times 47\times 313}{2^{17}3^9}\alpha^3 \\
&+ \frac{5^2\, 29\times 373}{2^{15}3^6}\alpha^4 + \frac{19\times 3581}{2^{12}3^6}\alpha^5 + \frac{11\times 2377}{2^{10}3^6}\alpha^6 + \frac{3\times 11\times 53}{2^9 3^4}\alpha^7 + \frac{11\times 17}{2^9 3^2}\alpha^8 + \frac{3^2\times 5\times 11}{2^6 3^5}\alpha^9 \\
&+ \frac{1}{2^4\, 3}\alpha^{10}\Bigg)\Bigg]^{1/12}, \quad \alpha \le \frac{121}{48},
\end{aligned}
$$

(B.5)

where the large coefficients have been decomposed into their prime factors to demonstrate that it is not possible to simplify the equations.

# Appendix C

# Enumerating Feedback Loops

The following computational algorithm generates a list of all the cycles of length $L \leq 5$ in a directed network, from which the number of cycles each node participates in can be computed.

1. Calculate the adjacency matrix $\boldsymbol{A}$.

2. Identify and record the nodes that participate in self–loops ($a_{ii} = 1$). These are cycles of length 1.

3. Remove self–loops from the adjacency matrix by setting the diagonal entries equal to zero.

4. Create a vector, $\boldsymbol{v}$, containing the node labels, denoted by the integers 1:9.

5. Remove from $\boldsymbol{v}$ the nodes that do not have both an incoming link and an outgoing link (these nodes cannot participate in a cycle). In this case, only 5 nodes remain in $\boldsymbol{v}$ as there are no cycles greater than length 5.

6. List all possible combinations of nodes of length 2 to 5, in vector form, and store in a list variable, $\boldsymbol{B}$.

7. Sort all combinations in $\boldsymbol{B}$ in ascending order.

8. Keeping the first element (the minimum) fixed, list all possible permutations of each combination in $\boldsymbol{B}$. Fixing the first element ensures that there are no duplicated cycles. Store these permutations (vectors) in a list variable, $\boldsymbol{C}$.

9. For each permutation in $\boldsymbol{C}$, repeat the first element at the end of the vector, e.g. $(2, 4, 5) \rightarrow (2, 4, 5, 2)$. This represents a complete cycle. $\boldsymbol{C}$ is now a list of all possible cycles.

10. Use the adjacency matrix to test if each cycle in $\boldsymbol{C}$ is a true cycle, as follows: A cycle $\boldsymbol{c}$ with entries $c_i$ is a valid cycle if all $a_{c_i c_{i+1}} = 1$. Thus, if $\prod_{i=1} a_{c_i c_{i+1}} = 1$, keep $\boldsymbol{c}$ in $\boldsymbol{C}$, and if $\prod_{i=1} a_{c_i c_{i+1}} = 0$, remove $\boldsymbol{c}$ from $\boldsymbol{C}$.

11. Use list of valid cycles, $\boldsymbol{C}$, to count the number of cycles of each length.

12. For each node, count the number of cycles of each length that it participates in.

# Bibliography

[1] S. J. Ridden, H. H. Chang, K. C. Zygalakis, and B. D. MacArthur. Entropy, ergodicity, and stem cell multipotency. *Physical Review Letters*, 115:208103, 2015.

[2] A. Ma'ayan and B. D. MacArthur. *New frontiers of network analysis in systems biology.* Springer Science & Business Media, 2012.

[3] B. D. MacArthur, A. Sevilla, M. Lenz, F. Müller, B. M. Schuldt, A. A. Schuppert, S. J. Ridden, P. S. Stumpf, A. Fidalgo, M.and Ma'ayan, J. Wang, and I. R. Lemischka. Nanog-dependent feedback loops regulate murine embryonic stem cell heterogeneity. *Nature cell biology*, 14(11):1139–1147, 2012.

[4] R. C. Smith, P. S. Stumpf, S. J. Ridden, A. Sim, S. Filippi, H. Harrington, and B. D. MacArthur. Nanog fluctuations in es cells highlight the problem of measurement in cell biology. *bioRxiv*, page 060558, 2016.

[5] S. T. Kosak and M. Groudine. Form follows function: the genomic organization of cellular differentiation. *Genes and development*, 18(12):1371–1384, 2004.

[6] B. John. The biology of heterochromatin. *Heterochromatin: molecular and structural aspects*, pages 1–147, 1988.

[7] B. Kholodenko. Cell-signalling dynamics in time and space. *Nature reviews Molecular cell biology*, 7(3):165–176, 2006.

[8] P. A. Jones and S. B. Baylin. The fundamental role of epigenetic events in cancer. *Nature reviews genetics*, 3(6):415–428, 2002.

[9] T. M. Geiman and K. D. Robertson. Chromatin remodeling, histone modifications, and dna methylation – how does it all fit together? *Journal of cellular biochemistry*, 87(2):117–125, 2002.

[10] D. Schübeler, C. Francastel, D. M. Cimbora, A. Reik, D. I. Martin, and M. Groudine. Nuclear localization and histone acetylation: a pathway for chromatin opening and transcriptional activation of the human $\beta$-globin locus. *Genes & Development*, 14(8): 940–950, 2000.

[11] A. Eberharter and P. B. Becker. Histone acetylation: a switch between repressive and permissive chromatin. *EMBO reports*, 3(3):224–229, 2002.

[12] L. Ivana, Y. Ohkawa, C. A. Berkes, D. A. Bergstrom, C. S. Dacwag, S. J. Tapscott, and A. N. Imbalzano. Myod targets chromatin remodeling complexes to the myogenin locus prior to forming a stable dna-bound complex. *Molecular and cellular biology*, 25 (10):3997–4009, 2005.

[13] J. A. Grass, M. E. Boyer, S. Pal, J. Wu, M. J. Weiss, and E. H. Bresnick. Gata-1-dependent transcriptional repression of gata-2 via disruption of positive autoregulation and domain-wide chromatin remodeling. *Proceedings of the National Academy of Sciences*, 100(15):8811–8816, 2003.

[14] E. Metzger, M. Wissmann, and R. Schule. Histone demethylation and androgen-dependent transcription. *Current opinion in genetics & development*, 16(5):513–517, 2006. ISSN 0959-437X.

[15] M. Suzuki, T. Yamada, F. Kihara-Negishi, T. Sakurai, E. Hara, D. Tenen, N. Hozumi, and T. Oikawa. Site-specific DNA methylation by a complex of PU. 1 and Dnmt3a/b. *Oncogene*, 25(17):2477–2488, 2005. ISSN 0950-9232.

