# Quantitative single-cell imaging reveals insulation of morphogenic signal transduction

Approved by Supervisory Committee

<u>                                          </u>

Lani F. Wu, Ph.D.

<u>                                          </u>

Steven J. Altschuler, Ph.D.

<u>                                          </u>

James F. Amatruda, M.D., Ph.D.

<u>                                          </u>

Rama Ranganathan, M.D., Ph.D.

<u>                                          </u>

Neal M. Alto, Ph.D.

For my family (no matter how distant our most recent common ancestor). You are my favorite people, and your strength and general awesomeness are my inspiration.

QUANTITATIVE SINGLE-CELL IMAGING REVEALS INSULATION OF MORPHOGENIC SIGNAL
TRANSDUCTION

by

Adam D. Coster

DISSERTATION

Presented to the Faculty of the Graduate School of Biomedical Sciences
The University of Texas Southwestern Medical Center at Dallas
In Partial Fulfillment of the Requirements
For the Degree of

DOCTOR OF PHILOSOPHY

The University of Texas Southwestern Medical Center
Dallas, TX

August, 2014

QUANTITATIVE SINGLE-CELL IMAGING REVEALS INSULATION OF MORPHOGENIC SIGNAL
TRANSDUCTION

Adam D. Coster, Ph.D.

The University of Texas Southwestern Medical Center at Dallas, 2014

Lani F. Wu, Ph.D.

Steven J. Altschuler, Ph.D.

How cells integrate external cues in order to make behavioral decisions is a central problem of cell biology. In development and in tissue-homeostasis, cell-fate decisions are made by the integration of multiple morphogenic signals, but how cells convert such combinations of signals into distinct behaviors is not well understood. A major complication is our incomplete knowledge of which signal properties encode the information that cells use for decision-making. A further complication is that the static networks we use to describe cellular signaling pathways are likely to be overly-complex; the true signaling network, in a given cellular context and at a particular point in time, may be much simpler. Using a rigorous and quantitative single-cell imaging approach, I find that such simplicity is present in the integration between Wnt and Transforming Growth Factor Beta (TGFB), which are key developmental pathways. Surprisingly, this insulation extends to the integration of signals within the TGFB superfamily, which are expected to compete for shared components and so interfere with one another during signal transduction. My results thus add clarity to and simplify our understanding of how cells integrate information from the Wnt and TGFB pathways, and further suggest that insulation of signal transduction may be a common feature of morphogenic pathways.

# Contents

# Prior Publications

1. *Coster*, Thorne, Wu, Altschuler. "Disentangling signaling from transcriptional crosstalk in morphogenic pathways." (*Manucript in preparation*).

2. Thorne, Wichaidit, **Coster**, Wu, Altschuler. "GSK-3 is a gatekeeper for kinase-targeted drug response." *Nature Chemical Biology.* (*in press*).

3. **Coster**, Wichaidit, Rajaram, Altschuler, Wu. "A simple image correction method for high-throughput microscopy." *Nature Methods*, 11(2014): 602.

4. Xiong, Brunson, Huh, Huang, **Coster**, Wendt, Fay, Qin. "The Role of Surface Chemistry on the Toxicity of Ag Nanoparticles." *Small* 9(2013): 2628-38.

5. Marjanovic, Chalupska, Patenode, **Coster**, Arnold, Ye, Anesi, Lu, Okun, Tkachenko, Haselkorn, Gornicki. "Recombinant yeast screen for new inhibitors of human acetyl-CoA carboxylase 2 identifies potential drugs to treat obesity." *PNAS* 107(2010): 9093-9088.

# List of Figures

# List of Tables

# List of Abbreviations

**TGFβ** Transforming Growth Factor Beta (ligand)

**ACVR1/2** Activin A receptors, types 1 and 2

**ALK** Activin receptor-Like kinase. Type I $TGFB_{sf}$ receptors

**arm** Armadillo. The *Drosophila* ortholog of β-catenin

**ATP** Adenosine tri-phospate

$B$ Mathematical symbol representing a background fluorescence layer in an image

**BMP** Bone Morphogenic Protein (ligand). A family within $TGFB_{sf}$

**BMPR1/2** BMP receptors, types 1 and 2

**bSmad** BMP-specific rSmad. Comprises Smad1/5/8

**co-IP** co-Immunopreciptiation. An IP, except that proteins different from the original target are then assayed. If a protein can co-IP with another protein, this is taken as evidence that they physically interact

**CTNNB1** Catenin, Beta 1 (β-catenin). Transcription factor activated by canonical Wnt signaling

$cv$ The coefficient of variation. Defined as standard deviation over the mean

$D$ Mathematical symbol representing the baseline detector value in a fluorescence image

**DKK1** Dikkopf-1. An extracellular Wnt antagonist

**DPP** Decapentaplegic. The *Drosophila* ortholog of BMP2/4

$EC_{50}$ Half-maximal effective concentration.

$F$ Mathematical symbol representing a foreground fluorescence layer in an image

**FZD** Frizzled. Receptor for Wnts

**G1** Growth phase 1. The phase of the cell cycle in which cells have normal diploid (2N) DNA content.

**G2** Growth phase 2. The phase of the cell cycle in which cells have twice-normal (4N) DNA content.

**GDF** Growth and Differentiation Factor (ligand). A family within $TGFB_{sf}$

**GPCR** G-protein coupled receptor. 7-pass transmembrane receptors that signal through hetero-trimeric G-proteins

$I$ Mathematical symbol representing an image

**IP** Immunopreciptiation. A standard biochemical technique in which an antibody is used to capture a target protein

**iSmad** inhibitory-Smad. Comprises Smad6/7

**JUP** Junctional Plakoglobin (γcatenin). Paralog of β-catenin

**LGR4/5** Leucine-rich repeat-containing G-protein coupled receptors 4 and 5.

RSPO1 receptors and Wnt co-receptors

**LRP5/6** Low-density lipoprotein receptor-related protein 5 and 6. Wnt ligand co-receptors

**MAD** Mothers Against DPP. The *Drosophila* ortholog of the rSmads

**MH1** MAD homology 1. DNA-binding domain of the rSmad proteins

**MH2** MAD homology 2. Protein-binding domain of the rSmad proteins

**mRNA** messenger RNA

**qPCR** quantitative PCR. Method for measuring relative RNA abundance

**ROR1/2** Receptor tyrosine kinase-like orphan receptors 1 and 2. Receptors of non-canonical Wnt5A

**rRNA** ribosomal RNA. E.g. 18S rRNA used in qPCR controls

**rSmad** receptor-Smad

**RSPO1** R-spondin 1. An extracellular Wnt co-factor

$S$ Mathematical symbol representing the shading term in a fluorescence image

**s.d.** The standard deviation of a set of values

**siRNA** small interfering RNA. Short RNAs (22-24 bases) used to reduce a target RNA abundance via the cellular RNA-interference machinery

**TCF/Lef** T-Cell Factor/Lymphoid enhancer-binding factor. β-catenin co-factors

**TGFB$_{sf}$** Transforming Growth Factor Beta superfamily

**TGFBR1-3** TGFB receptors, types 1 through 3

**tSmad** TGFβ-specific rSmad. Comprises Smad2/3

# Chapter 1

# On cellular signaling

## 1.1 Introduction

Both in the context of multi-cellular and single-celled organisms, cells are constantly challenged to stay alive and perform tasks in the face of unpredictable environmental changes. For single-cell organisms these changes can be particularly dramatic, as the external temperature, osmolarity, and other properties are outside of cellular control [1,2]. For cells within multi-cellular organisms, microenvironmental changes fluctuate much less due to controlled modification of the environment by neighboring cells. However, in order to exert control over the environment, cells must constantly communicate with one another. The messages sent from cell to cell are themselves a form of unpredictable environmental change that cells must deal with. Here, I focus on this latter problem. That is, how do cells within multi-cellular organisms accurately interpret messages sent from their neighbors?

The potential variety of cellular signals that cells face is explosively large [3], and yet cells must somehow be able to tell these signals apart. Mammalian cells must generally be able to respond to changes within a highly complex biochemical milieu that contains proteins, small molecules, and ions. Adult stem cells must be able to reliably divide and make differentiation decisions so as to recreate functional units of organs. Embryonic stem cells must be able to generate entire organisms, going from a single cell to billions that each have different functional and morphological properties. And those embryonic stem cells must perform this task with extreme accuracy, since even a small error at the early stages would be compounded through the developmental process [4].

It is amazing that cells can respond to such an unpredictable, complex, and ever-changing environment. Even more amazing is that they do so using the interactions between finite numbers of molecules, both in quantity and type, to perform computational tasks. In order for cells to be so responsive, they must first be able to recognize that the environment has changed: they must have sensors. In order for a cell to "understand" what has happened, it must convert the influx of sensory information into an internal model of its environment. Finally, cells must map that model onto a decision regarding what action to take in response. I refer to the first part of this process, the conversion of external information into an internal model, as "signal transduction" or, in short,

"signaling." The second part, the conversion of the internal model into a behavior, I refer to as "cellular decision-making."

Understanding how cells make decisions, as a consequence of environmental or pathological perturbation, is at the core of cell biology. In experimental cell biology, we purposely break the ability of a cell to accurately process information, or its ability to make a correct decision after processing that information, in order to understand the decision-making process. A cell, on the other hand, may "unintentionally" break those same processes, thus resulting in pathology. If we can understand the basis of cellular signaling and decision-making, then we can intervene to correct such pathologies. In this way, we hope that discoveries made in basic biology will eventually show utility in the clinical treatment of human patients.

Cellular signaling is difficult to study, and so the degree of uncertainty in even the best-studied systems is astonishing (as exemplified in Chapter 2). In this chapter, I outline an abstraction of the problem of cellular signaling to give some perspective on why it is so difficult to understand. This same abstract framework can be used to rigorously define cell biological problems, and thus serves as a tool for designing meaningful experiments. By approaching the problem of cellular signal processing in this way, we become more able to directly answer the most basic questions in cell biology: what signals do cells "listen to," how do they model these signals internally, and how do they use those models to make decisions?

## 1.2 Canon and crosstalk

In cell biology, the big questions that we are interested in are generally imprecise. For example, we may want to know: "How does a stem cell decide its fate?" We cannot answer such questions directly, as they are made up of an unknown number of sub-questions. Such sub-questions that must first be addressed include: What is a stem cell? What is a fate? And what does it mean for a cell to "decide"?

### Canon

In practice, therefore, we typically begin with simpler and more concrete questions, such as "What factors influence cellular response $R$?" where $R$ may be some property such as cell cycle arrest. We can then screen for mutations, growth factors, or small molecules that affect $R$, as measured using a convenient technology. Here we are already limited in the experimental design by the variety of the factors we have access to for testing against $R$. Further, and perhaps more importantly, we are limited by what aspects of $R$ our technology can measure, and by not knowing whether we should even be looking at $R$ at all. But we must start somewhere, and so we begin to collect relationships between experimental perturbations and measurements of $R$ for the biological system we care about.

Over time, and across many laboratories, we amass a library of knowledge consisting of these experimental relationships. The meaning of each relationship alone, especially at the beginning of the process, is fuzzy. In combination, however, we hope that we can begin to build a model

**Figure 1.1:** Apparent signaling complexity increases over time. **a**, Components of signaling pathways are typically discovered by genetic means or by treating cells with unknown, purified factors and observing the resulting phenotypes. **b**, These components are then organized into canonical signaling pathways based on epistasis experiments. **c**, Finally, canonical pathways are interconnected by new experiments whose results do not fit into the canonical framework. The network edges, as drawn here, may each have a different meaning and may be specific in time or to certain experimental contexts.

of how the biological system works. Unfortunately, we do not know which of the perturbation-measurement relationships are the most important, which are outright false, nor which are only true under a particular set of experimental or biological circumstances. We work these disparate pieces of data into a general model anyway, and allow that model to evolve along with our library of knowledge, and take note of exceptions to the rules of the model. If enough exceptions build up over time, models will sometimes emerge that can better explain more of the data.

In the context of cellular signaling, this gradual process typically leads to the development of so-called "canonical" signaling networks. These networks are often constructed through the use of genetic experiments and epistasis analysis, and then built upon by biochemical and other means. Resulting canonical networks are typically depicted as protein **nodes** connected by **edges** that carry some functional meaning, where that meaning might be anything from the generic "up-regulates" to something more specific, for example an edge may mean "phosphorylates residue $Y$, leading to ubiquitination and subsequent degradation."

Historical happenstance and available methodologies therefore play a large role in how we define canonical pathways. In cell and molecular biology education, we are often taught cellular signaling through these canonical pathways via the key experiments that laid their foundations. We are therefore trained from the beginning to see the pathways as non-overlapping, distinct channels of information within cells that each carry out prototypical functions.

**Crosstalk**

However, once a canonical framework is in place, and is generally accepted throughout the field, new findings must still be attached to that framework. In this way, canonical networks tend to continuously expand, eventually encroaching onto territory that once belonged to some other canonical pathway (see Fig. 1.1) [5].

**Figure 1.2:** Static networks may not represent true network behaviors. **a**, Networks collected from multiple experimental conditions may show a variety of topologies. **b**, The static network diagrams that we typically draw are maximum projects or averages across the various experimentally observed topologies.

As canonical pathways send more and more tendrils into the global network, we end up with models of cellular signaling wherein any perturbation to the system ends up reverberating throughout the entire web: everything is connected to everything. I refer to this phenomenon generally as "crosstalk." While one may wonder how we can begin to address such complexity, it is possible that the situation is less complex than we think.

An important aspect of these networks that is all too frequently ignored is **time**. The models that are sketched out in any good review are static approximations of the true signaling network. In effect, these static networks are the maximum projections of a set of networks that exist across time and across different experimental conditions. The actual network, evolving dynamically within a cell as it processes information, might not ever look like the static map (Fig. 1.2).

It is a rare experiment indeed that measures all of the network edges simultaneously, as the goal of most experiments is to flesh out a single edge or node. Therefore we do not generally know if a given edge or node exists at all times, in all systems, or if it is instead an ephemeral thing that comes and goes as needed. Indeed, experimental evidence has demonstrated that the topologies of signaling networks may not be constant, and that they may be relatively simple at any instance in time [6,7].

Aside from the missing temporal aspect in our static canonical networks and inter-networks, there is another important and oft-ignored aspect. That is, an edge is only useful for signaling if it somehow transfers **information**. The purpose of a biochemical signaling pathway is to carry information from one node to another. The fact that two nodes are connected by a link, one that perhaps indicates binding or phosphorylation, does not imply that information has been transferred. Some edges between nodes may be tangential to the signaling process being studied, or may be sending information into parallel signaling channels. This information content problem should become clear later in this chapter and in the particular case of signaling crosstalk reviewed in Chapter 2 and studied in Chapter 3.

Instead of continually adding inter-network edges, perhaps then we should carefully evaluate the edges that already exist. By testing those edges across systems, at different times, and explicitly verifying that they carry information, we may find that some of these edges carry little weight and that therefore our currently-complex view of signaling must be both simplified and made dynamic in order to reflect reality.

## 1.3    Cells as functions

Abstracting cell biological systems into mathematical or computational models can be a powerful way to learn about those systems, though the value of any given model can vary tremendously. If a model is as complicated as the system it represents, we have learned nothing; if the model is too simple, we have not captured the biology we are trying to understand. In any case, the process of modeling itself brings a rigor to and an awareness of the studied system that might not have otherwise been possible [8–10].

This section is primarily intended for classically trained biologists, like myself, who are not often exposed to this way of thinking. Those in fields that are historically more modeling-oriented, such as computer science and systems biology, will likely find the following to be familiar.

**Setting up an abstraction**

A useful abstraction when developing models of cell signaling is to think of cells as functions $f$ that convert sensory inputs $S$ into behavioral responses $R$, such that $f(S) = R$. For example, $S$ may be the concentration of an extracellular ligand, while $R$ may be the nuclear concentration of a downstream transcription factor, so that $f(S)$ models the behavior of a receptor that converts one to the other. The function, then, is typically the biological process that we wish to understand.