[16] S. Kubicek and T. Jenuwein. A crack in histone lysine methylation. *Cell*, 119(7): 903–906, 2004. ISSN 0092-8674.

[17] P. Trojer and D. Reinberg. Histone lysine demethylases and their impact on epigenetics. *Cell*, 125(2):213–217, 2006. ISSN 0092-8674.

[18] T. A. Brevini and P. GEORGIA. *Gametogenesis, Early Embryo Development and Stem Cell Derivation*. Springer Science & Business Media, 2012.

[19] L. Li and T. Xie. Stem cell niche: structure and function. *Annual Review of Cell and Developmental Biology*, 21:605–631, 2005.

[20] M. Pittenger, A. Mackay, S. Beck, R. Jaiswal, R. Douglas, J. Mosca, M. Moorman, D. Simonetti, S. Craig, and D. Marshak. Multilineage potential of adult human mesenchymal stem cells. *Science*, 284(5411):143, 1999.

[21] O. Naveiras and G. Daley. Stem cells and their niche: a matter of fate. *Cellular and molecular life sciences*, 63(7):760–766, 2006.

[22] M. Evans, M. Kaufman, et al. Establishment in culture of pluripotential cells from mouse embryos. *Nature*, 292(5819):154–156, 1981.

[23] G. R. Martin. Isolation of a pluripotent cell line from early mouse embryos cultured in medium conditioned by teratocarcinoma stem cells. *Proceedings of the National Academy of Sciences*, 78(12):7634–7638, 1981.

[24] J. Nichols and A. Smith. Naive and primed pluripotent states. *Cell stem cell*, 4(6): 487–492, 2009.

[25] Y. Suda, M. Suzuki, Y. Ikawa, and S. Aizawa. Mouse embryonic stem cells exhibit indefinite proliferative potential. *Journal of cellular physiology*, 133(1):197–201, 1987.

[26] A. P. Beltrami, L. Barlucchi, D. Torella, M. Baker, F. Limana, S. Chimenti, H. Kasahara, M. Rota, E. Musso, K. Urbanek, et al. Adult cardiac stem cells are multipotent and support myocardial regeneration. *Cell*, 114(6):763–776, 2003.

[27] S. J. Morrison, N. M. Shah, and D. J. Anderson. Regulatory mechanisms in stem cell biology. *Cell*, 88(3):287–298, 1997.

[28] S. J. Morrison and J. Kimble. Asymmetric and symmetric stem-cell divisions in development and cancer. *Nature*, 441(7097):1068–1074, 2006.

[29] F. M. Watt and B. L. Hogan. Out of eden: stem cells and their niches. *Science*, 287 (5457):1427, 2000.

[30] M. Ramalho-Santos, S. Yoon, Y. Matsuzaki, R. C. Mulligan, and D. A. Melton. " stemness": transcriptional profiling of embryonic and adult stem cells. *Science*, 298 (5593):597–600, 2002.

[31] W. Bloom. Cellular differentiation and tissue culture. *Physiological Reviews*, 17(4): 589–617, 1937.

[32] R. Jaenisch and R. Young. Stem cells, the molecular circuitry of pluripotency and nuclear reprogramming. *Cell*, 132(4):567–582, 2008.

[33] C. Starr, R. Taggart, and L. Starr. *Biology: The unity and diversity of life.* Wadsworth Publishing Company, 1987.

[34] P. R. Brauer. *Human embryology: the ultimate USMLE step 1 review.* Elsevier Health Sciences, 2003.

[35] F. M. Kamm. *Bioethical prescriptions: to create, end, choose, and improve lives.* Oxford University Press, 2013.

[36] I. Chambers and A. Smith. Self-renewal of teratocarcinoma and embryonic stem cells. *Oncogene*, 23(43):7150–7160, 2004.

[37] J. A. Thomson, J. Itskovitz-Eldor, S. S. Shapiro, M. A. Waknitz, J. J. Swiergiel, V. S. Marshall, and J. M. Jones. Embryonic stem cell lines derived from human blastocysts. *science*, 282(5391):1145–1147, 1998.

[38] A. Wagers and I. Weissman. Plasticity of adult stem cells. *Cell*, 116(5):639–648, 2004. ISSN 0092-8674.

[39] M. E. Mabon, X. Mao, Y. Jiao, B. A. Scott, and C. M. Crowder. Systematic identification of gene activities promoting hypoxic death. *Genetics*, 181(2):483–496, 2009.

[40] R. McKay. Stem cells in the central nervous system. *Science*, 276(5309):66–71, 1997.

[41] J. Ray, T. Palmer, J. Suhonen, J. Takahashi, and F. Gage. Neurogenesis in the adult brain: Lessons learned from the studies of progenitor cells from the embryonic and adult central nervous systems. In *Isolation, Characterization and Utilization of CNS Stem Cells*, pages 129–149. Springer, 1997.

[42] P. Quesenberry and L. Levitt. Hematopoietic stem cells. *New England Journal of Medicine*, 301(14):755–760, 1979.

[43] E. Gunsilius, G. Gastl, and A. Petzer. Hematopoietic stem cells. *Biomedicine & pharmacotherapy*, 55(4):186–194, 2001.

[44] J. J. Minguell, A. Erices, and P. Conget. Mesenchymal stem cells. *Experimental biology and medicine*, 226(6):507–520, 2001.

[45] A. L. Kierszenbaum and L. Tres. *Histology and cell biology: an introduction to pathology.* Elsevier Health Sciences, 2015.

[46] J. E. Morgan and T. A. Partridge. Muscle satellite cells. *The international journal of biochemistry & cell biology*, 35(8):1151–1156, 2003.

[47] A. Zhang. *Protein interaction networks: computational analysis.* Cambridge University Press, 2009.

[48] J. Zhu, H. B. Larman, G. Gao, R. Somwar, Z. Zhang, U. Laserson, A. Ciccia, N. Pavlova, G. Church, W. Zhang, et al. Protein interaction discovery using parallel analysis of translated orfs (plato). *Nature biotechnology*, 31(4):331–334, 2013.

[49] P. Uetz, L. Giot, G. Cagney, T. A. Mansfield, R. S. Judson, J. R. Knight, D. Lockshon, V. Narayan, M. Srinivasan, P. Pochart, et al. A comprehensive analysis of protein–protein interactions in saccharomyces cerevisiae. *Nature*, 403(6770):623–627, 2000.

[50] C. Von Mering, R. Krause, B. Snel, M. Cornell, S. G. Oliver, S. Fields, and P. Bork. Comparative assessment of large-scale data sets of protein–protein interactions. *Nature*, 417(6887):399–403, 2002.

[51] J. Rual, K. Venkatesan, T. Hao, T. Hirozane-Kishikawa, A. Dricot, N. Li, G. Berriz, F. Gibbons, M. Dreze, N. Ayivi-Guedehoussou, et al. Towards a proteome-scale map of the human protein–protein interaction network. *Nature*, 437(7062):1173–1178, 2005. ISSN 0028-0836.

[52] C. Landgraf, S. Panni, L. Montecchi-Palazzi, L. Castagnoli, J. Schneider-Mergener, R. Volkmer-Engert, and G. Cesareni. Protein interaction networks by proteome peptide scanning. *PLoS Biol.*, 2(1):e14, 2004.

[53] U. Stelzl, U. Worm, M. Lalowski, C. Haenig, F. Brembeck, H. Goehler, M. Stroedicke, M. Zenkner, A. Schoenherr, S. Koeppen, et al. A human protein-protein interaction network: a resource for annotating the proteome. *Cell*, 122(6):957–968, 2005.

[54] J. Bader, A. Chaudhuri, J. Rothberg, and J. Chant. Gaining confidence in high-throughput protein interaction networks. *Nature Biotechnology*, 22(1):78–85, 2004.

[55] L. Giot, J. Bader, C. Brouwer, A. Chaudhuri, B. Kuang, Y. Li, Y. Hao, C. Ooi, B. Godwin, E. Vitols, et al. A protein interaction map of drosophila melanogaster. *Science*, 302(5651):1727, 2003.

[56] O. Rinner, L. Mueller, M. Hubálek, M. Müller, M. Gstaiger, and R. Aebersold. An integrated mass spectrometric and computational framework for the analysis of protein interaction networks. *Nature Biotechnology*, 25(3):345–352, 2007.

[57] A. Tong, B. Drees, G. Nardelli, G. Bader, B. Brannetti, L. Castagnoli, M. Evangelista, S. Ferracuti, B. Nelson, S. Paoluzi, et al. A combined experimental and computational strategy to define protein interaction networks for peptide recognition modules. *Science*, 295(5553):321, 2002.

[58] L. Boyer, T. Lee, M. Cole, S. Johnstone, S. Levine, J. Zucker, M. Guenther, R. Kumar, H. Murray, R. Jenner, et al. Core transcriptional regulatory circuitry in human embryonic stem cells. *Cell*, 122(6):947–956, 2005. ISSN 0092-8674.

[59] S. Shen-Orr, R. Milo, S. Mangan, and U. Alon. Network motifs in the transcriptional regulation network of Escherichia coli. *Nature Genetics*, 31(1):64–68, 2002.

[60] A.-L. Barabasi and Z. N. Oltvai. Network biology: understanding the cell's functional organization. *Nature reviews genetics*, 5(2):101–113, 2004.

[61] M. M. Babu, N. M. Luscombe, L. Aravind, M. Gerstein, and S. A. Teichmann. Structure and evolution of transcriptional regulatory networks. *Current opinion in structural biology*, 14(3):283–291, 2004.

[62] J. Stelling, S. Klamt, K. Bettenbrock, S. Schuster, and E. D. Gilles. Metabolic network structure determines key aspects of functionality and regulation. *Nature*, 420(6912): 190–193, 2002.

[63] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network motifs: simple building blocks of complex networks. *Science*, 298(5594):824, 2002.

[64] A. Mogilner, R. Wollman, and W. F. Marshall. Quantitative modeling in cell biology: what is it good for? *Developmental cell*, 11(3):279–287, 2006.

[65] J. Hasty, D. McMillen, F. Isaacs, and J. Collins. Computational studies of gene regulatory networks: in numero molecular biology. *Nature Reviews Genetics*, 2(4):268–279, 2001. ISSN 1471-0056.

[66] G. Karlebach and R. Shamir. Modelling and analysis of gene regulatory networks. *Nature Reviews Molecular Cell Biology*, 9(10):770–780, 2008.

[67] H. Bolouri and E. Davidson. Modeling transcriptional regulatory networks. *BioEssays*, 24(12):1118–1129, 2002. ISSN 1521-1878.

[68] C. Waddington. *The strategy of the genes: a discussion of some aspects of theoretical biology.* Allen & Unwin, 1957.

[69] F. Crick. On protein synthesis. In *Symposia of the Society for Experimental Biology*, volume 12, page 138, 1958.

[70] J. Watson and F. Crick. Molecular structure of nucleic acids. *Nature*, 171(4356): 737–738, 1953.

[71] P. Andrews. From teratocarcinomas to embryonic stem cells. *Philosophical Transactions of the Royal Society B*, 357(1420):405, 2002.

[72] R. Thomas et al. Laws for the dynamics of regulatory networks. *The International Journal of Developmental Biology*, 42:479–485, 1998.

[73] M. Delbruck. Discussion to paper by sonneborn and beale. *Unites biolo*, pages 33–35, 1949.

[74] J. Monod and F. Jacob. General conclusions: teleonomic mechanisms in cellular metabolism, growth, and differentiation. In *Cold Spring Harbor Symposia on Quantitative Biology*, volume 26, page 389. Cold Spring Harbor Laboratory Press, 1961.

[75] S. Kauffman. Metabolic stability and epigenesis in randomly constructed genetic nets. *Journal of Theoretical Biology*, 22(3):437–467, 1969. ISSN 0022-5193.

[76] S. Kauffman. *The origins of order*, volume 209. Oxford University Press New York, 1993.

[77] S. Huang, G. Eichler, Y. Bar-Yam, and D. Ingber. Cell fates as high-dimensional attractor states of a complex gene regulatory network. *Physical Review Letters*, 94(12): 128701, 2005. ISSN 1079-7114.

[78] D. J. Higham. Modeling and simulating chemical reactions. *SIAM review*, 50(2):347–368, 2008.

[79] D. Gillespie. The chemical langevin equation. *The Journal of Chemical Physics*, 113: 297, 2000.

[80] D. T. Gillespie. A rigorous derivation of the chemical master equation. *Physica A*, 188: 404–425, 1992.

[81] D. T. Gillespie. Stochastic simulation of chemical kinetics. *Annual Review of Physical Chemistry*, 58:35–55, 2007.

[82] D. T. Gillespie and L. Petzold. Numerical simulation for biochemical kinetics. *Systems Modelling in Cellular Biology*, pages 331–354, 2006.

[83] M. Cox. *Molecular Biology*. Palgrave Macmillan, 2011.