I borrow the term "encoding" from computer science to refer to this relationship, because it carries with it the idea that it is **information** that is being converted from one form to another. Using this language, $f$ encodes $S$ into $R$. Take human speech as an example. In speech our brains generate words that are encoded into complex temporal patterns of vocal chord tension, lung contraction, tongue movement, and so on. Those patterns in turn encode temporal changes in air pressure that propagate away from the speaker. Cells in the ear of the listener encode those pressure changes into mechanical movements, which then encode those movements into neuronal activity that the listener finally decodes into the original spoken words.

For experiments in this framework, $S$ is whatever experimental perturbation is being applied, $R$ is the experimental readout, and $f$ is the biological encoding process that we are trying to understand. The assumption in our experiments, then, is that knowledge of $S$ and $R$ are sufficient to infer $f$. When we first dive into the complete unknowns of a phenomenon, a precise inference is essentially impossible. As a consequence, we find ourselves using vague words to describe $f$, such as "recruits," "activates," or "mediates"; the function remains a mystery.

For values of $S$ and $R$ that are closely connected by some process, experiments allow us to associate more precise mechanisms with $f$, so that we can describe the function with words like "binds" or "phosphorylates." However, for values of $S$ and $R$ that are distantly connected, perhaps through other nodes such as in the case of speech described above, clear definitions of $f$ become more difficult. This is a constant struggle in the study of developmental signaling pathways, as studied in this dissertation, because many of their interesting biological effects are far removed in time from the signaling event. Activation of these pathways may therefore induce many complex

**Figure 1.3:** Cells can be abstracted as functions that take sensory inputs $S$ and yield output responses $R$. However, interpretation of this abstraction suffers from unknowns between the initial input $S$ and the final output of interest $R$ (**a**). By adding more (**b**) and more (**c**) components the model becomes complex but each link gains functional insight.

layers of signal processing before finally affecting $R$ (Fig. 1.3): the function is a composition of functions.

To precisely model mechanism in such cases, collapsing the relationship between $S$ and $R$ into an understandable function is impossible. Instead we would need to break the original function into many, where the output of each function becomes the input for the next (Fig. 1.3), and study them individually. By breaking up a general, indirect model (e.g. "Wnt blocks stem cell differentiation") into a set of more specific models with direct relationships (e.g. "Wnt binds Frizzled, resulting in increased β-catenin, which in turn binds to the Myc promoter and increases its expression"), we obtain insight about mechanistic details at the cost of simplicity.

As an additional point, for a given biological function $f$ there can be multiple inputs and outputs, many (or most) of which are unknown. And so the inputs and outputs are better thought of as lists or vectors, as in equations 1.1-1.3 (bold face, non-italics indicate a vector). In this model, the true values of **S** are the known and experimental parameters as well as any unmeasured parameters (e.g. $S_1$ to $S_3$ may represent treatment duration, treatment concentration, and ambient temperature). The values of **R** are the measured responses as well as any unmeasured cellular parameters that change in response to **S**. To complicate matters, each parameter may have completely different units (e.g. concentration versus temperature). We must inevitably approximate the true biological parameter vectors **S** and **R** by choosing a small known subset of their values, and thus can only ever obtain estimates of the true function $f$.

$$\mathbf{S} = [S_1 \cdots S_n] \tag{1.1}$$

$$\mathbf{R} = [R_1 \cdots R_m] \tag{1.2}$$

$$f(\mathbf{S}) = \mathbf{R} \tag{1.3}$$

Aside from the obvious issue that our approximate models can only include things that we know

about, cellular functions are often highly non-linear. That is, it need not be true that $f(S_{1a} + S_{1b}) = f(S_{1a}) + f(S_{1b})$, where $S_{1a}$ and $S_{1b}$ are different values for the same parameter, nor that $f(aS) = af(S)$, where $a$ is a constant. A typical dose-response curve makes a good example, since doubling the dose (by setting $a = 2$) does not necessarily double the output (i.e. $f(2S) \neq 2f(S)$). Temporal feedback makes the models even more complex, since this allows a function to eventually modify its own input.

**Making use of the abstraction**

To summarize, we can think of cells as hierarchies of functions, where the true signals $\mathbf{S}$ and responses $\mathbf{R}$ make up the nodes of a signaling network, and the functions are the mechanistic activities that connect them in time. The functions we study may take into account many parameters of which we are completely unaware, and so we only obtain estimates of $f$. Our goal in studying cell signaling is to estimate $f$ as accurately as possible, so that $f(S) \approx f(\mathbf{S})$. Finally, to be able to assign a clear functional meaning to a particular $f(S) = R$ relationship, $S$ and $R$ must be closely connected.

A common approach to modeling signaling pathways, that allows for both non-linearity and temporal feedback, is to assemble systems of ordinary differential equations for every known node and edge in the network and to explore the behavior of the system computationally. Even a small non-linear network can generate a wide variety of outcomes given different parameters or small changes in topology [5, 11–13]. Therefore, we can test the completeness of our understanding of a signaling pathway by, for example, testing the robustness of the model's output against a wide array of biologically reasonable parameters (analogous to experimental perturbations). Such approaches can uncover deficiencies in our knowledge, as the fragility or failures of the model may indicate a missing function or node. Unfortunately, it is not true that successful recapitulation of a biological behavior by a mathematical or computational model implies that we have captured the true biology with that model.

Given the potential complexity of biological signaling models, and the high likelihood that any given model is incomplete or flat-out wrong, what is the value in developing models at all? I have already noted that models give us the potential to identify deficiencies in our knowledge, and that building models forces us to rigorously define the questions we seek to address. There is an additional important benefit, which is that models may allow us to uncover underlying simplicity.

If a model is built that recapitulates a biological phenomenon, then parts of the model can be removed in order to determine the minimal set of components that could create the observed biological behaviors. Having identified these nodes one can convert a complex mechanistic model into a simple conceptual one [7, 14]. Additionally, if a behavior can be represented by a highly simplified version of the functions that cause it, we can begin to ask why the system is more complicated than it needs to be. We must be careful with such "why" questions, since biological functions are created via a random evolutionary process. However, such questions can lead to biological insights about benefits to regulation or signal processing of the more complex observed signaling network.

## 1.4    The information encoding problem

An analogous problem to the one we face when trying to understand cell signaling is the following. Going back to the example of human speech, how would alien scientists with no concept of sound think that we communicate? Perhaps they would observe us over long periods of time, and eventually correlate certain patterns of mouth movements by one human to some behavioral task performed by another. These aliens might quite reasonably infer that we encode communicated messages into mouth movements that are then decoded visually by the recipient.

Indeed, that putative encoding is a reasonable approximation of the true encoding, largely because mouth movements are part of the encoding process that is used by the full sound-based encoding. This is, of course, why humans can learn to accurately "read lips." Some of these alien scientists would eventually notice that we can still communicate in the dark, and claim that this discovery shows that the original encoding model was incorrect. Importantly, breakdown of the model in this context did not mean that it was incorrect, it simply meant that some part of the message was being carried through an additional, unobserved channel (sound).

We face the same issue when trying to identify $\mathbf{S}$ and $\mathbf{R}$ for studying cell signaling. Communication via molecules is so outside of our experience that we have no choice but to make educated guesses as to how it might work in each context. There are many approaches for converting biological data into models, but how do we know what the relevant data are? In the broad sense, the relevant data are whatever the cell "cares about." Unfortunately we neither know what signals a cell is listening to, nor into what form it encodes this information. If we don't know $\mathbf{S}$, and we don't know $\mathbf{R}$, how can we possibly determine $f$? I refer to this generally as the "encoding problem."

### 1.4.1    Identifying inputs and outputs

The first difficulty we face is the determination of $\mathbf{S}$. We frequently assume that the property of signal that the cell cares about is its concentration (as for a ligand or drug) [15, 16]. From a biochemical perspective this is a sensible guess for what is being encoded, since we understand biology as a collection of intermolecular interactions that have binding constants, may show cooperativity, and that have behaviors that fit onto Hill curves. This leads to the further expectation that the input concentration is saturable, following some form of a sigmoid curve. The assumption that cells sense absolute concentrations need not necessarily hold true, however, as various pathways can instead show fold-change detection [13, 17, 18]. Further, responses need not follow typical sigmoid curves, as they can sometimes show linear responses over a large range of concentrations [15].

For $\mathbf{R}$, we often make a similar assumption that cells encode the received signal into intracellular concentrations of some factor. Indeed, this concentration-based encoding defines morphogenic signaling, which is typically thought to convert external ligand concentrations into active transcription factor concentrations. But again the concentration property need not be the value of $R$ that encodes the sensory input.

For example, the Tumor Necrosis Factor/Nuclear Factor kappa B (TNF/NF-$\kappa$B) pathway does

show ligand concentration-dependent increases in its nuclear transcription factor accumulation, but in such a noisy way that single cells may not be able to accurately sense the absolute ligand concentration [16]. This implies either that single cells are poor signal processors or that, alternatively, the absolute ligand concentration only partially encodes the information that the cells are using. The latter may be the case, as recent work indicates that the information content of TNF concentrations is more accurately encoded into the fold-change of transcription factor activity [18].

Another alternative to encoding to or from molecule concentration is the use of temporal information, such as integration over time or oscillatory behavior [19,20]. Therefore, care should be taken with the assumption that absolute concentration is of utmost informational value to the cell. This assumption is difficult to test, however, as there may be a concentration dependence even if this is not the primary encoding that the cell uses, as is the case for TNF/NF-$\kappa$B signaling and for the analogy of mouth movements in human speech.

### 1.4.2   Cellular variability

For a given approximation of $\mathbf{S}$ and $\mathbf{R}$, is it fair to assume that this approximation is equally meaningful for all cells in the population? Most of our knowledge of signaling stems from population-based measurements of cellular responses, for example from Western blots, microarrays, and other common cellular lysate-based methods. Such methods yield averaged cellular behaviors, therefore making the implicit assumption that this average reflects individual cell behavior [21–23]. If this assumption is incorrect, such that our measured values of $R$ do not reflect any real cellular behaviors, then the properties of $f$ that we infer will be incorrect.

Indeed, many studies have shown that this assumption of cellular homogeneity is unjustified. Dramatic examples include the classic demonstration that single *Xenopus laevis* embryos have switch-like instead of graded behavior [24], the finding that various factors thought to be correlated during adipocyte differentiation were only correlated in a small subset of cells while being anti-correlated in others [25], and the discovery that population-averaged measurements were hiding the ultra-sensitivity of temporally asynchronous bacterial motors [26].

How can single cells display behaviors different from the population average? Take the trivial case: a tissue sample may include many different cell types that have quite different properties. For example, the intestinal epithelium contains highly-secretory goblet cells and highly-absorptive enterocytes, but the average behavior of these cells might be neither secretory nor absorptive. In a less trivial case, single cells within an apparently homogeneous cultured cell line can also exhibit cell-to-cell differences, even if they are derived from the same clone [27]. Such within-cell type differences could be due to asynchrony in cell cycle position [28] or to asynchrony of other phenotypic states [19,25,29,30].

Cultured "homogeneous" populations can thus show single-cell variation due to an asynchrony in temporal movement between stable phenotypic types, but they can also vary more stochastically. Randomness in cellular phenotypes might stem from asymmetry in inherited properties after stem cell division, for example [30–33]. Additionally, because genes are typically present in only two

copies, and a finite number of molecules mediate the process of transcription, it is necessarily a noisy process [34–36] that can also generate cellular variability. Note that, without careful temporal studies, variation due to asynchrony in stable states cannot be easily differentiated from that due to rapid movement between unstable states.

Outside of fully distinct cell types, why do cells display such variability? There is no clear answer to this question, though there are many potential explanations. Perhaps the use of molecules to process information is simply so inaccurate that it must generate extensive noise. Maybe cells can work around the noise we that see, so that their behaviors are more precise than they appear. Or perhaps cells have adapted to deal with such noise by either suppressing or making use of it. Indeed, in some cases variability can be useful to the population [37], as single cells may generate subpopulations with useful functions. An example is cellular differentiation in mammals, where stem cells need to choose whether to differentiate, and to which fate [30, 32, 33]. In other cases subpopulations may have resistance to toxic environmental stresses [27, 38, 39], though we should be cautious with just-so explanations for these kinds of links.

### 1.4.3   Context-dependency

Transcriptional networks and chromatin state are highly cell type- and environment-dependent. As a consequence, properties that are essential to signaling may vary between experimental systems. Examples include concentrations of cellular receptors, pathway modulators, and effectors. Such differences in cellular properties and in the microenvironments in which cells live are often collapsed into the term "cellular context."

Cellular context also includes properties of other signaling pathways, which is important because canonical pathways may be not be isolated information channels. Thus, knowing the likely extent of crosstalk is important, though general pathway interconnectedness is difficult to measure. Some types of signaling are particularly interconnected, such as for the growth factors that modulate downstream kinase cascades due to the use of higly overlapping downstream components [40]. On the other hand, other pathways may be much less interconnected, as I show in this dissertation for several key morphogenic signals (Chapter 3).

Attempts to quantify pathway interconnectedness are few and are necessarily limited by what can be practically measured, along with the issues I have already noted in this section. Computational work seems to show that signaling through one pathway can be broadly modulated by properties of the entire cell signaling network [41], while experimental work shows that non-additive inter-pathway crosstalk is a sparse phenomenon even for high-order combinations of up to five signaling pathways [3, 42].

## 1.5   Solving the encoding problem

I have painted an admittedly bleak picture of the difficulties in studying cell signaling. It should be clear from this discussion that experimental design in cell signaling is quite difficult, and inter-

pretation of experiments must be done with extreme care. But can we remove, or at least reduce, some of the difficulties discussed above?

### Obtaining meaningful $S$ and $R$

Perhaps the most problematic of the issues discussed is that of not knowing which signals **S** a cell is interested in nor which responses **R** encode those signals. One can begin to address this by considering different variations of the inputs (e.g. changing concentrations or treatment durations) and outputs (e.g. phosphorylation states, live-cell measurements of the same cell over time, fold-change in concentration, total change in concentration). Further discussion of different output measurements, in the particular context of single-cell fluorescence microscopy, can be found in Chapter 4.

Assuming that one could identify a series of potential input signals and response metrics, it is not obvious how to go about identifying which are the more accurate estimates of **S** and **R**. An interesting approach would be to perform measurements of information content between putative combinations of $S$ and $R$, for example using the mutual information metric [16]. Under the assumption that the cell is encoding signals in the most informative way possible, the $R$ and $S$ choices that maximizes mutual information can then be considered to be the best approximation of the encoding that the cell uses. This would require either high accuracy in measurement or precise knowledge of measurement error. In Chapter 3 I test multiple reagents and readouts to ensure that they have similar information content, and I discuss this approach in more detail in Chapter 4.

In an example of this approach, research groups measured the information content of transcription factor gradients across *Drosophila* embryos, with the question of whether enough information was present to specify the location of all nuclei along the embryo. While each transcription factor had low positional information content when taken alone, in combination the factors did have enough information to specify each nucleus position with high accuracy [43–45].

As with the example of human speech above, it is important to remember that low information content of a single $S$ does not necessarily mean that it is an *incorrect* encoding. It may alternatively signify an *incomplete* encoding, and that that other unmeasured signal properties need to also be considered. Importantly, incomplete encodings can be good enough for many experimental biology questions.

### Compensating for cellular variability

To address issues stemming from cellular variability, the straightforward solution is to directly measure its effects on the experimental relationship between the chosen $S$ and $R$. Measurements of $R$ can be performed on a single-cell basis, for example by microscopy (as in this dissertation) or by flow sorting. The distributions of single-cell values can then be checked for properties, such as multi-modality, that would suggest the presence of multiple cellular states. (In my work (Chapter 3), I verified that each measurement generated unimodal single-cell distributions.)

**Figure 1.4:** The presence of subpopulations in one measurement dimension does not imply subpopulations in another dimension; subpopulations are a phenotype-dependent property. (**a**) Cells show bi-modality in their total DNA content (as measured by total Hoechst fluorescence), but (**b**) not in the coefficient of variation in DNA content. These measurements are addressed fully in Chapter 4.