[84] D. Gillespie. Deterministic limit of stochastic chemical kinetics. *The Journal of Physical Chemistry B*, 113(6):1640–1644, 2009.

[85] A. Golightly and C. S. Gillespie. Simulation of stochastic kinetic models. *In Silico Systems Biology*, pages 169–187, 2013.

[86] D. Gillespie. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of computational physics*, 22(4):403–434, 1976.

[87] D. Gillespie. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry*, 81(25):2340–2361, 1977.

[88] D. Gillespie. Approximate accelerated stochastic simulation of chemically reacting systems. *The Journal of Chemical Physics*, 115:1716, 2001.

[89] E. W. J. Wallace, D. T. Gillespie, K. R. Sanft, and L. R. Petzold. Linear noise approximation is valid over limited times for any chemical system that is sufficiently large. *IET systems biology*, 6(4):102–115, 2012.

[90] C. Gardiner. *Handbook of stochastic methods for physics, chemistry, and the natural sciences*. Springer-Verlag (Berlin and New York), 1985.

[91] H. McAdams and A. Arkin. Stochastic mechanisms in gene expression. *Proceedings of the National Academy of Sciences*, 94(3):814, 1997.

[92] N. van Kampen. *Stochastic processes in physics and chemistry*, volume 1. North Holland, 1992.

[93] P. Swain, M. Elowitz, and E. Siggia. Intrinsic and extrinsic contributions to stochasticity in gene expression. *Proceedings of the National Academy of Sciences*, 99(20):12795, 2002.

[94] D. Cook, A. Gerber, and S. Tapscott. Modeling stochastic gene expression: implications for haploinsufficiency. *Proceedings of the National Academy of Sciences*, 95(26):15641, 1998.

[95] P. Guptasarma. Does replication-induced transcription regulate synthesis of the myriad low copy number proteins of escherichia coli? *Bioessays*, 17(11):987–997, 1995.

[96] J. Spudich, D. Koshland Jr, et al. Non-genetic individuality: chance in the single cell. *Nature*, 262(5568):467, 1976.

[97] D. Austin, M. Allen, J. McCollum, R. Dar, J. Wilgus, G. Sayler, N. Samatova, C. Cox, and M. Simpson. Gene network shaping of inherent noise spectra. *Nature*, 439(7076): 608–611, 2006.

[98] A. Arias and P. Hayward. Filtering transcriptional noise during development: concepts and mechanisms. *Nature Reviews Genetics*, 7(1):34–44, 2006. ISSN 1471-0056.

[99] M. Kærn, T. Elston, W. Blake, and J. Collins. Stochasticity in gene expression: from theories to phenotypes. *Nature Reviews Genetics*, 6(6):451–464, 2005. ISSN 1471-0056.

[100] M. Thattai and A. Van Oudenaarden. Stochastic gene expression in fluctuating environments. *Genetics*, 167(1):523, 2004.

[101] A. Arkin, J. Ross, and H. McAdams. Stochastic kinetic analysis of developmental pathway bifurcation in phage $\lambda$-infected escherichia coli cells. *Genetics*, 149(4):1633–1648, 1998.

[102] H. McAdams and A. Arkin. It's a noisy business! genetic regulation at the nanomolar scale. *Trends in Genetics*, 15(2):65–69, 1999.

[103] H. Chang, M. Hemberg, M. Barahona, D. Ingber, and S. Huang. Transcriptome-wide noise controls lineage choice in mammalian progenitor cells. *Nature*, 453(7194):544–547, 2008. ISSN 0028-0836.

[104] J. , O. Berg, and M. Ehrenberg. Stochastic focusing: fluctuation-enhanced sensitivity of intracellular regulation. *Proceedings of the National Academy of Sciences*, 97(13): 7148, 2000.

[105] H. El Samad, M. Khammash, L. Petzold, and D. Gillespie. Stochastic modelling of gene regulatory networks. *International Journal of Robust and Nonlinear Control*, 15 (15):691–711, 2005.

[106] J. Vilar, H. Kueh, N. Barkai, and S. Leibler. Mechanisms of noise-resistance in genetic oscillators. *Proceedings of the National Academy of Sciences*, 99(9):5988, 2002.

[107] V. Karwacki-Neisius, J. Göke, R. Osorno, F. Halbritter, J. H. Ng, A. Y. Weiße, F. C. Wong, A. Gagliardi, N. P. Mullin, N. Festuccia, et al. Reduced oct4 expression directs a robust pluripotent state with distinct signaling activity and increased enhancer occupancy by oct4 and nanog. *Cell Stem Cell*, 12(5):531–545, 2013.

[108] T. Kalmar, C. Lim, P. Hayward, S. Muñoz-Descalzo, J. Nichols, J. Garcia-Ojalvo, and A. M. Arias. Regulated fluctuations in nanog expression mediate cell fate decisions in embryonic stem cells. *PLoS Biol.*, 7(7):e1000149, 2009.

[109] H. Niwa, J. Miyazaki, A. Smith, et al. Quantitative expression of oct-3/4 defines differentiation, dedifferentiation or self-renewal of es cells. *Nature genetics*, 24(4):372–376, 2000.

[110] K. Mitsui, Y. Tokuzawa, H. Itoh, K. Segawa, M. Murakami, K. Takahashi, M. Maruyama, M. Maeda, and S. Yamanaka. The homeoprotein Nanog is required for maintenance of pluripotency in mouse epiblast and ES cells. *Cell*, 113(5):631–642, 2003. ISSN 0092-8674.

[111] J. Chew, Y. Loh, W. Zhang, X. Chen, W. Tam, L. Yeap, P. Li, Y. Ang, B. Lim, P. Robson, et al. Reciprocal transcriptional regulation of pou5f1 and sox2 via the oct4/sox2 complex in embryonic stem cells. *Molecular and cellular biology*, 25(14):6031–6046, 2005.

[112] D. J. Rodda, J.-L. Chew, L.-H. Lim, Y.-H. Loh, B. Wang, H.-H. Ng, and P. Robson. Transcriptional regulation of nanog by oct4 and sox2. *Journal of Biological Chemistry*, 280(26):24731–24737, 2005.

[113] M. Bosnali, B. Münst, M. Thier, and F. Edenhofer. Deciphering the stem cell machinery as a basis for understanding the molecular mechanism underlying reprogramming. *Cellular and Molecular Life Sciences*, 66(21):3403–3420, 2009.

[114] N. Ivanova, R. Dobrin, R. Lu, I. Kotenko, J. Levorse, C. DeCoste, X. Schafer, Y. Lun, and I. Lemischka. Dissecting self-renewal in stem cells with RNA interference. *Nature*, 442(7102):533–538, 2006. ISSN 0028-0836.

[115] Y. Loh, Q. Wu, J. Chew, V. Vega, W. Zhang, X. Chen, G. Bourque, J. George, B. Leong, J. Liu, et al. The Oct4 and Nanog transcription network regulates pluripotency in mouse embryonic stem cells. *Nature genetics*, 38(4):431–440, 2006.

[116] A. Sharov, S. Masui, L. Sharova, Y. Piao, K. Aiba, R. Matoba, L. Xin, H. Niwa, and M. Ko. Identification of pou5f1, sox2, and nanog downstream target genes with statistical confidence by applying a novel algorithm to time course microarray and genome-wide chromatin immunoprecipitation data. *BMC genomics*, 9(1):269, 2008.

[117] I. Chambers, J. Silva, D. Colby, J. Nichols, B. Nijmeijer, M. Robertson, J. Vrana, K. Jones, L. Grotewold, and A. Smith. Nanog safeguards pluripotency and mediates germline development. *Nature*, 450(7173):1230–1234, 2007. ISSN 0028-0836.

[118] S.-y. Hatano, M. Tada, H. Kimura, S. Yamaguchi, T. Kono, T. Nakano, H. Suemori, N. Nakatsuji, and T. Tada. Pluripotential competence of cells associated with nanog activity. *Mechanisms of development*, 122(1):67–79, 2005.

[119] A. M. Singh, T. Hamazaki, K. E. Hankowski, and N. Terada. A heterogeneous expression pattern for nanog in embryonic stem cells. *Stem cells*, 25(10):2534–2542, 2007.

[120] M. A. Canham, A. A. Sharov, M. Ko, and J. M. Brickman. Functional heterogeneity of embryonic stem cells revealed through translational amplification of an early endodermal transcript. *PLoS Biol*, 8(5):e1000379, 2010.

[121] Y. Toyooka, D. Shimosato, K. Murakami, K. Takahashi, and H. Niwa. Identification and characterization of subpopulations in undifferentiated es cell culture. *Development*, 135(5):909–918, 2008.

[122] K. Hayashi, S. M. C. de Sousa Lopes, F. Tang, and M. A. Surani. Dynamic equilibrium and heterogeneity of mouse pluripotent stem cells with distinct functional and epigenetic states. *Cell stem cell*, 3(4):391–401, 2008.

[123] J. Trott, K. Hayashi, A. Surani, M. Babu, and A. Martinez-Arias. Dissecting ensemble networks in es cell populations reveals micro-heterogeneity underlying pluripotency. *Molecular BioSystems*, 8(3):744–752, 2012.

[124] T. S. Macfarlan, W. D. Gifford, S. Driscoll, K. Lettieri, H. M. Rowe, D. Bonanomi, A. Firth, O. Singer, D. Trono, and S. L. Pfaff. Embryonic stem cell potency fluctuates with endogenous retrovirus activity. *Nature*, 487(7405):57–63, 2012.

[125] A. M. Arias and J. M. Brickman. Gene expression heterogeneities in embryonic stem cell populations: origin and function. *Current opinion in cell biology*, 23(6):650–656, 2011.

[126] M. H. Stewart, S. C. Bendall, M. Levadoux-Martin, and M. Bhatia. Clonal tracking of hescs reveals differential contribution to functional assays. *Nature methods*, 7(11): 917–922, 2010.

[127] H. Niwa. How is pluripotency determined and maintained? *Development*, 134(4): 635–646, 2007.

[128] J. Silva and A. Smith. Capturing pluripotency. *Cell*, 132(4):532–536, 2008.

[129] J. Wang, D. Levasseur, and S. Orkin. Requirement of nanog dimerization for stem cell self-renewal and pluripotency. *Proceedings of the National Academy of Sciences*, 105 (17):6326, 2008.

[130] Q. Ying, J. Wray, J. Nichols, L. Batlle-Morera, B. Doble, J. Woodgett, P. Cohen, and A. Smith. The ground state of embryonic stem cell self-renewal. *Nature*, 453(7194): 519–523, 2008. ISSN 0028-0836.

[131] R. Lu, F. Markowetz, R. D. Unwin, J. T. Leek, E. M. Airoldi, B. D. MacArthur, A. Lachmann, R. Rozov, A. Ma?ayan, L. A. Boyer, et al. Systems-level dynamic analyses of fate change in murine embryonic stem cells. *Nature*, 462(7271):358–362, 2009.

[132] J. Tyson, K. Chen, and B. Novak. Sniffers, buzzers, toggles and blinkers: dynamics of regulatory and signaling pathways in the cell. *Current Opinion in Cell Biology*, 15(2): 221–231, 2003. ISSN 0955-0674.

[133] W. K. Smits, O. P. Kuipers, and J.-W. Veening. Phenotypic variation in bacteria: the role of feedback regulation. *Nature Reviews Microbiology*, 4(4):259–271, 2006.

[134] J. E. Ferrell. Self-perpetuating states in signal transduction: positive feedback, double-negative feedback and bistability. *Current opinion in cell biology*, 14(2):140–148, 2002.

[135] A. Becskei, B. Séraphin, and L. Serrano. Positive feedback in eukaryotic gene networks: cell differentiation by graded to binary response conversion. *EMBO J.*, 20(10):2528–2535, 2001.

[136] B. MacArthur, A. Ma'ayan, and I. Lemischka. Systems biology of stem cell fate and cellular reprogramming. *Nature Reviews Molecular Cell Biology*, 10(10):672–681, 2009. ISSN 1471-0072.

[137] B. MacArthur, C. Please, and R. Oreffo. Stochasticity and the molecular mechanisms of induced pluripotency. *PLoS One*, 3(8):3086, 2008. ISSN 1932-6203.

[138] B. D. MacArthur, A. Ma?ayan, and I. Lemischka. Toward stem cell systems biology: from molecules to networks and landscapes. In *Cold Spring Harbor symposia on quantitative biology*, pages sqb–2008. Cold Spring Harbor Laboratory Press, 2008.

[139] I. Glauche, M. Herberg, and I. Roeder. Nanog variability and pluripotency regulation of embryonic stem cells-insights from a mathematical model analysis. *PloS one*, 5(6): e11238, 2010.

[140] M. Tigges, T. T. Marquez-Lago, J. Stelling, and M. Fussenegger. A tunable synthetic mammalian oscillator. *Nature*, 457(7227):309–312, 2009.