In the case that different cellular states do exists, each cell can be grouped into statistically distinct subpopulations by measured phenotype [27, 46]. Each subpopulation can then be tested separately to see if each has the same $f(S) = R$ relationship. For example, in Chapter 3 I test how cell cycle phase affects measurements of inter-pathway crosstalk. Further, if live-cell markers are able differentiate subpopulations, then cells can be physically sorted and experimented on separately.

Importantly, the absence of subpopulations along the dimension of the measured response $R$ does not imply that all cells are the same. Rather, it implies that we cannot claim that they are different. Conversely, the presence of subpopulations in one dimension does not imply existence of the same subpopulations with respect to other dimensions (see Fig. 1.4). In other words, the presence or absence of cellular subpopulations as measured by one readout is insufficient evidence to make a claim about whether $R$ is being distorted by the presence of subpopulations. Such a connection must be explicitly tested.

**Determining context-dependency**

Finally, how can we deal with the issue of context-dependency? First, it is important to verify that the context-dependency truly exists. As I discuss in Chapter 2 and implied in this chapter, context-dependency of biological phenomena is often inferred by the fact that different labs produce different results when asking the same questions. However, interpretation of cell signaling results are incredibly complicated, which may simply mean that the labs were not, in fact, asking the same questions. Because cells may be encoding information differently than we expect, we should take care when comparing interpretations of results obtained by different experimental methods, as each method will approximate **S** and **R** differently.

However, some (perhaps much) of context-dependency is undoubtedly real, and can be absorbed into the simple model of cells as functions. While we could allow the function to vary from context to context, it is more useful to say that $f$ does not change but that subsets of the inputs **S** and outputs **R** can vary.

To get around this parameter variation, we can first make sure that all controllable conditions are kept the same and that all measurements are the same from experiment to experiment. Thus, the experimentally-defined subset of values in **S** and **R** do not change. Experiments can then be repeated identically across multiple cell types, so that the only varying parameters are those inherent

to cell type differences. Any consistent aspects of the relationship between experimental $S$ and $R$ across diverse cell types can then be used to infer the general properties of $f$. Indeed, this approach is common in cell biology, as it is widely believed that any given cell line may have a myriad of idiosyncratic properties.

When using such a multiple cell-type approach, one may find a case where context-dependency is so dramatic that no general properties of $f$ can be uncovered. The first aspect of this problem to tackle would be to carefully ask if the cell types are truly being treated "identically." As suggested above, experiments typically use an absolute set of conditions across cell types (e.g. identical ligand or drug concentrations). But it may be the case that two cell types simply vary in sensitivity to the conditions, such that one type is effectively receiving a half-maximal dose while the other is saturated. Because we do not know which property $R$ encodes the treatment condition, it is also difficult to know if we are measuring an "identical" readout. Perhaps some of the apparent context-dependency of signaling is due to incorrectly interpreting what it means to treat different cell types identically.

Instead of relying on constant treatment concentrations derived from the literature and applying such conditions generally across experiments, another approach would be to measure dose-response and time-response curves for all cell types that are under experimentation. Conditions could then be calibrated on a cell type-specific basis so that, for example, all cell lines receive a half-maximal input concentration.

**The encoding problem is unsolved**

Part of the intention of this chapter was to make it clear that cellular signaling is an incredibly difficult phenomenon to understand, and that experimental designs are making many assumptions that are either going unnoticed or are not being made explicit. Some of these assumptions, if made explicit, might dramatically affect how we interpret our experimental results.

There is no general solution to the encoding problem but, as I have outlined here, steps can be taken to minimize its effects. Perhaps more importantly, an awareness of the assumptions allows for them to be tested in some cases or, at minimum, allows for results to be interpreted cautiously in the light of those assumptions.

## 1.6   Dissertation aims

This chapter provided an abstract foundation on the problems faced in the study of cellular signaling. In particular, I focused on our lack of clear knowledge about how cells encode signals into intracellular models, and how this lack of clarity may be leading us to unnecessarily complex signaling pathways. In this dissertation I present a case study of one such apparently-complex signaling phenomenon, that of cross-pathway integration between Wnt and Transforming Growth Factor Beta signaling, wherein I demonstrate that the interactions are simpler than is currently believed.

In Chapter 2 I review the literature on the classic developmental signaling pathways that are the

focus of my case study, and the claimed mechanisms off crosstalk between them. These are the Wnt and Transforming Growth Factor Beta pathways. I chose these signaling networks because they are highly studied, and so have well-established approximations of what cells care about both for inputs **S** and outputs **R**. Further, both Wnt and TGFB have relatively clean canonical forms that do not share any core components, and yet there is a large body of work that ties these pathways together.

In Chapter 4 I establish rigorous, quantitative methods for fluorescence microscopy image analysis that I use to study crosstalk between the Wnt and TGFB pathways. The study of crosstalk requires a large number of experimental conditions, and the literature on crosstalk generally lacks single-cell resolution. I therefore chose high-throughput immunofluorescence microscopy as my primary experimental method. Precise single-cell measurements are essential to quantifying single-cell phenotypes, and so I focus on the discussion on how experimental error can be removed or measured.

In Chapter 3 I make use of the conceptual approach to cellular signaling described in Chapter 1, the body of literature about the Wnt and TGFB signaling pathways reviewed in Chapter 2, and the quantitative methodologies established in Chapter 4, to experimentally determine the degree of Wnt/TGFB crosstalk during signal transduction. There, I demonstrate the finding that these pathways are in fact insulated from one another, thus making a case for simplicity in morphogenic signal integration.

**Reading this dissertation**

Each chapter is relatively independent, though the reader may find some points confusing without reading earlier chapters. Chapter 4, on quantitative single-cell imaging, in particular can stand alone and so I have placed it near the end so as to not disrupt the flow of the biological content of Chapters 2 and 3. The imaging chapter should be a useful guide to any biologist or analyst in need of a conceptual and practical reference for rigorous image analysis. Those readers primarily interested in the biology of Wnt and TGFB signaling crosstalk can skip Chapter 4 without significant loss of coherence.

# Chapter 2

# On Wnt and the Transforming Growth Factor Beta superfamily

## 2.1 Introduction

The signaling components of the Transforming Growth Factor Beta superfamily (TGFB$_\text{sf}$) and Wnt pathways are deeply conserved across metazoa. Further, they are often active in the same tissue compartments at the same time, yielding ample opportunity for these pathways to interact. Indeed, stem cells in many systems integrate signals from these two major pathways to make fate decisions. In the language of Chapter 1, cells must encode some properties of the external Wnt and TGFB$_\text{sf}$ ligands into internal models that can be subsequently mapped onto a cellular decision (such as differentiation). For the Wnt and TGFB$_\text{sf}$ pathways the encoded property is thought to be concentration. In other words, it is the ligand concentration that carries information; it is this aspect of the signal that the cells eventually convert into a decision. This dose-dependent encoding mechanism is what classifies the TGFB$_\text{sf}$ and Wnt ligands as classical "morphogens."

Wnt and TGFB$_\text{sf}$ have been extensively studied for decades, yet many uncertainties remain [47–49]. The uncertainties stem in part from the extensive redundancy found in these pathways: each signaling component is represented by anywhere from one to over twenty distinct gene products. This redundancy makes classical analysis by genetic ablation difficult, as many genes would have to be ablated simultaneously. Additionally, these pathways are central to mammalian development, and so to study them in the adult requires inducible genetic constructs. Further, each pathway shows extensive context-dependency in its behavioral output, so that finding general signaling principles has been no trivial task. Finally, many of the most-studied output behaviors are temporally far removed from the initial signaling events. As discussed in Chapter 1, such distant temporal connections between network nodes makes the assignment of meaningful cellular functions quite difficult.

Even less understood is how cells combine information from the Wnt and TGFB$_\text{sf}$ pathways to make fate decisions. While a body of literature exists on this topic, no context-independent

mechanisms of pathway integration have been uncovered. Is it true then that "the context, more than the proteins … is what shapes the response [49]"? Or, as I suggested in Chapter 1, are there simpler principles of crosstalk that are being hidden by overly-complex models of signaling crosstalk?

Before answering that question in Chapter 3, in this chapter I review the TGFB$_{sf}$ (Section 2.2) and Wnt pathways (Section 2.3). This review should provide the background necessary to understand how signals and responses are thought to be encoded by these pathways. I also discuss the putative mechanisms of cross-pathway signal integration so far discovered in the field (Section 2.5), and how to make sense of those results in the conceptual context of Chapter 1.

As a reminder, I narrowly define "signaling" as the process of converting an extracellular signal into an internal model of that signal, and I define "cellular decision-making" as the use of that internal model to affect a behavioral change. Bear in mind that this distinction is not used in the majority of the literature, and so much of the work presented here conflates these two processes. This is important, since I believe that this conflation has led to inaccurate inferences regarding the integration of signaling events. All together, this chapter paves the way for the primary claim of this dissertation, that the Wnt and TGFB$_{sf}$ pathways are insulated from one another during signal transduction.

### Nomenclature

I refer to genes from multiple organisms throughout this chapter. There are differences in standardized writing conventions for gene and protein names between organisms, so to minimize confusion I adopt one convention for all species. I refer to the protein products of genes using all-uppercase for symbols or initial capitals for full protein names. To refer to a protein family, I drop the alphanumeric identifier associated with individual members. For example, the Frizzled (FZD) protein family contains the member Frizzled-1 (FZD1). Finally, note that the TGFB superfamily shares its name with one of the families contained within it, the prototypical TGFβ family (as discussed below). For clarity, I add a subscript and use the Latin 'B' when referring to the TGFB superfamily (TGFB$_{sf}$); I drop the subscript and use the Greek 'β' when referring specifically to the TGFβ family.

## 2.2   TGFB superfamily signaling

### 2.2.1   Brief overview of the TGFB$_{sf}$ signaling network

TGFB$_{sf}$ represents a broad array of morphogenic signals [50] that are deeply conserved across the metazoa. Orthologs of each pathway component are found in nematodes, flies, mammals, and even the basal metazoan *Trichoplax adhaerens* [51–53]. These pathways are functionally essential to organismal development and adult tissue homeostasis, and are therefore frequently mis-regulated in cancer and other pathologies. Despite their general importance, activity of these pathways yields broad context-dependency in phenotypic outcomes [49–51, 54–57].

**Figure 2.1:** Structure of the TGFB$_{\text{sf}}$ signaling pathway. ∼30 homodimer ligands (circles) are classified into ∼20 BMPs (green), 3 TGFβs (blue), and others (outlines). BMP2/4 and TGFβ1-3 have distinct sets of receptors. Upon ligand binding, Type II receptors activate Type I receptors. In turn, active Type I receptors phosphorylate downstream Smads, specifically the TGFβ-rSmads (tSmads) or the BMP-rSmads (bSmads) (see Section 2.2.4). Upon phosphorylation the rSmads associate with Smad4, translocate to the nucleus, and bind promoters to affect transcription.

The network structure of the canonical TGFB$_{\text{sf}}$ pathway is simple enough to prompt the statement that it has been "solved, to a first approximation [49]." It includes only three primary nodes: homodimeric ligands (Section 2.2.2) that bind to heterodimeric receptors (Section 2.2.3) which, in turn, activate the downstream Smad family of transcription factors (Section 2.2.4). Figure 2.1 shows the structure of the TGFB$_{\text{sf}}$ pathway and the diversity of its components as discussed below.

Throughout this section, it is important to be aware that the degree of functional overlap between the various TGFB superfamily members is poorly established. Most studies of these pathways include only one or a few families at a time and, likely due to historical reasons, each family has been primarily studied in the context of a handful of tissues or diseases. Therefore, having a function ascribed to one TGFB$_{\text{sf}}$ member does not at all imply that other members do not have that function. For this reason I often refer to the superfamily in general even for functions that are known only to a subset of its members.

## 2.2.2  The TGFB$_{\text{sf}}$ ligands

The more than thirty diverse ligands of the TGFB superfamily are spread across several families, each having a different degree of homology within their ranks [58–60]. The families include three prototypical TGFβs and over twenty Bone Morphogenic Proteins (BMPs), as well as various Activins, Inhibins, Growth/differentiation Factors (GDFs) and others [50,54,61] (see the phylogenetic

tree in Fig. 2.2). In this review I focus on the TGFβs and BMPs. These two families are well-studied in the context of mammalian development and disease and, as I describe below, they are selective for distinct receptors as well as downstream transcription factors.

TGFB$_{sf}$ ligands are initially made as large precursor proteins. These precursors are cleaved intracellularly, allowing the non-ligand portion of the precursor, the so-called "Latency-associated peptide (LAP)," to remain non-covalently associated [62]. This interaction is inhibitory and so the resulting inactive complex is secreted into the extracellular space. Within this non-signaling complex resides the mature homodimeric TGFB$_{sf}$ ligand.

Within the mature ligand dimer, each monomer forms a "cysteine knot" composed of internal di-sulfide bridges between three cysteine-cysteine pairs. The two monomers are covalently linked by an additional intramolecular cysteine-cysteine bridge [50,63]. The active homodimer is revealed, and then able to bind to its receptors, by separation from the inhibitory complex. Experimentally, this separation can be induced by various conditions (e.g. low pH or the addition of chaotropic salts or certain proteases). In cells, the evidence suggests that separation is a mechanical process performed by integrins [54,62,64].

The three prototypical TGFβs are highly homologous, though they have small differences in their tertiary structures that allow for some divergence in affinities for binding partners. [65]. For example, TGFβ2 requires a co-receptor for binding to its cognate receptors, while TGFβ1/3 do not. Despite these differences, the functions of the three TGFβs are thought to be essentially the same. Therefore, their effects in cells are expected to be a property of *where* and *when* a member is expressed, not *which* is expressed [66]. In my own experiments, I find that TGFβ1 and TGFβ3 do indeed produce the same phenotypic effects, though I observe ∼10-fold differences in potency between these two ligands (data not shown).

The BMPs show a much larger degree of evolutionary diversity than do the TGFβs, and can be broken into several subfamilies by both homology and function (see Fig. 2.2). The BMPs have differential specificity to several receptors, and have relatively low receptor affinity in comparison to the TGFβs (nanomolar versus picomolar [59]). Along with tissue-specific expression patterns, this differential affinity may account for the putative functional specificities attributed to each BMP. Further, there are a large number of extracellular secreted proteins that can differentially antagonize BMP signaling (e.g. Noggin, Chordin, Gremlin, and Cerberus).

The differential receptor- and antagonist-binding affinities of the BMPs are frequently cited as responsible for the idiosyncratic signaling outcomes across this ligand family [51,67]. Importantly, this is suggestive that each of these ligands may then carry the same information. The context-dependency may then simply stem from differences in the effective concentrations of the ligands. There is evidence for this, since there is broad functional redundancy across the entire BMP family. For example, knockout experiments in mice show that ablation of individual BMPs yields developmental defects but only rarely lethality (BMP2/4 being the notable exceptions) [68,69].

For simplicity, in the rest of this dissertation I focus on only one of the BMP subfamilies. This family is composed of BMP2 and BMP4, orthologs of *Drosophila* Decapentaplegic (DPP) [61].

**Figure 2.2:** Phylogeny of the TGFβ and BMP pathway ligands (**a**, Section 2.2.2), receptors (**b**, Section 2.2.3), and transcription factors (**c**, Section 2.2.4). The *T.a.* prefix indicates putative *Trichoplax adhaerens* proteins (gray). DPP and MAD are *Drosophila* orthologs of BMP2/4 and Smad1/5/8. TGFβ-specific nodes (blue) and BMP2/4-specific nodes (green) are highlighted. See [53] for a deeper discussion of the phylogeny of receptors and Smads in bilateria. Distances are approximate, and each tree has a different scale (see Methods).

This family is highly studied in the context of mammalian stem cell differentiation and represents a distinct information channel from the TGFβ family also studied in this dissertation (i.e. the receptors and downstream effectors are mostly non-overlapping) .

### 2.2.3 The TGFB$_{\text{sf}}$ receptors

Receptors of morphogenic ligands must encode extracellular ligand concentrations into some intracellular property. The TGFB$_{\text{sf}}$ receptors do this by converting ligand concentration into intracellular kinase activity. Specifically, these single-pass receptors are heterodimeric serine/threonine kinases [58]. The heterodimers consist of a Type I and a Type II receptor that initially have no association with one another. The two receptor types are brought together by ligand binding, which causes a large conformational change of the receptors [70]. This conformational change is needed to bring the receptor kinase domains together, though even after binding the receptor-receptor interactions are minimal [71] implying a high dependence on the ligand for receptor activity.

Once brought together, the Type II receptor activates the Type I receptor by phosphorylation. The activated Type I receptor can then, in turn, phosphorylate the Smad transcription factors (see the pathway structure in Fig. 2.1). Receptor activity is likely modulated, to some degree, by endocytosis via clathrin-coated pits, though the functional consequences of this to signaling are not

**Table 2.1:** Sequence sources for the TGFB$_{sf}$ ligand alignments in Fig. 2.2. Note that only the TGFβ and BMP families are represented.

| Symbol | Species | | NCBI GI | Accession |
|---|---|---|---|---|
| TGFβ1 | *Homo* | *sapiens* | 63025222 | NP_000651.3 |
| TGFβ2 | *Homo* | *sapiens* | 208022653 | NP_001129071.1 |
| TGFβ3 | *Homo* | *sapiens* | 4507465 | NP_003230.1 |
| BMP1 | *Homo* | *sapiens* | 4502421 | NP_001190.1 |
| BMP2 | *Homo* | *sapiens* | 4557369 | NP_001191.1 |
| BMP3 | *Homo* | *sapiens* | 126507087 | NP_001192.2 |
| BMP3B | *Homo* | *sapiens* | 4826740 | NP_004953.1 |
| BMP4 | *Homo* | *sapiens* | 157276593 | NP_001193.2 |
| BMP5 | *Homo* | *sapiens* | 10835091 | NP_066551.1 |
| BMP6 | *Homo* | *sapiens* | 4502425 | NP_001709.1 |
| BMP7 | *Homo* | *sapiens* | 4502427 | NP_001710.1 |
| BMP8A | *Homo* | *sapiens* | 145611428 | NP_861525.2 |
| BMP8B | *Homo* | *sapiens* | 29571106 | NP_001711.2 |
| BMP10 | *Homo* | *sapiens* | 7656928 | NP_055297.1 |
| BMP15 | *Homo* | *sapiens* | 257743454 | NP_005439.2 |
| GDF11 | *Homo* | *sapiens* | 5031613 | NP_005802.1 |
| GDF2 | *Homo* | *sapiens* | 7705308 | NP_057288.1 |
| GDF5 | *Homo* | *sapiens* | 611435007 | NP_000548.2 |
| GDF6 | *Homo* | *sapiens* | 48475062 | NP_001001557.1 |
| GDF15 | *Homo* | *sapiens* | 153792495 | NP_004855.2 |
| DPP | *Drosophila* | *melanogaster* | 17137468 | NP_477311.1 |
| T.a.57057 | *Trichoplax* | *adhaerens* | 196006614 | XP_002113173.1 |
| T.a.58663 | *Trichoplax* | *adhaerens* | 196009532 | XP_002114631.1 |
| T.a.57877 | *Trichoplax* | *adhaerens* | 196008151 | XP_002113941.1 |

**Table 2.2:** Sequence sources for the TGFB$_{sf}$ receptor alignments in Fig. 2.2.

| Symbol | | Type | NCBI GI | Accession |
|---|---|---|---|---|
| TGFBR2 | | II | 67782326 | NP_001020018.1 |
| BMPR2 | | II | 15451916 | NP_001195.2 |
| ACVR2A | | II | 518828583 | NP_001265508.1 |
| ACVR2B | | II | 116734708 | NP_001097.2 |
| AMHR2 | | II | 257743467 | NP_001158162.1 |
| ACVRL1 | (ALK1) | I | 116734712 | NP_000011.2 |
| ACVR1 | (ALK2) | I | 4501895 | NP_001096.1 |
| BMPR1A | (ALK3) | I | 41349437 | NP_004320.2 |
| ACVR1B | (ALK4) | I | 4757720 | NP_004293.1 |
| TGFBR1 | (ALK5) | I | 195963412 | NP_001124388.1 |
| BMPR1B | (ALK6) | I | 4502431 | NP_001194.1 |
| ACVR1C | (ALK7) | I | 161333835 | NP_001104501.1 |

well established [50, 58, 68, 72].

There are fewer TGFB$_{sf}$ receptor types than there are distinct ligands, implying that there must be a large degree of promiscuity in receptor-ligand binding specificity. On the other hand, the heterodimeric nature of these receptors could in principle generate as many distinct signaling complexes as there are distinct TGFB$_{sf}$ ligands, though it is unlikely that all combinations are used in signaling. Indeed, both receptor types do show some degree of specificity in ligand binding [72, 73]. The five Type II receptors and the seven type I receptors are named after their prototypical ligands, though the Type I receptors are commonly referred to as Activin receptor-Like Kinases (ALKs) 1-7 (see Table 2.2).

Receptor-ligand specificity is particularly sharp between the TGFβ and BMP2/4 subfamilies that

are the focus of this dissertation: TGFβ1-3 preferentially bind to TGFBR1 (TGFβ receptor, type 1) and the Type II receptor TGFBR2 [74], while BMP2/4 preferentially bind to BMPR1A/B (BMP receptor, type 1 A/B) and the Type II receptors BMPR2 and ACVR2A/B (Activin A receptor, type 2 A/B). [61]. This specificity has an important consequence, that BMP2/4 and TGFβ1-3 signaling can be considered separate information channels at the level of receptor activation. I note however, that this separation is not complete: cross-pathway activation has been reported in the literature [75–78] and I observe it myself (Section 3.3).

The channel specificity allows for the blocking of one pathway or the other using receptor-specific inhibitors. Several small-molecule Type I receptor inhibitors have been discovered [79,80], though the molecule SB431542 [74] is probably the most widely used. SB431542 specifically inhibits the Activin- and TGFB-specific Type I receptors, and can thus be used to block TGFB signaling while leaving BMP signaling intact. The source of specificity of this binding was demonstrated by the structure of the TGFBR1 intracellular domain bound to SB431542. A single amino acid difference was predicted and then experimentally shown to be able to generate a functional but SB431542-sensitive BMPR1, or to generate a SB431542-insensitive TGFBR1 [81].

In addition to differences in receptor-ligand affinities, further receptor-ligand specificity is gained through interactions with co-receptors. A particularly important co-receptor, betaglycan (also called TGFBR3), increases binding of TGFβ1/3, is required for TGFβ2 binding [72,73], and also generally increases BMP signaling [82]. At the same time, it seems to inhibit Activin signaling [83]. Betaglycan is a large protein with ∼800 extracellular amino acids collectively containing two known TGFB binding sites [84,85]. It has no known intracellular signaling mechanism [83], and so the primary known role of betaglycan is to increase TGFB$_{sf}$ ligand-receptor binding affinity. Endoglin, a distinct co-receptor that is related to betaglycan, appears to have the opposite role: it binds to TGFβ1/3 (not TGFβ2) and acts negatively on TGFB signaling [86–88], though its effects on other TGFB$_{sf}$ members are less studied.

## 2.2.4   The Smad transcription factors

The information from ligand concentrations in morphogenic pathways must be encoded by the receptors into some property of intracellular effectors. In the case of TGFB$_{sf}$ signaling, these effectors are the Smad transcription factors, named after the orthologous *Drosphila* MAD protein (standing for Mothers Against DPP, DPP being the *Drosophila* BMP2/4 ortholog mentioned previously). Ligand concentrations are believed to be encoded into Smad nuclear concentrations and phosphorylation states.

Smads have two functional domains, the N-terminal MH1 (MAD homology 1) and the C-terminal MH2. The MH2 domain is considerably more conserved across the Smads. The conservation of MH2 is likely due to its functional importance in mediating most protein-protein interactions, though both domains have been found to interact with various transcription factors. In particular, the C-terminus of the MH2 domain is phosphorylated by TGFB$_{sf}$ receptors, which creates a binding site for Smad-Smad interactions. The MH1 domain, on the other hand, binds DNA [50,58,72].

**Table 2.3:** Sequence sources for the Smad alignments in Fig. 2.2. Abbreviations: rSmad = receptor-Smad, iSmad = inhibitory Smad, bSmad = BMP-specific rSmad, tSmad = TGFβ-specific rSmad.

| Symbol | Function | | Species | NCBI GI | Accession |
|---|---|---|---|---|---|
| SMAD1 | rSmad | (bSmad) | *H. sapiens* | 51173727 | NP_001003688.1 |
| SMAD2 | rSmad | (tSmad) | *H. sapiens* | 51173730 | NP_001003652.1 |
| SMAD3 | rSmad | (tSmad) | *H. sapiens* | 223029440 | NP_001138574.1 |
| SMAD4 | co-Smad | | *H. sapiens* | 4885457 | NP_005350.1 |
| SMAD5 | rSmad | (bSmad) | *H. sapiens* | 47778929 | NP_001001419.1 |
| SMAD6 | iSmad | | *H. sapiens* | 218749837 | NP_001136333.1 |
| SMAD7 | iSmad | | *H. sapiens* | 299890805 | NP_001177750.1 |
| SMAD8 | rSmad | (bSmad) | *H. sapiens* | 187828357 | NP_001120689.1 |
| MAD | rSmad | | *D. melanogaster* | 442625684 | NP_001259992.1 |
| T.a.50301 | unknown | | *T. adhaerens* | 196005967 | XP_002112850.1 |
| T.a.49742 | unknown | | *T. adhaerens* | 195998077 | XP_002108907.1 |
| T.a.30731 | unknown | | *T. adhaerens* | 196012704 | XP_002116214.1 |

The MH1/2 domains are connected by a linker that can be phosphorylated by several regulatory proteins including GSK3β (this enzyme is central to Wnt signaling, as discussed in Section 2.3.5) and Mitogen-activated protein kinases, which typically results in Smad degradation [89,90].

There are eight mammalian Smads, Smad1-8. (The official name of Smad8 is, unfortunately, Smad9. I use the more common and common-sense designation Smad8 in this dissertation.) This transcription factor family is deeply conserved and consists of several subfamilies (see Fig. 2.2), each subfamily having distinct functions described below. In particular, Smad2/3 are nearly identical (though Smad2 lacks a DNA binding domain), Smad1/5/8 are quite similar to one another, and the other Smads are considerably more diverse. Additionally, Smad4 and Smad1/5/8 have remarkably homologous orthologs in the basal metazoan *Trichoplax adhaerens*.

Functionally, the eight Smads fall into distinct groups (see the summary Table 2.3): the receptor-Smads (rSmads), the inhibitory-Smads (iSmads), and the only co-Smad, Smad4. In brief, the functions are as follows. The rSmads are phosphorylated by the TGFB$_{sf}$ receptors, which creates a new binding surface for interaction with the co-Smad. It is then the rSmads, in conjunction with Smad4, that mediate the downstream transcriptional functions of TGFB$_{sf}$ signaling. The iSmads, on the other hand, contain MH2 domains and are similar in length to rSmads but lack the MH1 domain. This difference is consistent with the primary function of the iSmads, which is to compete for interactions with receptors, the rSmads, Smad4, and other factors. The iSmads thus behave like dominant-negative rSmads [72].

**Mechanisms of Smad activity**

Active Smads are thought to exist as heterotrimers of two rSmads and Smad4, though the evidence for this is indirect except in rare cases [91]. While the simple signaling model typically presented in the literature describes cytosolic heterotrimerization of Smad4 and phosphorylated rSmads, followed by transport into the nucleus, several aspects of this model either lack explicit support or are likely incorrect.

First, it is unclear whether these heteromeric complexes are created in the cytosol or in the nucleus, as nuclear import of rSmads does not require the presence of Smad4 in some cases [72].

Further, the rSmads and Smad4 constantly shuttle between the nucleus and cytosol even in the absence of active signaling [92], implying that neither phosphorylation nor trimerization are pre-requisites to nuclear localization. Alternatively, if trimerization is required then it must be able to occur in the absence of active signaling. In that case it would be possible that this shuttling is allowed by transient trimerization of inactive Smads. To my knowledge this idea has not been tested, though an interpretation of my own results is suggestive that this may be the case (Section 3.3).

Nuclear-cytosolic Smad shuttling was first inferred from localization studies, from the discovery of nuclear export signals in Smad4 [93, 94], and from interactions with nuclear pore proteins [95]. Shuttling was first shown directly by the use of live-cell experiments with exogenously-expressed fluorescent Smads [92]. In those experiments, photobleaching of either the nucleus or the cytosol depleted fluorescence in both compartments over short time scales (tens of minutes), implying that the labeled Smads were constantly shuttling between compartments. Further, both the co-Smad and the rSmad had significant decreases in mobility after phosphorylation. A better model of Smad activity, then, is that activation by the $TGFB_{sf}$ receptors stabilizes the nuclear fraction after import, though the mechanism of stabilization (e.g. anchoring to nuclear proteins or reduced export rate) is still under debate.

The dynamics of $TGFB_{sf}$ signaling in cells have not been intensively explored [96], therefore it is not known how these dynamics vary across cell types, signaling subfamilies, or experimental conditions. Nor is it established what the parameters are that define the kinetics (though there have been efforts to mathematically model these pathways [57, 97]). There is broad consensus on the qualitative behavior of these pathways, however. In response to a $TGFB_{sf}$ ligand, the dimerized receptors begin phosphorylating the rSmads. Phospho-rSmads and Smad4 then accumulate in the nucleus, reaching saturation on the order of 40-60 minutes. Importantly, it is widely believed that it is the phosphorylated forms of the rSmads that are responsible for downstream transcription, and yet there is little direct evidence of this. The rate of decay of nuclear localization after the peak response seems to be highly context-dependent, and is likely regulated in large part by still-mysterious nuclear phosphatases that leave total protein levels intact in the short-term (hours) [49].

Constitutive de-phosphorylation of nuclear Smads is argued to allow for continuous re-sampling of receptor activity. In combination with endocytosis of active receptors, this constant re-sampling may explain the apparent temporal "memory" of $TGFB_{sf}$ signaling. For example, washout of the ligand shortly after treatment results in a slow decay of the Smad response while direct receptor inhbition by small molecules results in rapid decay of the Smad signal [98]. However, this model is somewhat inaccurate since extracellular antagonists (such as Noggin) can more rapidly switch off signaling than can simple removal of ligands from the media [50].

$TGFB_{sf}$ responses are also heavily regulated in the long term (hours and days) by a number of processes that change total levels of Smads, regulatory proteins, and other pathway components [99]. Importantly, the TGFβ and BMP pathways commonly upregulate their own antagonists, the iSmads, in essentially all cell lines tested, so that these iSmads can be considered conserved targets of $TGFB_{sf}$ across all or most mammalian cell types [49, 72]. (I use expression of one of the iSmads, Smad7, to

confirm TGFB$_{\text{sf}}$ pathway activity in Section 3.1.)

**Specificity of Smad activity**

TGFB$_{\text{sf}}$ signaling through the many ligands and receptors is eventually funneled through the five rSmads described above, representing a dramatic collapse in the amount of potential information carried by the large diversity of upstream components. In fact, the reduction of information content is even more dramatic, as the rSmads are further grouped into only two distinct information channels. I refer to these as the bSmads (BMP-responsive rSmads) and tSmads (TGFβ-responsive rSmads) as shown in Table 2.3. As a consequence of sending all ligand stimulation through these two pathways, the cell is effectively ignorant about which of many specific ligands has triggered an intracellular response.