[141] M. Elowitz and S. Leibler. A synthetic oscillatory network of transcriptional regulators. *Nature*, 403(6767):335–338, 2000. ISSN 0028-0836.

[142] J. Paulsson. Summing up the noise in gene networks. *Nature*, 427(6973):415–418, 2004.

[143] M. Chalfie, Y. Tu, G. Euskirchen, W. W. Ward, and D. C. Prasher. Green fluorescent protein as a marker for gene expression. *Science*, 263(5148):802–805, 1994.

[144] Y. Miyanari and M.-E. Torres-Padilla. Control of ground-state pluripotency by allelic regulation of nanog. *Nature*, 483(7390):470–473, 2012.

[145] I. Chambers, D. Colby, M. Robertson, J. Nichols, S. Lee, S. Tweedie, and A. Smith. Functional expression cloning of Nanog, a pluripotency sustaining factor in embryonic stem cells. *Cell*, 113(5):643–655, 2003. ISSN 0092-8674.

[146] H. Shapiro. *Practical Flow Cytometry*. Wiley, 2005.

[147] M. R. Soboleski, J. Oaks, and W. P. Halford. Green fluorescent protein is a quantitative reporter of gene expression in individual eukaryotic cells. *The FASEB journal*, 19(3): 440–442, 2005.

[148] P. Ray and S. S. Gambhir. Noninvasive imaging of molecular events with bioluminescent reporter genes in living subjects. *Reporter Genes: A Practical Guide*, pages 131–144, 2007.

[149] D. A. Faddah, H. Wang, A. W. Cheng, Y. Katz, Y. Buganim, and R. Jaenisch. Single-cell analysis reveals that expression of nanog is biallelic and equally variable as that of other pluripotency factors in mouse escs. *Cell Stem Cell*, 13(1):23–29, 2013.

[150] G. Dooner, G. Colvin, M. Dooner, K. Johnson, and P. Quesenberry. Gene expression fluctuations in murine hematopoietic stem cells with cell cycle progression. *Journal of cellular physiology*, 214(3):786–795, 2008.

[151] H. Habibian, S. Peters, C. Hsieh, J. Wuu, K. Vergilis, C. Grimaldi, J. Reilly, J. Carlson, A. Frimberger, F. Stewart, et al. The fluctuating phenotype of the lymphohematopoietic stem cell with cell cycle transit. *The Journal of experimental medicine*, 188(2):393–398, 1998.

[152] J. Lambert, M. Liu, G. Colvin, M. Dooner, C. McAuliffe, P. Becker, B. Forget, S. Weissman, and P. Quesenberry. Marrow stem cells shift gene expression and engraftment phenotype with cell cycle transit. *The Journal of experimental medicine*, 197(11):1563–1572, 2003.

[153] E. Davidson. Gene Regulation: gene control network in development. *Annual Review of Biophysics and Biomolecular Structure*, 36:191–212, 2007. ISSN 1056-8700.

[154] E. Segal, M. Shapira, A. Regev, D. Pe'er, D. Botstein, D. Koller, and N. Friedman. Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. *Nature Genetics*, 34(2):166–176, 2003.

[155] R. Thomas. On the relation between the logical structure of systems and their ability to generate multiple steady states or sustained oscillations. *Springer Series Syne.*, 9:180–193, 1981.

[156] W. Xiong and J. E. Ferrell. A positive-feedback-based bistable 'memory module' that governs a cell fate decision. *Nature*, 426(6965):460–465, 2003.

[157] P. Zhang, G. Behre, J. Pan, A. Iwama, N. Wara-Aswapati, H. Radomska, P. Auron, D. Tenen, and Z. Sun. Negative cross-talk between hematopoietic regulators: GATA

proteins repress PU. 1. *Proceedings of the National Academy of Sciences*, 96(15):8705, 1999.

[158] M. Laurent and N. Kellershohn. Multistability: a major means of differentiation and evolution in biological systems. *Trends in Biochemical Sciences*, 24(11):418–422, 1999.

[159] L. Niswander, S. Jeffrey, G. Martin, and C. Tickle. A positive feedback loop coordinates growth and patterning in the vertebrate limb. *Nature*, 371(6498):609–612, 1994.

[160] T. Gardner, C. Cantor, and J. Collins. Construction of a genetic toggle switch in Escherichia coli. *Nature*, 403(6767):339–342, 2000. ISSN 0028-0836.

[161] M. Savageau. Comparison of classical and autogenous systems of regulation in inducible operons. *Nature*, 252(5484):546–549, 1974.

[162] U. Heiden. Delays in physiological systems. *Journal of Mathematical Biology*, 8(4): 345–364, 1979. ISSN 0303-6812.

[163] S. Mangan and U. Alon. Structure and function of the feed-forward loop network motif. *Proceedings of the National Academy of Sciences*, 100(21):11980, 2003.

[164] N. Iranfar, D. Fuller, and W. Loomis. Transcriptional regulation of post-aggregation genes in dictyostelium by a feed-forward loop involving gbf and lagc. *Developmental Biology*, 290(2):460–469, 2006.

[165] R. Milo, S. Itzkovitz, N. Kashtan, R. Levitt, S. Shen-Orr, I. Ayzenshtat, M. Sheffer, and U. Alon. Superfamilies of evolved and designed networks. *Science*, 303(5663):1538, 2004.

[166] U. Alon. Network motifs: theory and experimental approaches. *Nature Reviews Genetics*, 8(6):450–461, 2007.

[167] P. Ingram, M. Stumpf, and J. Stark. Network motifs: structure does not determine function. *BMC Genomics*, 7(1):108, 2006.

[168] M. Kaufman, C. Soule, and R. Thomas. A new necessary condition on interaction graphs for multistationarity. *Journal of Theoretical Biology*, 248(4):675–685, 2007.

[169] D. Dubnau and R. Losick. Bistability in bacteria. *Molecular microbiology*, 61(3):564–572, 2006.

[170] D. Longo and J. Hasty. Dynamics of single-cell gene expression. *Molecular systems biology*, 2(1), 2006.

[171] A. Raj and A. Van Oudenaarden. Nature, nurture, or chance: stochastic gene expression and its consequences. *Cell*, 135(2):216–226, 2008.

[172] F. Harary and B. Manvel. On the number of cycles in a graph. *Mathematica Slovaca*, 21(1):55–63, 1971.

[173] F. Harary, R. Norman, and D. Cartwright. *Structural models: an introduction to the theory of directed graphs*. Wiley, 1965.

[174] E. Estrada and J. A. Rodríguez-Velázquez. Subgraph centrality in complex networks. *Phys. Rev. E*, 71:056103, May 2005.