The Type I TGFB$_{\text{sf}}$ receptors that phosphorylate the rSmads have high specificity for either the bSmads or tSmads. In combinatation with the receptor-ligand specificity discussed above, all ligands eventually signal through either the tSmads or the bSmads (though the separation is not perfect, as I observe in Section 3.3). For example, the TGFβ and Activin Type I receptors TGFBR1 and ACVR1B/C specifically phosphorylate Smad2/3 [73], while the BMP Type I receptors BMPR1A/B specifically phosphorylate Smad1/5/8. This dramatic reduction of extracellular information into only two channels is still a point of confusion in the literature, as this idea stands opposed to broad differences in observed consequences of TGFB$_{\text{sf}}$ signaling, especially with respect to effects on the transcriptional network.

As discussed in Chapter 1, long temporal distances between pathway stimulation and observation of phenotypic responses make it difficult to assign clear functional relationships between these events. This is a problem for studies of TGFB$_{\text{sf}}$ function, since its most highly studied outcomes are slow processes like epithelial-to-mesenchymal transition and cellular growth arrest and differentiation.

Adding to the complication, each rSmad has a different target DNA sequence, and binding affinity to DNA is low for all Smads. The rSmads then have different, but overlapping, sets of transcriptional targets that are dramatically affected by which co-factors are present in the nucleus. For transduced signals that carry the same information content, then, the nucleus may use these co-factors to decide on a completely different set of outputs [55]. There are therefore only a small number of conserved transcriptional targets for these pathways, with iSmads being perhaps the only targets present across nearly all experimental observations [49, 61].

## 2.2.5   Signaling crosstalk within the TGFB superfamily

As the preceding sections show, there is a tremendous diversity of TGFB$_{\text{sf}}$ components used throughout mammalian signaling. Importantly, it is likely that for a given cell type, in a given microenvironment, there are multiple members of each component acting at once. How cells integrate simultaneous TGFB$_{\text{sf}}$ signals has not been extensively studied, though the consensus opinion seems to be that signaling crosstalk between these pathways is mostly a product of competition for signaling components. Given that a cell can only tell the difference between two primary arms of TGFB$_{\text{sf}}$

signaling, it is worth asking what information content can exist in the combination of input signals.

Research on the problem of intra-TGFB$_{sf}$ signaling crosstalk has been primarily performed in developmental biology systems, especially in *Xenopus laevis*. In this system distinct TGFB$_{sf}$ ligands set up partially overlapping gradients in portions of the developing embryo. Experimental manipulation of these "morphogenetic fields" causes TGFB$_{sf}$ type-specific defects. For example, Activin and BMPs have opposing roles in overlapping portions of the embryos, so that an an artificial abundance of signaling through one pathway ablates signaling through the other [100]. The ablation occurred even when extracellular aspects of signaling interaction were removed by expression of constitutive receptors. It was proposed then, without explicit evidence, that the Activin/BMP cross- inhibition could be due to competition for the co-Smad, Smad4. This claim is frequently cited in the TGFB$_{sf}$ literature, though to my knowledge it has not been directly tested. In fact, my own results in Section 3.3 suggest that the bSmads and tSmads do not compete for the co-Smad.

Similarly, competition for receptors, co-receptors (e.g. betaglycan [83] and endoglin), extra-cellular antagonists, and downstream transcriptional co-factors are all cited as sources of potential crosstalk between the TGFB$_{sf}$ members. Such competition could be modified by differential binding affinities of receptors to each ligand [59,71], and thus be further amplified by differential expression of those same receptors. There are no well-accepted models of intra-TGFB$_{sf}$ signaling crosstalk besides the competition-based mechanisms above, and explicit evidence that these mechanisms are in use by cells is currently lacking.

Finally, these pathways are also expected to cross-talk at the level of transcription, and could do so by regulating many cellular components. Because of the high context-dependency of target gene experession, the most obvious method of inter-TGFB$_{sf}$ transcriptional crosstalk is through the canonical expression of the iSmads. While expression of these proteins will clearly inhibit TGFB$_{sf}$ signaling, it is not at all obvious that such a mechanism could discriminate between TGFB$_{sf}$ family members.

## 2.3    Review of Wnt signaling

### 2.3.1    Brief overview of the Wnt signaling network

Just as with TGFB$_{sf}$, Wnt signaling is morphogenic and is deeply conserved, having orthologs in the basal metazoan *Trichoplax adhaerens* [52]. Also like TGFB$_{sf}$, the Wnt pathway is essential to development, results in highly context-dependent phenotypic outcomes, and is often disregulated in disease. In most other respects, however, these two signaling pathways are quite different. Wnt shares no core components with the TGFB$_{sf}$ pathway, has many more components overall (and is thus more complex), and has different kinetics. Finally, the mechanisms of Wnt signal transduction are less understood and more contentious than are the mechanisms of TGFB$_{sf}$ signaling.

The "canonical" Wnt signaling pathway is shown schematically in Fig. 2.3. This pathway consists of diverse extracellular Wnt ligands (Section 2.3.2) that bind to FZD receptors (Section 2.3.3) and

**Figure 2.3:** A simplified structure of the canonical Wnt signaling pathway. Of 19 ligands (circles), Wnt3A is the prototypical canonical ligand (filled red circle). Wnts binds to a subset of 10 FZDs, with a generally unknown degree of specificity. In the absence of Wnt (inset gray box), β-catenin is proteosomally degraded after phosphorylation and ubiquitination by the destruction complex. This complex includes the kinase GSK3β, the E3 ubiquitin ligase βTrCP, and the scaffold Axin. In response to ligand, the destruction complex is disrupted in a Dishevelled (DVL)-dependent manner, allowing β-catenin levels to accumulate. Nuclear β-catenin binds to the co-factor TCF/Lef and together these factors modulate the transcriptional network of the cell.

subsequently cause stabilization of the transcription factor β-catenin (Section 2.3.4). As a result, concentrations of nuclear β-catenin increase and this protein can then take part in highly context-dependent changes to the cellular transcriptional program. In addition to this canonical signaling pathway, there are many "non-canonical" Wnt pathways (as many as ten! [101]) that are poorly understood. My focus in this dissertation is on canonical Wnt signaling, though I briefly review non-canonical Wnt signaling below (Section 2.3.6).

I note that the term "canonical" is falling out of favor in the Wnt field, so that "canonical Wnt pathway" is being replaced by "Wnt/β-catenin pathway." I defined the term "canonical" somewhat differently in Chapter 1 to refer to a signaling network that is relatively distinct from other networks. Indeed, such signaling insulation is one of the aspects of Wnt/β-catenin signaling that has made it easier to study than other Wnt pathways. I therefore maintain the use of "canonical" to refer to the Wnt/β-catenin pathway in this dissertation, as it carries with it the important connotation of independence. For simplicity, by the shorthand "Wnt signaling" I always refer to the canonical variant unless otherwise specified.

**Figure 2.4:** Phylogeny of the Wnt ligands (**a**, Section 2.3.2), receptors (**b**, Section 2.3.3), and transcription factors (**c**, Section 2.3.4). The *T.a.* prefix indicates putative *Trichoplax adhaerens* proteins (gray). WG (Wingless), FZ (Frizzled), and ARM (Armadillo) are *Drosophila* orthologs of mammalian Wnt1, FZDs, and β-catenin. JUP, Junctional Plakoglobin (also known as γ-catenin) is a paralog of β-catenin. CTNNB1, gene symbol for β-catenin. 4F0A.B/A, identifiers for crystal structures of *Xenopus laevis* Wnt8A and FZD8 used in the alignments. Wnt3A and Wnt5A (reds) are the prototypical canonical and non-canonical Wnt ligands, respectively. As with the TGFB$_{sf}$ pathway, note the high degree of diversity for each node when considering the basal *T. adhaerens* proteins, and that the Wnts and FZDs each fall into multiple distinct subfamilies. Distances are approximate, in arbitrary length units (see Methods).

## 2.3.2   The Wnt ligands

The mammalian Wnt ligands were initially identified with the discovery of mouse Int-1, which was later found to be homologous to the *Drosophila* protein Wingless (WG). Mice are, of course, wingless by default, and so in mammals these two gene names were concatenated into the meaningless "Wnt1." All other Wnt ligands are named similarly [102]. In mammals, there are a total of 19 Wnts that fall into multiple subfamilies by sequence homology (Fig. 2.4) [101, 103].

These small ligands (~350 amino acids) are highly cysteine-rich, with 22 cysteine residues that have conserved spacing across all Wnts. The evolutionary maintenance of these cysteines was taken to imply the formation of intramolecular cysteine bridges, which suggested that Wnt ligands should be stable proteins [104]. Further, aspects of the primary sequence, including many charged residues, were suggestive of a protein that should be highly soluble. Despite the expected stability and solubility of the Wnt ligands, they proved to be quite problematic to work with [105].

In overexpression systems, intracellular Wnts were found to have many glycosylated forms and tended to associate with chaperones, implying some difficulty in properly folding these proteins. Extracellular Wnts, in contrast, were found to have fewer glycosylated forms. Further, the bulk of overexpressed Wnt products tended to remain in the cell [102], and the Wnt proteins that left the cell tended to stay associated with membranes or with the extracellular matrix [106]. All of this data suggested that the Wnt ligands were in fact neither stable nor soluble, explaining in part the

**Table 2.4:** Sequence sources for the Wnt alignments in Fig. 2.4.

| Symbol | Species | NCBI GI | Accession |
|--------|---------|---------|-----------|
| Wnt1 | *H. sapiens* | 4885655 | NP_005421.1 |
| Wnt2 | *H. sapiens* | 4507927 | NP_003382.1 |
| Wnt2B | *H. sapiens* | 630044901 | NP_001278809.1 |
| Wnt3 | *H. sapiens* | 13540477 | NP_110380.1 |
| Wnt3A | *H. sapiens* | 14916475 | NP_149122.1 |
| Wnt4 | *H. sapiens* | 17402922 | NP_110388.2 |
| Wnt5A | *H. sapiens* | 371502087 | NP_001243034.1 |
| Wnt5B | *H. sapiens* | 17402919 | NP_110402.2 |
| Wnt6 | *H. sapiens* | 16507239 | NP_006513.1 |
| Wnt7A | *H. sapiens* | 17505191 | NP_004616.2 |
| Wnt7B | *H. sapiens* | 17505193 | NP_478679.1 |
| Wnt8A | *H. sapiens* | 17505195 | NP_490645.1 |
| Wnt8B | *H. sapiens* | 110735437 | NP_003384.2 |
| Wnt9A | *H. sapiens* | 15082261 | NP_003386.1 |
| Wnt9B | *H. sapiens* | 17017976 | NP_003387.1 |
| Wnt10A | *H. sapiens* | 16936520 | NP_079492.2 |
| Wnt10B | *H. sapiens* | 16936522 | NP_003385.2 |
| Wnt11 | *H. sapiens* | 17017974 | NP_004617.2 |
| Wnt16 | *H. sapiens* | 17402914 | NP_057171.2 |
| WG | *D. melanogaster* | 17648113 | NP_523502.1 |
| T.a.30370 | *T. adhaerens* | 196012489 | XP_002116107.1 |
| T.a.52489 | *T. adhaerens* | 195996709 | XP_002108223.1 |

difficulty in purifying these proteins.

With the first successful purification of a functional Wnt in 2003, the explanation for the lack of solubility became clear: Wnts have an absolutely conserved palmitoylated cysteine at the N-terminus. This lipidation is essential, as mutant proteins lacking the lipidated cysteine are incapable of signaling [107]. The relative insolubility of Wnt ligands have made them quite difficult to purify. Indeed, purified Wnts are commercially available primarily through a single vendor (R&D Biosystems), and at relatively low purity. As I show later in Section 3.1 (Fig. 3.6), this low purity of the common reagent has important consequences to interpretation of some experimental results.

To reduce experimental costs, studies therefore frequently use Wnt-conditioned media instead of purified Wnts. This approach suffers in that conditioned media contains a large amount of unknown secreted cellular products. Cell lines that lack the Wnt overexpression are often used as negative controls. However, given the degree of transcriptional remodeling that the Wnt pathway can cause, there are likely many unknown differences between Wnt-conditioned and control-conditioned media.

The solubility problem for the Wnt ligands led to an additional difficulty, that of determination of the tertiary structure of these proteins. Indeed, the first crystal structure of a Wnt was obtained only recently [108]. The primary sequence of Wnts are not related to any known protein fold, so prior to the crystal structure there were many unknowns regarding how Wnts bind to their various receptors and co-receptors. I touch on this more below.

### 2.3.3 The Frizzled receptors

There are 10 human Wnt receptors, the Frizzleds (FZDs). These are a diverse group of seven-pass G-coupled protein receptors (GPCRs), though canonical signaling is mediated primarily by non-

G-protein mechanisms [109]. Upon binding to Wnt and various cofactors the Wnt/FZD complex is likely endocytosed and sequestered into multivesicular bodies within the cell. This receptor internalization is essential to proper signaling [110, 111]. As a consequence of Wnt binding to the Frizzled receptors, the transcription factor β-catenin is stabilized. However, to date there are many competing models and no clear front-runner for how FZD causes this outcome (see Section 2.3.4).

Essential for Wnt signaling through FZDs are the single-pass transmembrane co-receptors, LRP5 and LRP6 (tongue-twistingly expanding to "Low-density lipoprotein receptor-related protein" 5 and 6) [103]. The ectodomain structure of LRP6 was recently solved, revealing four tandem beta-propellar-EGF-like domains that provide a large interface for interacting proteins. Prior mutational studies of these domains showed that they are essential for Wnt binding. Further, Dickkopf-1 (DKK1), a classical extracellular Wnt antagonist, was shown to bind to these same domains, which would likely occlude the Wnt binding interface [112–114]. (I use purified DKK1 in Section 3.1 to demonstrate specificity of Wnt responses.)

Additional co-factors that are not essential but that dramatically amplify Wnt signaling are the soluble protein R-spondin 1 (RSPO1) [115] and another subfamily of GPCRs, LGR4 and LGR5 (short for "Leucine-rich repeat-containing G-protein coupled receptor 4 and 5"). It was recently discovered that these two co-factors are themselves likely a ligand-receptor pair [115–117], though the mechanism by which these co-factors enhance Wnt signaling is still a mystery.

With ten receptors and multiple co-factors, we are left with the issue of signaling specificity. Wnts appear to have high promiscuity for the Frizzleds, with little certainty in the field regarding which, if any, receptor-ligand pairings are excluded. It is still unknown which Wnts bind to which Frizzleds, or if Frizzleds can generally signal in both canonical and non-canonical ways [109, 118]. Unfortunately, while Wnt3A and Wnt5A are generally considered to be the prototypical ligands for canonical and non-canonical Wnt signaling, respectively [119], there is evidence that each can transduce signals through a non-prototypical pathway under certain conditions [118]. The recent crystal structure of *Xenopus laevis* Wnt8 with FZD8 did little to enhance our understanding of the origins of receptor-ligand specificity, as the intra-protein contacts occurred primarily at highly conserved Wnt residues [108]. This is suggestive that there is either little specificity between ligands and receptors, or that specificity is a more complicated outcome of interactions with other factors.

### 2.3.4   β-catenin, the canonical Wnt effector

β-catenin is the main effector of canonical Wnt signaling. This large protein is one of several catenins that are used in cellular adhesion. β-catenin itself has a large membrane-associated pool dedicated to this task [120], while only a "vanishingly small" baseline cytosolic component is found in Wnt un-stimulated cells [47].

In the classic model of Wnt signaling, basal β-catenin is constitutively transcribed and translated, but then degraded just as quickly in the absence of Wnt. Wnt stimulation stabilizes β-catenin, thus allowing cytosolic levels to build. In some systems, cytosolic β-catenin accumulation is measureable in as little as 15 minutes, and reaches a stable peak response by 2 hours that can last longer than

**Table 2.5:** Sequence sources for the Frizzled (FZD) alignments in Fig. 2.4.

| Symbol | Species | NCBI GI | Accession |
|---|---|---|---|
| FZD1 | *H. sapiens* | 4503825 | NP_003496.1 |
| FZD2 | *H. sapiens* | 4503827 | NP_001457.1 |
| FZD3 | *H. sapiens* | 8393378 | NP_059108.1 |
| FZD4 | *H. sapiens* | 22547161 | NP_036325.2 |
| FZD5 | *H. sapiens* | 27894385 | NP_003459.2 |
| FZD6 | *H. sapiens* | 257470999 | NP_001158087.1 |
| FZD7 | *H. sapiens* | 4503833 | NP_003498.1 |
| FZD8 | *H. sapiens* | 13994190 | NP_114072.1 |
| FZD9 | *H. sapiens* | 4503835 | NP_003499.1 |
| FZD10 | *H. sapiens* | 6005762 | NP_009128.1 |
| FZ | *D. melanogaster* | 17864440 | NP_524812.1 |
| T.a.12196 | *T. adhaerens* | 196002269 | XP_002111002.1 |
| T.a.31674 | *T. adhaerens* | 196014261 | XP_002116990.1 |

24 hours [47, 48]. The nuclear import/export rates for β-catenin are unaffected by signaling, and so it appears that nuclear accumulation resulting from Wnt signaling is due to an overall change in protein abundance throughout the entire cell [103, 110].

Importantly, β-catenin appears to encode information about extracellular Wnt concentrations primarily in its nuclear concentration, as opposed to the concentration of particular phosphorylation states. While various β-catenin phospho-states have been discovered and claimed to be required for the transcriptional activity of this protein, careful quantification of post-signaling protein levels revealed that the vast majority of β-catenin is not phosphorylated in the presence of Wnt. Instead, the phospho-state ends up at similar concentrations to the basal state [48], implying that phospho-β-catenin does not encode Wnt ligand concentrations.

It is interesting to note that, like with the TGFB$_{sf}$ pathways, the large number of ligands and receptors end up bottlenecking to a small number of effectors. This bottlenecking is more dramatic in the case of Wnt, as there is only one β-catenin gene to which all canonical Wnt signaling leads. There is a highly homologous paralog, γ-catenin (also called Junctional Plakoglobin, JUP), that differs from β-catenin to a similar degree as does the functionally identical *Drosophila* ortholog, Armadillo (see Fig. 2.4). γ-catenin appears to have some functional redundancy to β-catenin, and even after deletion of both catenins some Wnt signaling has still been shown [101]. However, it appears that β-catenin is the primary mediator of canonical Wnt signaling. Interestingly, similarly to the TGFB$_{sf}$ pathways described earlier, this severe bottleneck implies that cells cannot know which Wnt or Frizzled activated the pathway, unless an additional channel of signaling carries this information.

To enact its transcriptional functions, accumulated nuclear β-catenin partners with one of a set of transcription factors in the TCF/Lef family (standing for T-Cell Factor/Lymphoid enhancer-binding factor). There are four such proteins in mammals that appear to be functionally redundant but do have different expression patterns [103]. These are TCF7, TCF7L1, TCF7L2, and LEF1. TCF1 was first discovered as a factor involved with T-cell differentiation (hence the name) [121], and later connected to Wnt signaling after finding that a *Xenopus laevis* ortholog, Xtcf-3, could cause β-catenin nuclear translocation [122]. TCF/Lef is normally bound to the protein Groucho,

which prevents association of TCF/Lef with β-catenin. This block is lifted in response to Wnt signaling after ubiquitination of Groucho by XIAP (X-linked inhibitor of apoptosis) [123]. It has been proposed that TCF/Lef binding to β-catenin is stable enough to effectively sequester it in the nucleus and so allow for long-term Wnt signaling [120], though to my knowledge this has not been explicitly tested.

In complex with TCF/Lef, nuclear β-catenin is able to bind to a diverse set of promoters and modulate transcription. These targets are highly context-dependent, though the negative auto-regulator Axin2 is considered to be a conserved target of β-catenin [103, 124]. All canonical Wnt signals are encoded into nuclear β-catenin concentrations, therefore ligand- and receptor-specific information must be lost during signal transduction. This implies that the context-dependency of Wnt signaling responses should be the result of transcriptional idiosyncrasies or of additional information channels (e.g. non-canonical Wnt pathways).

### 2.3.5   The destruction complex

The functional effect of FZD, after binding Wnt, is to relieve the otherwise constitutive degradation of β-catenin. The set of proteins that are involved with β-catenin degradation is referred to, ominously, as the "destruction complex." The mechanisms by which this complex degrades β-catenin are fairly well understood, but how active FZD causes an end to the destruction is still highly contentious [47, 48]. Here, I discuss the components of this complex that are best understood, and that will become important later in the discussion on crosstalk between Wnt and TGFB$_{sf}$ (Section 2.5).

**GSK3β**

The protein Glycogen synthase kinase 3 beta (GSK3β) maintains a central role in the destruction complex. In fact, GSK3β is a signaling hub for many pathways, making it a prime node for signal integration. GSK3β is inhibited by Insulin and growth factor signaling, which lead to phosphorylation of its N-terminal tail. The phosphorylated tail becomes a "pseudo-substrate" that competes with the actual substrates of GSK3β. However, while this mechanism of inhibition is well established, it is not clear if it is used by the Wnt pathway. Nor is it clear whether the Wnt pathway maintains a distinct pool of GSK3β from growth factor pathways that can act as an insulated information channel. If not, then modulation of a common pool of GSK3β by Wnt or other pathways would necessarily create inter-pathway crosstalk.

What is clear is that this kinase requires "priming" phosphorylation by another member of the β-catenin destruction complex, Casein Kinase 1 alpha (CK1α). This priming stems from phosphorylation of β-catenin by CK1α, which produces a recognition site for subsequent phosphorylation of β-catenin by GSK3β. Phosphorylation by GSK3β in turn creates a recognition site for the the ubiquitin E3 ligase BTRC (beta-transducin repeat-containing protein, also known as βTrCP), which subsequently ubiquitinates β-catenin and sends it off to the grinding mill of the proteosome [48,125]. The widely described model of Wnt signaling has this GSK3β-mediated phosphorylation being somehow prevented in the presence of Wnt, though exactly how is a major subject of debate. In any

case, pharmacological inhibition of GSK3β by LiCl [126, 127] or the ATP analog BIO [128] is suffi-cient to induce large increases in β-catenin levels. Inhibition of GSK3β is therefore commonly used to experimentally simulate Wnt treatment, as it gets around the earlier-mentioned difficulties of working with purified Wnt ligands. Care should be taken when interpreting results of these experi-ments, as disentangling GSK3β/β-catenin effects from the myriad of other GSK3β-mediated effects is non-trivial.

**Axin and APC**

There are two large proteins that act as scaffolds on which the drama of β-catenin destruction unfolds. These are Adenomatous Polyposis Coli (APC) and Axin. Despite general agreement that APC is involved in this process, its effects are often incongruous between experiments and systems, and so there is a lack of consensus for what APC actually does in the complex [47]. The roles of Axin are better understood, and it is generally believed that this protein forms the primary scaffold of the β-catenin degradation complex [129].

Until recently, a prevailing model was that APC, which is known to be a mobile protein, would drag the destruction complex to an unknown subcellular compartment in order for β-catenin degra-dation to occur. However, a careful study found that APC could effectively block Wnt signaling no matter its localization. This was done by expressing APC constructs, modified with various localization signals to specific compartments, in the background of APC-null cell types [129]. All variants allowed continued β-catenin destruction, which convincingly disproves the APC localization model.

There exists some structural work showing the binding of APC and Axin to one another. Intrigu-ingly, these studies reveal that Axin binds using a region with homology to Regulator of G-Protein Signaling (RGS) sequences [130]. Because FZDs are GPCRs (and GPCRs activate G-proteins), it is intriguing that Axin has an RGS domain, though to my knowledge no functional link has been demonstrated. More important is the fact that Axin and APC appear to bind to the same inter-face of β-catenin [110], implying that they should compete for β-catenin binding. This is especially important because Axin levels are widely-cited to be extremely low in cells [131], such that even a small change in its concentration (e.g. via overexpression) would then affect the balance of APC versus Axin binding to β-catenin.

**Dishevelled**

Immediately downstream of FZD, before the destruction complex, is the protein Dishevelled (DVL), named for the phenotype of the *Drosophila* ortholog. This component is absolutely required for Wnt signaling in general, including non-canonical, though its precise role is something of a mystery [109, 132]. There are three paralogs in mammals, though they have a high degree of functional redundancy and so likely differ primarily in expression pattern (DVL1 or DVL2-null mice are viable, though DVL3-null mice die of developmental heart defects [132]).

**Mechanisms of destruction**

Having completed a brief tour of the destruction complex components, I turn to the mechanisms by which Wnt signaling modulates activity of the complex. The favored models of Wnt-induced β-catenin stabilization involve disabling the destruction complex so that β-catenin phosphorylation, and subsequent degradation, is reduced. How it does so is still contentious; I review a few competing models below.

DVL has a structural PDZ domain that binds to Frizzled, and DVL and Axin both contain DIX domains. The homologous DIX domains of these two proteins have been proposed to allow them to polymerize, yielding a model wherein FZD activation pulls Axin to the membrane, through a DVL bridge, thus disrupting the destruction complex [110,132]. No definitive evidence for this model has been established, though the components can be found together in internalized endosomes during Wnt signaling events. This suggests sequestration as a possible mechanism for DVL-mediated Wnt signaling [109]. LRP6, a Wnt co-receptor, has been similarly proposed to pull Axin to the membrane. This is consistent with overexpression of either DVL or an LRP6 endodomain fragment being sufficient to induce increased β-catenin levels [110, 111]. This is the model that I show in Fig. 2.3.

Sequestration of the Wnt signaling complex into internal compartments, causing isolation from interactions with cytosolic proteins like β-catenin, is a mechanism with strong support [111]. However, there are some inconsistencies with this model. The primary mediators of the destruction complex, Axin and GSK3β, are constitutively expressed, such that they would have to be continuously sequestered to prevent new protein products from slowly causing a decay in β-catenin. Axin is eventually destabilized by Wnt signaling, which gets around this caveat, but while this occurs Axin2 levels are increased to compensate. Therefore the long-term maintenance of Wnt activity becomes hard to explain [48].

Further, the key study showing this receptor-Axin complex sequestration also found that global protein levels increased in cells after Wnt treatment, which they interpreted to be due to removal of a significant enough fraction of GSK3β to also decrease its suppression of non-Wnt proteins [111]. Given the extremely low quantities of Axin in cells [131], and the fact that its interaction with the much more abundant GSK3β is stoichiometric, it not clear how a Wnt signal could so significantly impact global GSK3β signaling by sequestration of the destruction complex-associated pool. Recent work is suggestive that the global increase in protein levels after after a Wnt response is instead due to a cell-cycle dependent form of non-canonical signaling [133].

In addition to the sequestration model, active FZD complexes could negatively affect GSK3β, Axin, the priming kinase CK1α, or some combination of these. All models have support, though the experimental bases for these models nearly universally depend on overexpression or knockdown of pathway components [47,48]. While use of such dramatic pathway modulation does not negate the results, it does complicate the interpretation of those results. Recent work has begun to make use of endogenous protein levels, which has led to two strong but incompatible models.

In the first model, the authors take an Axin-centric focus, making the assumption that all Wnt-

**Table 2.6:** Sequence sources for the Wnt-related catenin alignments in Fig. 2.4.

| Symbol | Name | Species | NCBI GI | Accession |
|---|---|---|---|---|
| CTNNB1 | β-catenin | *H. sapiens* | 148233338 | NP_001091679.1 |
| JUP | γ-catenin | *H. sapiens* | 4504811 | NP_002221.1 |
| ARM | Armadillo | *D. melanogaster* | 24639204 | NP_476665.2 |
| T.a.22780 | | *T. adhaerens* | 196000871 | XP_002110303.1 |

relevant β-catenin destruction is mediated by Axin [47]. They therefore rely on co-immunoprecipitation (co-IP) of endogenous Axin for all experiments. Using this approach, they show that endogenous Axin levels remain mostly unchanged during the early Wnt response, implying that modulation of Axin levels is not the primary mechanism of β-catenin stabilization (though this does not rule out non-degradative forms of modulation). Instead, the authors were surprised to find that Axin pulled down minimal β-catenin in the absence of Wnt, and that Wnt signaling caused more β-catenin to associate with Axin. They interpreted this to mean that interactions on this scaffold are normally transient, but that Wnt signaling blocks some step and so renders the destruction complex inert by saturating it with non-transient β-catenin binding. With the additional finding that Axin-bound βTrCP dropped rapidly after Wnt treatment, the authors concluded that Wnt activity somehow blocks ubiquitination of β-catenin instead of phosphorylation.

In opposition to the saturation model, another careful study used kinetic arguments and careful measurements of protein quantities to show that it is indeed the phosphorylation steps that control the Wnt response [48]. The authors start from the position that the current models, including the saturation model, are insufficient to explain the kinetics of the β-catenin response to Wnt. In particular, the models fail to explain the maintained high level of β-catenin after ∼2 hours. They show experimentally that the destruction complex is functional whether or not Wnt is present, which directly argues against the saturation model. Further, the authors show mathematically that a saturation-based mechanism would have a limited range for regulation. Finally, their kinetic models demonstrated that all observed dynamics of the pathway response could be explained by the degree of phosphorylation at two sites within β-catenin. As I note in Section 1.3, however, agreement with an abstract model does not imply that the model is correct.

What can we make of this general confusion in the field for how Wnt causes an increase in its effector, β-catenin? First, recent studies that focus on endogeous protein levels are clearly moving in the right direction, and I think that the contradictions between the current best studies will begin to be resolved by future use of similarly clean, endogenous approaches. For my own work, I decided that it was best to avoid favoring one model over another, since defining experiments with respect to such tenuous models would make the future of their interpretation uncertain. For this reason, I take a simple, mechanism-independent approach to Wnt signaling that only requires that extracellular Wnt concentrations are encoded as intracellular β-catenin concentrations (Chapter 3).

### 2.3.6 Non-canonical Wnt signaling

Our understanding of the best-understood Wnt signaling pathway, the canonical pathway, is still incomplete. The various non-canonical Wnt signaling pathways are even less understood. This in

large part because the readouts of these pathways are either hard to measure or are affected by other factors whose contributions are difficult to disentangle [109, 118]. Though the focus of this dissertation is on canonical signaling, here I briefly review the non-canonical pathways to provide background for an instance of inter-pathway crossalk between TGFB$_{sf}$ and non-canonical Wnt5A signaling.

There are many different non-canonical pathways, with some authors counting up to ten [101], though the true number is certainly unknown. Perhaps the best-studied non-canonical pathways are those of Planar Cell Polarity (PCP) and Convergent Extension (CE) in *Drosophila*. However, it is important to note that these pathways are studied in different systems (the wing and developing embryos, respectively) and share core components, so it is unclear how distinct these pathways truly are [118]. This ambiguity is common among the non-canonical pathways. An additional non-canonical pathway studied in mammalian systems activates $Ca^{2+}$ signaling, though the consequences of such signaling are mostly unknown [134].

While the mechanisms, readouts, and functional consequences of non-canonical signaling are poorly understood, there are several key components that are well-established. At the level of receptors, ROR2 (expanding to "Receptor tyrosine kinase-like orphan receptor") has an extracellular Wnt-binding domain homologous to that of FZD. This receptor activates $Ca^{2+}$ signaling in response to Wnt5A. The ROR2 homolog ROR1 also has this domain, but may be a pseudo-kinase [118]. Importantly, Wnt5A signaling through ROR2 is well established to inhibit long-term canonical Wnt/β-catenin signaling, though how it does so is still unknown [101, 109, 118, 119].

Downstream of the receptors, DVL is frequently (but not always) required for non-canonical signaling. Interestingly, DVL seems to use different protein domains for canonical and non-canonical signaling, so that mutation of one domain can preferentially block one type of signaling [118, 132].

Finally, there is no good reason to believe that canonical and non-canonical signaling must occur in isolation. Because of the promiscuity of Wnt-FZD binding and the unknown contribution of each FZD to the various Wnt signaling pathways, it is reasonable to expect that a Wnt signal may trigger multiple simultaneous downstream pathways. Perhaps this is a mechanism by which the cell obtains more information about the original signal, since β-catenin concentrations alone can only encode information about the overall extent of canonical Wnt activity, not which Wnts triggered that activity.