[175] E. Estrada and N. Hatano. Returnability in complex directed networks (digraphs). *Linear Algebra and Its Applications*, 430(8):1886–1896, 2009.

[176] Y. A. Kuznetsov. *Elements of applied bifurcation theory*, volume 112. Springer Science & Business Media, 2013.

[177] A. Bar-Even, J. Paulsson, N. Maheshri, M. Carmi, O. E., Y. Pilpel, and N. Barkai. Noise in protein expression scales with natural protein abundance. *Nature Genetics*, 38 (6):636–643, 2006.

[178] J. W. Veening, E. J. Stewart, T. W. Berngruber, F. Taddei, O. P. Kuipers, and L. W. Hamoen. Bet-hedging and epigenetic inheritance in bacterial development. *Proceedings of the National Academy of Sciences USA*, 105:4393–4398, 2008.

[179] J. W. Veening, W. K. Smits, and O. P. Kuipers. Bistability, epigenetics, and bet-hedging in bacteria. *Annual Review of Microbiology*, 62:193–210, 2008.

[180] T. Graf and M. Stadtfeld. Heterogeneity of embryonic and adult stem cells. *Cell Stem Cell*, 3(5):480–483, 2008.

[181] P. Cahan and G. Q. Daley. Origins and implications of pluripotent stem cell variability and heterogeneity. *Nature Reviews Molecular Cell Biology*, 14:357–368, 2013.

[182] N. Barker, J. H. van Es, J. Kuipers, P. Kujala, M. van den Born, M. Cozijnsen, A. Haegebarth, J. Korving, H. Begthel, P. J. Peters, et al. Identification of stem

cells in small intestine and colon by marker gene lgr5. *Nature*, 449(7165):1003–1007, 2007.

[183] J. Lei, S. Levin, and Q. Nie. Mathematical model of adult stem cell regeneration with cross-talk between genetic and epigenetic regulation. *Proceedings of the National Academy of Sciences*, 111(10):E880–E887, 2014.

[184] B. MacArthur and I. R. Lemischka. Statistical mechanics of pluripotency. *Cell*, 154(3): 484–489, 2013.

[185] J. Garcia-Ojalvo and A. M. Arias. Towards a statistical mechanics of cell fate decisions. *Current Opinion in Genetics & Development*, 22(6):619–626, 2012.

[186] M. Osawa, K. I. Hanada, H. Hamada, and H. Nakauchi. Long-term lymphohematopoietic reconstitution by a single cd34-low/negative hematopoietic stem cell. *Science*, 273 (5272):242–245, 1996.

[187] Z. Wang and R. Jaenisch. At most three es cells contribute to the somatic lineages of chimeric mice and of mice produced by es-tetraploid complementation. *Developmental Biology*, 275(1):192–201, 2004.

[188] C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:623, 1948.

[189] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

[190] C. Gardiner. *Stochastic methods*. Springer Berlin, 2009.

[191] E. T. Jaynes. Information theory and statistical mechanics. *Physical Review*, 106(4): 620, 1957.

[192] S. D. Pietra, V. D. Pietra, and J. Lafferty. Inducing features of random fields. *IEEE transactions on pattern analysis and machine intelligence*, 19(4):380–393, 1997.

[193] D. J. MacKay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.

[194] E. L. Lehmann, G. Casella, and G. Casella. *Theory of point estimation*. Wadsworth & Brooks/Cole Advanced Books & Software, 1991.

[195] P. B. Gupta, C. M. Fillmore, G. Jiang, S. D. Shapira, K. Tao, C. Kuperwasser, and E. S. Lander. Stochastic state transitions give rise to phenotypic equilibrium in populations of cancer cells. *Cell*, 146(4):633–644, 2011.

[196] H. Chang, M. Hemberg, M. Barahona, D. Ingber, and S. Huang. Transcriptome-wide noise controls lineage choice in mammalian progenitor cells. *Nature*, 453(7194):544–547, 2008. ISSN 0028-0836.

[197] M. Van De Rijn, S. Heimfeld, G. J. Spangrude, and I. L. Weissman. Mouse hematopoietic stem-cell antigen sca-1 is a member of the ly-6 antigen family. *Proceedings of the National Academy of Sciences*, 86(12):4634–4638, 1989.

[198] S. Tsai, S. Bartelmez, E. Sitnicka, and S. Collins. Lymphohematopoietic progenitors immortalized by a retroviral vector harboring a dominant-negative retinoic acid receptor can recapitulate lymphoid, myeloid, and erythroid development. *Genes & Development*, 8(23):2831–2841, 1994.

[199] C. Holmes and W. L. Stanford. Concise review: stem cell antigen-1: expression, function, and enigma. *Stem cells*, 25(6):1339–1347, 2007.

[200] S. H. Orkin and L. I. Zon. Hematopoiesis: an evolving paradigm for stem cell biology. *Cell*, 132(4):631–644, 2008.

[201] S. Massberg, P. Schaerli, I. Knezevic-Maramica, M. Köllnberger, N. Tubo, E. A. Moseman, I. V. Huff, T. Junt, A. J. Wagers, I. B. Mazo, et al. Immunosurveillance by hematopoietic progenitor cells trafficking through blood, lymph, and peripheral tissues. *Cell*, 131(5):994–1008, 2007.

[202] S. Biechonski and M. Milyavsky. Differences between human and rodent dna-damage response in hematopoietic stem cells: at the crossroads of self-renewal, aging and leukemogenesis. *Translational Cancer Research*, 2(5):372–383, 2013.

[203] A. Nebel, E. Schaffitzel, and M. Hertweck. Aging at the interface of stem cell renewal, apoptosis, senescence, and cancer. *Science's SAGE KE*, 2006(9):pe14, 2006.

[204] S. H. Orkin. Diversification of haematopoietic stem cells to specific lineages. *Nature Reviews Genetics*, 1(1):57–64, 2000.

[205] T. Reya, S. J. Morrison, M. F. Clarke, and I. L. Weissman. Stem cells, cancer, and cancer stem cells. *Nature*, 414(6859):105–111, 2001.

[206] D. Yu, D. Allman, M. H. Goldschmidt, M. L. Atchison, J. G. Monroe, and A. Thomas-Tikhonenko. Oscillation between b-lymphoid and myeloid lineages in myc-induced hematopoietic tumors following spontaneous silencing/reactivation of the ebf/pax5 pathway. *Blood*, 101(5):1950–1955, 2003.