## 2.4 Functional importance of TGFB$_{sf}$ and Wnt signal integration

Having covered the canonical TGFB$_{sf}$ and Wnt signaling pathways, we can now move forward to the topic of this dissertation: inter-pathway crosstalk. The goal of this section is to both motivate the study of Wnt/TGFB$_{sf}$ crosstalk and to review the work that has been done in this field. There are many systems in which both of these pathways are intensively studied (indeed entire textbooks have been written on these pathways). However, most of this work studies each pathway indepenently; there exists a much smaller body of work on Wnt/TGFB$_{sf}$ inter-pathway crosstalk. One of the more

experimentally-tractable of such systems is the adult mammalian gut, and so I focus my review on this system.

The study of adult stem cells has advanced rapidly in recent years, especially with the identification of resident stem cells for several tissues [135] and with the discovery that differentiated cells can be turned pluripotent by expression of a handful of transcription factors [136]. However, there is still much that we do not understand about stem cell biology in general, especially because *in vitro* studies are difficult to map back to *in vivo* phenomena, and because each stem cell system has proven to have its own quirks including differential use the same genetic pathways.

In particular, work in this field has pointed toward $TGFB_{sf}$ signaling being a major factor in stem cell differentiation, while canonical Wnt signaling seems to control stem cell maintenance. There are cases where this oppositional role assignment does not occur [51, 101, 137–142], though it is fairly common across stem cell systems [105, 143–147]. These roles are particularly well-studied in the context of the mammalian intestine.

Our understanding of the intestinal stem cell system has increased significantly in recent years, due in no small part to the work of Hans Clevers and colleagues. Clevers' group has dissected the function of many genes involved with intestinal stem cell development (with a focus on canonical Wnt), has unambiguously identified the stem-cell population [148], and has developed an *in vitro* system for culturing intestinal stem cells in such a way that they behave like homoestatic *in vivo* crypts [149, 150]. This work has even led to successful engraftment of intestinal stem cells isolated from one mouse into the damaged intestine of another [151]. Given the potential value to stem cell therapy, and the connection of stem cells with colon cancer [152, 153], a deeper understanding of intestinal stem cells is highly sought after in hopes of obtaining new therapeutic strategies for many bowel diseases.

While the gut stem cell system is highly studied, it is poorly understood how intestinal stem cells integrate distinct signals from their environment, such as Wnt and $TGFB_{sf}$, and thus choose one of several differentiated outcomes. Given the oppositional roles of $TGFB_{sf}$ and Wnt in the gut, and their various functional interactions in other systems [139, 140, 154–162], it is important to understand generally how cells integrate these two signals.

### 2.4.1   $TGFB_{sf}$ and Wnt in the intestinal crypt

**Crypts are the functional unit of intestinal epithelium**

The small and large intestine both contain similar stem cell systems. Between these two organs, morphological differences in stem cell systems are revealing in terms of the distinct organ functions. The small intestine contains a high fraction of absorptive cells (enterocytes) and an increased epithelial surface area due to relatively large lumenal projections termed villi, consistent with its role in digestion and absorption of food. The colon is predominantly made up of mucus-secreting goblet cells and completely lacks villi, consistent with its need to carry potentially abrasive and toxic waste [163]. The protective role of goblet cells is highlighted by the connection between goblet cell loss and rapid tumor formation in mice [164]. The colon is of particular clinical interest due to the

high rates of cancer and inflammatory diseases of this organ [165]. Despite macroscopic differences in morphology, many aspects of of the stem cell biology seem to be similar in the small and large intestine.

A cross-section of the intestine would reveal a hollow tube whose wall consists of several tissue layers. The most-lumenal layer consists of the mucosa, a single sheet of epithelium over stromal tissue (fibroblasts, collagen, blood/lymphatic vessels, etc.), together forming crypts and their support structures. The epithelium sits atop a basement membrane [166] which in turn is immediately adjacent to myofibroblasts such that each crypt is essentially sheathed in these cells. Though it is likely that important signaling takes place between the epithelium and stroma [166–168], these signals are difficult to study and their importance is unclear given the fact that *ex vivo* isolated crypts maintain a differentiating, homeostatic structure even in the absence of a macroscopically polarized stroma [149, 169].

The crypts are the functional unit of the intestine. The single epithelial layer of each crypt contains only ∼10 stem cells at the base that divide on average once per day, generating a conveyor belt of rapidly dividing and differentiating cells [170] (Fig. 2.5). These "transit amplifying" cells stop dividing after full differentiation and are eventually lost to the lumen of the gut [171]. This process is fast, such that the entire non-stem cell population of a crypt is turned over every 3-5 days [172].

In the small intestine, the top of the crypt (which is the base of the villus) is roughly the point at which cells are completely differentiated. Note that this generally accepted model of crypt turnover is based on studies of the small intestine and is thought to be similar in the colon, though there are important differences between the organs that should be considered. For example, the large intestine completely lacks the paneth cells that have recently been reported to be required for stem cell maintenance in the small bowel. However, colonic crypts do contain cells expressing similar surface proteins [150]. Additionally, the entirety of the colonic crypt contains apparently-differentiated cells, suggesting that the transit amplifying population is much smaller than in the small bowel. On the other hand, the literature generally supports that the population dynamics and signaling factors are similar between the small and large intestine, and between mouse and human.

**Opposing gradients of TGFB$_{sf}$ and Wnt in crypts**

Within the crypt epithelium, there are distinct patterns of expression for the Wnts and the TGFB$_{sf}$ ligands that together seem to control cell fate in a crypt-axis positional manner. Using fluorescence in situ hybridization (FISH), Clevers' group showed that canonical Wnts and various FZDs were preferentially expressed in the crypt base, while non-canonical Wnts were expressed throughout the crypt or only at the top [144]. The same group later found that LGR5, a Wnt co-receptor, is a highly specific marker for intestinal stem cells [148]. This discovery allowed for labeling and purification of this cellular population, and the demonstration that this cell type alone could re-create crypt-like structure *in vitro*.

There is evidence that the paneth cell population is the source of canonical Wnt3A in the crypt

**Figure 2.5:** Overview of TGFB$_{sf}$ and Wnt signaling patterns in the colonic crypt. **a**, Immunostained crypts from a CPC;APC mouse [173] showing two opposing sides pinched together, with the dark lumen in between. Blue, DNA; green, β-catenin, red, E-cadherin. The bottom set of crypts are tumorous; note the global increase in β-catenin. Image curtesy Michael Ramirez and Curtis Thorne (Altschuler & Wu lab, Univeristy of Texas Southwestern). **b**, Cartoon of colonic crypt structure, showing key cell types. **c**, Cartoon of signaling gradients in the crypt, showing high Wnt at the crypt base and high TGFB$_{sf}$ at the crypt tops.

base [150], though this is confounded by a study showing that this cell type can be ablated *in vivo* without loss of stem cell maintentance [174]. In contrast to canonical Wnt, but similarly to non-canonical Wnt, TGFB$_{sf}$ ligands and receptors are expressed exclusively in the the differentiated epithelium at the tops of the crypts [54,175] and villi [176]. Smad activity is restricted to the same parts of the crypt [143,177].

How the gradients are established and maintained, especially in the face of a constanty turning over cell population, is not well understood. The relatively long lifespans of the paneth cell population [178] hint that this cell type may serve as a sort of anchor for these gradients. Another potential source of gradient maintenance are the transmembrane ephrin receptors and their ligands, which seem to be required for proper cell sorting: their loss results in disorganized cellular positions within colonic crypts. Intriguingly, these ephrins and receptors are regulated by canonical Wnt signaling, though how these two phenomena interact is unclear [179].

Importantly, it is unclear whether these gradients actually function as such. While the gradients could be formed by secreted molecules diffusing over multiple cell lengths, the relative insolubility of Wnts argues against this. Indeed, in the developing fly wing and in the vertebrate notochord, there is recent evidence that diffusion of Wnt is not the mechanism by which it creates gradients [180,181]. Additionally, because Wnt and TGFB$_{sf}$ have so many extracellular antagonists, the presence of a ligand gradient is neither required nor sufficient to have a functional gradient [60].

### 2.4.2 Wnt drives stem cell maintenance; TGFB$_{sf}$ drives differentiation

While the opposing gradients of these two pathways hint at opposing functional roles, they do not imply it. The gradients could be a macroscopic consequence of cell sorting. Or, the gradients could be a marker of crypt position without exerting any concentration-dependent influence. The biological argument that most favors truly oppositional roles is that modulation of one pathway,

experimentally or in the context of disease, is generally paired with opposite deviations of the other pathway. I reviews examples of this below.

### Increased Wnt signaling is a driver of colon cancer

It is well established, and has been for some time, that over-activation of canonical Wnt is a primary force in colon cancer. Multiple components of the pathway have been implicated in this disease, though clearly some of the signaling nodes are easier to hijack than others. One could imagine upregulation of Wnts, receptors, or β-catenin as possible mechanisms. However, the pathway has negative feedback, via β-catenin-mediated expression of Axin2, and so constitutive upstream activity would be insufficient for maintenance of high β-catenin. Therefore constitutively high β-catenin activation is more easily obtained by blocking the function of the destruction complex. Perhaps for this reason, the most commonly modulated signaling nodes are components of the complex itself. There is at least one example of receptor-level modulation, however, which is that RSPO1 (the soluble Wnt co-factor) is frequently upregulated in colon cancer due to genomic translocation [182]. Whether this upregulation is functionally important to cancer, however, is still speculative.

By far, the most common Wnt pathway modification in colon cancer is mutation of APC, the destruction complex component with perhaps the most mysterious functional role. In fact, APC is nearly always mutated in this type of cancer (>80% of cases) [103, 129, 183]. Interestingly, APC is rarely completely lost, instead being frequently truncated to its N-terminal half. The reason for this truncation preference is unclear and has been the cause of much speculation [110]. However, it is worth noting that there need not be a reason. Perhaps there are more ways to mutate APC such that it is truncated instead of ablated, in which case truncation would simply be the more likely evolutionary step. However, there is some cell culture evidence that APC-truncated and APC-null cells have different phenotypes, especially with respect to cell adhesion and migration [184].

While it is not clear what exactly APC is doing in the destruction complex, the effects of APC loss on Wnt signaling are well-established. APC knockdown results in rapid nuclear accumulation of β-catenin, followed by increasingly dense cell packing due to overgrowth and then, eventually, a less-differentiated cellular phenotype [185]. Remarkably, stem cell-specific deletion of APC leads to adenoma formation in mere days and results in macroscopic tumor development in only 3-5 weeks [153]. Importantly, this effect requires β-catenin-dependent expression of Myc [183], a classic oncogene, as ablation of Myc can rescue the APC mutant phenotype [186].

### Decreased TGFB$_{sf}$ is a driver of crypt dysplasia

Given my earlier description of the TGFB$_{sf}$ pathways as inducers of differentiation, it is perhaps no surprise that these pathways tend to be lost in colon cancer and in other dysplastic diseases. As with the Wnt pathway, certain components of the TGFB$_{sf}$ pathways are more prone to being hijacked by disease than others. Loss of any particular TGFB$_{sf}$ ligand would likely be insufficient to ablate pathway activity, as would loss of several of the receptors, given the redundancy at this

level of signal transduction. The obvious targets then are the Smads, which are indeed mutated or lost in the context of several dysplasias.

Of the Smads, the most efficient target of ablation in disease would be Smad4, since it is the bottleneck for all upstream TGFB$_{sf}$ pathways. Indeed, Smad4 is commonly mutated in colon cancer (>15% of cases) and in pancreatic cancer (>50% of cases) [182, 187]. Additionally, next-generation sequencing of >70 human colon tumors revealed a high prevalence of Smad2 (one of the TGFβ-specific rSmads) mutations [182], allowing speculation on the importance of TGFβ signaling in colon cancer. In another dysplastic disease, Juvenile Polyposis, BMP signaling loss is found in >50% of cases [61, 177, 188]. Though it is unknown if the BMP loss in these cases was sufficient to cause the human phenotype, overexpression of the BMP-inhibitor Noggin is sufficient to generate ectopic crypts within the villi of mouse models [188].

The above data suggest that gain of Wnt signaling has a more potent effect than does loss of TGFB$_{sf}$ signaling, at least in the context of colon cancer. However, changes to one of these pathways is always accompanied by changes in the other. For example, BMP signaling has been found to be reduced in >70% of colon cancers [177], though only ~15% of cases have mutations in the pathway. Further, in a study of microarrays from 250 colorectal tumor samples it was found that Smad4 and β-catenin levels were generally anti-correlated [189]. Taken together, these correlative results are suggestive that these two pathways are in some way regulating one another, and that for Wnt signaling to become tumorigenic it has to suppress the TGFB$_{sf}$-mediated drive towards cellular differentiation.

## 2.5 Nodes of crosstalk between TGFB$_{sf}$ and Wnt

In this dissertation I am interested in understanding the inter-pathway crosstalk between TGFB$_{sf}$ and Wnt. As discussed in Section 1.4, one of the most problematic aspects of studying cell signaling is the determination of which input signals a cell cares about and into what intracellular property that information is encoded. As this chapter has so far detailed, the field consensus for both TGFB$_{sf}$ and Wnt is that it is the concentration of the ligands that carries the information that cells care about (these are morphogenic signals) and this information is encoded into nuclear concentrations of canonical transcription factors.

Further, the literature reviewed in the previous section strongly suggest that the TGFB$_{sf}$ and Wnt pathways have many opportunities for crosstalk, especially in the mammalian gut. With established input/output relationships in hand, and reason to suspect that the pathways in question modulate one another, we can begin to ask how these pathways integrate information. This section provides a review of the signaling nodes at which inter-pathway crosstalk is thought to occur (see the graphical summary in Fig. 2.6).

**Figure 2.6:** Overview of Wnt and TGFB$_{sf}$ pathway crosstalk. This dissertation focuses on canonical Wnt, BMP, and TGFβ as mediated by Wnt3A, BMP4, and TGFβ1/3. Dashed lines indicate literature-established links between pathways. References: *a* [190–192]; *b* [193]; *c* [89,90]; *d* [90,194,195]; *e* [100]; *f* [196].

## 2.5.1   Smad and DVL

The DVL proteins, being key post-receptor mediators of both canonical and non-canonical signaling, occupy an important position for potential crosstalk with the TGFB$_{sf}$ pathway. Interactions of DVL with the MH2 domain of Smads were first identified in a yeast two-hybrid screen [190]. The same group and others confirmed the possibility of such interactions in mammalian cells using co-immunoprecipitation (co-IP) after overexpression of Smads, with the conclusion that all three DVLs could bind to most of the Smads (Fig. 2.6a) [190–192]. The Smad-DVL interaction was also shown for non-canonical Wnt signaling, with a downstream mediator of Wnt5A, PAR1B, precipitating with both Smad and DVL (Fig. 2.6b) [193].

While co-IP of overexpressed proteins shows the possibility of interaction, they neither guarantee that it occurs under normal conditions nor imply that the interaction has a functional outcome. In some of the studies cited above, attempts were made to demonstrate functional consequences of Smad-DVL interaction. This was done by treating cells with ligands for one or both pathways after ablation or overexpression of other signaling components. Using this approach, it was found that BMP2 treatment could increase complexed DVL1/Smad1, while Wnt3A-conditioned media had the opposite effect [192], suggesting that pathway interactions are sensitive to signaling activity. Additionally, RNAi of all three DVLs, PAR1B, or ROR2 can attenuate TGFB signaling, while overexpression of DVL3 or FZD2 can enhance it [193,197]. The authors of these studies interpreted the data to mean that the TGFB$_{sf}$ pathway is directly modulated by the Wnt pathways during

signal transduction, though the mechanisms are unclear and seemingly complex. Importantly, the observed cross-pathway modulation seemed to be mediated by interactions between DVL and Smad.