[207] D. D. Chaplin. Overview of the immune response. *Journal of Allergy and Clinical Immunology*, 125(2):S3–S23, 2010.

[208] D. D. Chaplin. 1. overview of the immune response. *Journal of Allergy and Clinical Immunology*, 111(2):S442–S459, 2003.

[209] S. Nelson. Role of granulocyte colony-stimulating factor in the immune response to acute bacterial infection in the nonneutropenic host: an overview. *Clinical infectious diseases*, 18(Supplement 2):S197–S204, 1994.

[210] K. R. Machlus and J. E. Italiano. The incredible journey: From megakaryocyte development to platelet formation. *The Journal of cell biology*, 201(6):785–796, 2013.

[211] P. J. Murray and T. A. Wynn. Protective and pathogenic functions of macrophage subsets. *Nature reviews immunology*, 11(11):723–737, 2011.

[212] K. Amin. The role of mast cells in allergic inflammation. *Respiratory medicine*, 106(1):9–14, 2012.

[213] G. Lippi, G. Cervellin, E. Favaloro, and M. Plebani. *In Vitro and In Vivo Hemolysis: An Unresolved Dispute in Laboratory Medicine*. Patient Safety. De Gruyter, 2012.

[214] O. Leo, A. Cunningham, and P. L. Stern. Vaccine immunology. *Perspectives in Vaccinology*, 1(1):25–59, 2011.

[215] R. G. Ramsay and T. J. Gonda. Myb function in normal and cancer cells. *Nature Reviews Cancer*, 8(7):523–534, 2008.

[216] H. Oguro and A. Iwama. Life and death in hematopoietic stem cells. *Current opinion in immunology*, 19(5):503–509, 2007.

[217] A. L. Pecora. Progress in clinical application of use of progenitor cells expanded with hematopoietic growth factors. *Current opinion in hematology*, 8(3):142–148, 2001.

[218] A. C. Drake, M. Khoury, I. Leskov, B. P. Iliopoulou, M. Fragoso, H. Lodish, and J. Chen. Human cd34+ cd133+ hematopoietic stem cells cultured with growth factors

including angptl5 efficiently engraft adult nod-scid il2rgamma-/-(nsg) mice. *PloS one*, 6(4):e18382, 2011.

[219] G. Sauvageau, N. N. Iscove, and R. K. Humphries. In vitro and in vivo expansion of hematopoietic stem cells. *Oncogene*, 23(43):7223–7232, 2004.

[220] C. C. Zhang and H. F. Lodish. Cytokines regulating hematopoietic stem cell function. *Current opinion in hematology*, 15(4):307, 2008.

[221] A. Bamezai, D. Palliser, A. Berezovskaya, J. McGrew, K. Higgins, E. Lacy, and K. L. Rock. Regulated expression of ly-6a. 2 is important for t cell development. *The Journal of Immunology*, 154(9):4233–4239, 1995.

[222] S. B. Bradfute, T. A. Graubert, and M. A. Goodell. Roles of sca-1 in hematopoietic stem/progenitor cell function. *Experimental hematology*, 33(7):836–843, 2005.

[223] C. Y. Ito, C. Y. J. Li, A. Bernstein, J. E. Dick, and W. L. Stanford. Hematopoietic stem cell and progenitor defects in sca-1/ly-6a–null mice. *Blood*, 101(2):517–523, 2003.

[224] D. R. Sisan, M. Halter, J. B. Hubbard, and A. L. Plant. Predicting rates of cell state change caused by stochastic fluctuations using a data-driven landscape model. *Proceedings of the National Academy of Sciences*, 109(47):19262–19267, 2012.

[225] P. Hänggi, P. Talkner, and M. Borkovec. Reaction-rate theory: fifty years after kramers. *Reviews of Modern Physics*, 62(2):251, 1990.

[226] P. A. Markowich and C. Villani. On the trend to equilibrium for the fokker-planck equation: an interplay between physics and functional analysis. *Matematica Contemporanea*, 19:1–29, 2000.

[227] R. Jordan, D. Kinderlehrer, and F. Otto. The variational formulation of the fokker–planck equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998.

[228] P. Bak, C. Tang, and K. Wiesenfeld. Self-organized criticality: An explanation of the 1/f noise. *Physical review letters*, 59(4):381, 1987.

[229] P. Bak and M. Paczuski. Complexity, contingency, and criticality. *Proceedings of the National Academy of Sciences*, 92(15):6689–6696, 1995.

[230] B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walter. Molecular biology of the cell (garland science, new york, 2002). *Bioinorganic Chemistry: Inorganic Elements in the Chemistry of Life*, 1997.

[231] D. Grün, L. Kester, and A. van Oudenaarden. Validation of noise models for single-cell transcriptomics. *Nature Methods*, 11(6):637–640, 2014.

[232] R. M. Kumar, P. Cahan, A. K. Shalek, R. Satija, A. J. DaleyKeyser, H. Li, J. Zhang, K. Pardee, D. Gennert, J. J. Trombetta, et al. Deconstructing transcriptional heterogeneity in pluripotent stem cells. *Nature*, 516(7529):56–61, 2014.

[233] G. Guo, L. Pinello, X. Han, S. Lai, L. Shen, T.-W. Lin, K. Zou, G.-C. Yuan, and S. H. Orkin. Serum-based culture conditions provoke gene expression variability in mouse embryonic stem cells as revealed by single-cell analysis. *Cell Reports*, 2016.

[234] T. Kobayashi, H. Mizuno, I. Imayoshi, C. Furusawa, K. Shirahige, and R. Kageyama. The cyclic gene hes1 contributes to diverse differentiation responses of embryonic stem cells. *Genes & development*, 23(16):1870–1875, 2009.

[235] C. Pina, C. Fugazza, A. J. Tipping, J. Brown, S. Soneji, J. Teles, C. Peterson, and T. Enver. Inferring rules of lineage commitment in haematopoiesis. *Nature Cell Biology*, 14(3):287–294, 2012.

[236] P. Rué and A. Martinez Arias. Cell dynamics and gene expression control in tissue homeostasis and development. *Molecular Systems Biology*, 11(2), 2015.

[237] A. Filipczyk, K. Gkatzis, J. Fu, P. S. Hoppe, H. Lickert, K. Anastassiadis, and T. Schroeder. Biallelic expression of nanog protein in mouse embryonic stem cells. *Cell stem cell*, 13(1):12–13, 2013.

[238] E. Abranches, E. Bekman, and D. Henrique. Generation and characterization of a novel mouse embryonic stem cell line with a dynamic reporter of nanog expression. *PLoS One*, 8(3):e59928, 2013.