A functional consequence of TGFB and non-canonical Wnt5A integration was recently reported in the context of wounded colonic epithelium [197]. In that study, the authors found that mice lacking Wnt5A were less able to repair colonic epithelial lesions, seemingly due to an inability of progenitor cells to differentiate in the wound. Given the use of TGFB$_{sf}$ pathways in gut stem cell differentiation, it is perhaps unsurprising that the authors could link a TGFβ signaling deficiency to this Wnt5A phenotype. Specifically, they found that treatment of *ex vivo* colonic crypts with high concentrations of Wnt5A yielded downstream TGFβ responses. The mechanisms for this were not clear, except for a dependence on the TGFBR2 and ROR2 receptors, but are consistent with DVL-mediated interactions. Importantly, my own work is suggestive that this Wnt5A/TGFβ connection is due to an artifact (Section 3.1).

### 2.5.2 Smad and Axin

Axin, being a large protein that is central to canonical Wnt signaling, is another good candidate for interactions with TGFB$_{sf}$ components (Fig. 2.6d). Indeed, in a co-overexpression assay of Axin and Smad3, the two proteins were found to co-IP. This interaction was dependent upon a region of Axin between the β-catenin and DVL binding domains, though how binding to this site might affect Wnt signaling is unclear. Further, TGFβ responsiveness was increased when Axin was overexpressed [194]. This was interpreted to mean that Axin can stabilize Smad activity, either by preventing Smad degradation or by blocking some other form of Smad interference.

Two additional studies looked into the consequences of Smad-Axin interactions, but with contradictory results. In one case, Axin was found to bind to an iSmad (Smad7) and to mediate degradation of that Smad by Arkadia, an E3 ubiquitin ligase. As a consequence, overexpression of both Axin and Arkadia led to enhanced TGFB signaling [195]. While the outcome is the same as that described above, the mechanism is essentially the opposite. In a different study, Axin was again found to destabilize a Smad, but this time an rSmad via GSK3β instead of an iSmad via ubiquitination. As a result, TGFB signaling was instead attenuated by increased Axin, and Axin depletion by RNAi amplified TGFB signaling [90].

How can we make sense of the studies of Axin-Smad interaction, that are all mutually incompatible? It is important here to recall that basal Axin levels in cells are extremely low [131]. Because this protein acts as a scaffold, its overexpression can have a few obvious non-physiological effects. The first is that Axin is a big protein, with many binding surfaces, and so the observed Smad-Axin interactions may be due to what would normally be extremely rare binding events. The other is that while Smad and Axin may interact, overexpression of the scaffold causes interactions between Smad and other Axin-bound proteins to become at first enhanced and then diluted as the amount of Axin increases. All of the cited studies use Axin overexpression, and so the apparent contradictions between them may be due to either of these phenomena.

**Smad and GSK3β**

GSK3β, which primes β-catenin for ubiquitination, is able to phosphorylate a large fraction of the proteome. Recently, this was found include the bSmads (experimentally) and the tSmads (based on computational prediction). The Smads have a canonical GSK3β recognition sequence in the linker region between the functional MH1 and MH2 domains. Phosphorylation of these sites in bSmads was found to be GSK3β-dependent and could be down-regulated by Wnt3A treatment (Fig. 2.6c). Importantly, the study found that total Smad levels were not affected by Wnt treatment, but that the active phospho-state did show measurable decreases [89]. This data was therefore interpreted to mean that GSK3β could modulate the long-term duration of Smad activity.

**Smad and β-catenin**

Finally, we turn to interactions with the final effector of the Wnt pathway, β-catenin (Fig. 2.6f). Like Axin, the large size of β-catenin provides many potential binding interfaces. Additionally, its role as the bottleneck of all canonical Wnt signaling makes it a good candidate for cross-pathway interaction. Even so, evidence for interactions between Smad and β-catenin are quite limited. There is some evidence that *Drosophila* MAD and Armadillo (Smad and β-catenin orthologs) compete with one another for binding to TCF, such that DPP (the BMP2/4 ortholog) can cause a downstream block of β-catenin transcriptional output [196]. Similarly, though with different functional consequences, in mammals an iSmad (Smad7) can co-IP with β-catenin and overexpressed TCF/Lef [198]. It is therefore unclear whether a Smad/β-catenin interaction is important to signaling through either pathway.

**Transcriptional crosstalk**

The nodes of putative crosstalk described above occur prior to transcription factor entry into the nucleus. I therefore classify these as nodes of "signaling crosstalk." What about at the level of transcription, after the nucleus has received the transcription factor output of each signaling pathway? Our knowledge of $TGFB_{sf}$/Wnt crosstalk is almost entirely due to studies at this level of interaction, however studies that look at both pathways simultaneously are rare. Instead, transcriptional crosstalk is often inferred by studies finding that activity of one pathway modulates transcriptional output commonly associated with the other.

A few direct crosstalk studies have been performed, though their results are not easily comparable. In the case of TGFB and Wnt, activation of both pathways was found to increase output of a Wnt transcriptional reporter. The lack of Smad-response elements in this promoter suggested that the co-activation was due to an interaction occurring prior to promoter binding [199, 200]. Microarrays from similar co-treatment experiments showed that a set of genes were regulated differently in the context of both inputs than with either input alone, though no obvious patterns were found (e.g. genes that were increased by Wnt or TGFB were not necessarily further increased by both) [155, 201].

Transcriptional crosstalk is most frequently inferred by the presence of consensus binding motifs for β-catenin/TCF and the Smads within the same target gene promoter. In this way, instances of direct transcriptional co-regulation by TGFB/Wnt have been found in several systems [159,202–204]. For the focal case of the gut in this section, the co-regulation of the Myc promoter by these two sets of transcription factors is of obvious relevance [205] .

## 2.6   Discussion

The preceding chapter provided a glimpse of the complexity and uncertainty in the properties of Wnt and TGFB$_{sf}$ signaling and inter-pathway crosstalk. Here, I summarize the salient points that lead to the approaches and hypotheses of my experimental work in Chapter 3.

### 2.6.1   Wnt and TGFB$_{sf}$ use morphogenic encoding systems

The consensus in the literature is that Wnt and TGFB$_{sf}$ signaling are morphogenic, meaning that they encode information about ligand concentrations. As I have noted throughout this text, the fact that these pathways yield concentration-dependent effects does not imply that this is the only information about the signals that cells are using.

In Section 4.4.4 I discuss how the information content of biological outputs can be quantified using the mutual information metric. Published reports indicate that the information content of many biological input/output relationships, when measured by imaging, is quite low [16]. My own measurements, using the careful image correction discussed in Chapter 4, yield only slightly higher information content. This low mutual information between ligand concentrations and the outputs of signaling imply that cells can only effectively determine whether a signal is present or not; they cannot determine a precise absolute concentration. In effect, these morphogens are not morphogenic at the single-cell level! This limited information content of the concentration-based encoding system is especially interesting in light of the dramatic network bottlenecking of both pathways. Together, the implication is that cells can neither determine which of many different ligands is present in the environment nor the precise concentrations of those ligands. How much information about these pathways, then, do single cells have access to?

On the other hand, we may be neglecting additional information channels that cells use to encode more accurate models of their extracellular environments. "Non-canonical" information channels may provide just such content. Or, just like the fly embryo syncytium that makes use of multiple noisy transcription factor gradients to obtain accurate positional information (see Section 1.4) [43, 44], cellularized systems such as the intestinal crypt may integrate Wnt and TGFB$_{sf}$ concentrations in order to more accurately define fates or positions along the crypt axis.

### 2.6.2   Signaling crosstalk reduces information content

In the preceding sections and in my primary experimental work in Chapter 3, I take care to classify pathway crosstalk into distinct types: pre-nuclear signaling crosstalk versus nuclear crosstalk at

the level of transcription. This is for an important reason, which is that integrating pathways at the level of signal transduction causes a loss of information content available to the nucleus. In other words, the nucleus loses knowledge about the environment when pathways intersect during signaling.

As an example, the cell makes an internal model of environmental $TGFB_{sf}$ using concentrations of active Smad. If $TGFB_{sf}$ were the only thing that could activate Smad, then the cell would "know" that $TGFB_{sf}$ was present in the environment any time that active Smad levels increased. But what if Wnt could also modulate Smad? In this case, when the nucleus sees changes to active Smad levels, it cannot know whether $TGFB_{sf}$, Wnt, or some combination of the two ligands are present in the environment. This contrasts to the case where $TGFB_{sf}$ ligands only modulate Smad, and Wnt only modulates β-catenin, so that the internal model within the cell models the more complex reality of the external environment.

### 2.6.3   Informational asymmetry between the cell and its environment

Perhaps cells only need limited information about their environments, such that the reduction of information caused by complex inter-pathway processing during signal transduction is of no consequence. This would be surprising, however, as these very same cells are what create the information-rich extracellular milieu in the first place (e.g. by secreting signals and otherwise modifying the environments of their neighbors). How could information-poor intracellular networks create information-rich extracellular environments?

One possibility is that each cell type present in a complex environment does indeed have limited information about that environment. It could then be the case that the environmental complexity is simply due to the combination of low-information outputs from many distinct cell types. In effect, individual cell types would be insulated from the complexity of their environments, and would only have to worry about processing the few signals that they are capable of "understanding."

### 2.6.4   Uncertainty in the nodes of $TGFB_{sf}$/Wnt signaling crosstalk

Crosstalk at the level of signaling reduces what a cell can know about its environment. It is therefore important to determine whether such crosstalk truly exists before speculating on why a cell would need so little environmental information, as I so prematurely began to do above.

Cellular signaling is extremely difficult to study, for the reasons outlined in Chapter 1, and the studies cited in the current chapter are illustrative of this fact. While it is possible that the results of all of the cited studies are accurate, they must be interpreted with care. In particular, all of the cited studies that identified nodes of $TGFB_{sf}$/Wnt crosstalk relied on some combination of overexpression, RNAi knockdown, or pharmacological inhibition of pathway components. All of these methods can push the global cellular signaling network into states that normal cells cannot inhabit. This is especially true for studies involving β-catenin and Axin, as these proteins are scaffolds that are normally present at extremely low concentrations.

**Table 2.7:** Sources of structures used in alignments for Fig. 2.2 and Fig. 2.4.

| Symbol | Species | PBD ID | Source |
|--------|---------|--------|--------|
| TGFB1 | *Homo sapiens* | 3KFD.A | [66] |
| TGFB2 | *Homo sapiens* | 2TGI.A | [206] |
| TGFB3 | *Homo sapiens* | 2PJY.A | [71] |
| BMP2 | *Homo sapiens* | 1REW.A | [207] |
| SMAD3 | *Homo sapiens* | 1U7F.A | [208] |
| SMAD4 | *Homo sapiens* | 1U7F.B | [208] |
| SMAD7 | *Homo sapiens* | 3KMP.A | [209] |
| Wnt8 | *Xenopus laevis* | 4F0A.B | [108] |
| Fz8 | *Xenopus laevis* | 4F0A.A | [108] |

Also, the further apart experimental inputs and outputs are in time, the less accurate our inferences can be about the mechanisms connecting them. The cited studies rely heavily on co-immunoprecipitation of overexpressed pathway components to demonstrate the possibility of interactions. They then use transcriptional reporters and other long-term readouts to measure the consequences of these interactions. However, a relationship between physical interactions and downstream consequences can only be demonstrated by specifically blocking that interaction. This has not been done for any node of putative TGFB$_{sf}$/Wnt signaling crosstalk, and for good reason: such an experiment is exceedingly difficult to design.

In short, current studies on TGFB$_{sf}$/Wnt signaling crosstalk do not show either that such crosstalk definitively exists, nor that the crosstalk has functional consequences if it does exist. Indeed, in Chapter 3 I show strong evidence under endogenous signaling conditions that these pathways are insulated from one another during signaling. This implies that nuclei maintain more accurate models of their environments by integrating TGFB$_{sf}$ and Wnt primarily at the level of transcription.

It will be important for future studies to accurately determine the information content of the TGFB$_{sf}$ and Wnt signaling pathways, and how that information is encoded by the cell. Such studies may reveal that cells encode more about the signal than just its concentration. In combination with the results of Chapter 3, that the integration of these pathways gives cells more information, we may find that cells create more accurate internal models of their extracellular environments than is currently believed.

## 2.7  Methods

**Sequence alignments.** I obtained sequence data from the National Center for Bioinformatics (NCBI) servers, choosing each time the top listed isoform from the Genes database Protein structures are from the Research Collaboratory for Structural Bioinformatics (RCSB) protein database (PDB) (Table 2.7). I used subsets of these sequences and structures for alignment in Promals3D [210] to obtain phylogenetic trees. This multi-sequence alignment algorithm takes advantage of structural information for improved alignments, however the distances in the phylogenies should be considered approximate since the method for calculating genetic distance is more naive than the alignment method. The crystal structures were included in the Wnt and Frizzled trees because of species

differences from the primary sequence data. Structures for the TGFB$_{sf}$ and Smad trees were dropped because their sequences identically match a subset of the primary sequences. The trees were drawn with the R package 'ape' [211].

For predicting *Trichoplax adhaerens* orthologs of genes, I used PSI-BLAST [212] with some combination of human and *Drosophila melanogaster* reference proteins. I ran one or more iterations of the algorithm on the seed sequences, and then chose the few top hits that were dramatically better matches by both query coverage and identity. Each putative *Trichoplax adhaerens* ortholog is listed in the appropriate table in this chapter. These orthologs are in agreement with the literature [52].

Specifically, to identify TGFβ/BMP orthologs I used the TGFβ1-3, BMP2/4, and DPP sequences as seeds, resulting in 3 putative orthologs (>50% coverage). For the TGFB receptors, I used the Type I and Type II receptors separately, but each yielded a large number (>200) of putative matches that had reasonable identity (>30%) but low coverage (<50%) or vice versa. I therefore did not include these in the phylogenetic tree. For the Smads, I used the eight human genes and the *Drosophila* MAD as seeds, yielding three good hits (>90% coverage, >30% identity). For the Wnts, I used all 19 human genes and fly WG, yielding 2 hits (≥68% coverage, >30% identity). For the Frizzleds, I used all 10 human genes and fly FZ, yielding 2 hits (>50% coverage, ≥30% identity). For β-catenin I used human CTNNB1 with Junctional Plakoglobin (JUP, also known as gamma-catenin) and fly arm, yielding 1 hit (>70% coverage, >70% identity).

# Chapter 3

# On the insulation of morphogenic signaling

Morphogenic signals are frequently found in apparent gradients within developmental systems and stem cell niches, and the set of concentrations of these morphogens at a given point in space and time is thought to provide the information needed for cell fate specification (see Section 2.4.1). How cells integrate the concentrations of distinct morphogenic signals is an unsolved problem. Cells could integrate this information during the process of signal transduction to the nucleus (e.g. by direct protein-protein interactions between pathways), after the transduced signals reach the nucleus (e.g. by co-regulation of transcription targets), or at both of these levels.

Importantly, as discussed in Section 2.6, integration at the level of signal transduction (hereafter "signaling") may lead to a decrease in nuclear "knowledge" of the original signals. Integration at the level of transcription, however, can allow cells to maintain a more accurate internal model of the extracellular environment. Therefore, in order to understand how cells make decisions in the context of multiple extracellular information sources it is important to identify the points at which pathway integration occurs.

The Wnt and $TGFB_{sf}$ pathways provide highly-studied systems for understanding how cells integrate morphogenic signals. As discussed in Chapter 2, the Transforming Growth Factor Beta superfamily ($TGFB_{sf}$) and Wnt/β-catenin (hereafter simply "Wnt") signaling pathways are deeply conserved across metazoans, are essential to development, and are disrupted in many pathologies. These pathways are tightly intertwined, frequently being used within the same tissue compartments to coordinate cell fate decisions. This coordination occurs despite an absence of shared core pathway components, which suggests that it is primarily mediated by long-term transcriptional interactions. However, a number of studies have also identified putative nodes of short-term signaling interaction (see Section 2.5 and Fig. 2.6), though the generality and importance of these interactions remain unclear.

Wnt and TGFB are morphogenic: their extracellular ligands lead to concentration-dependent increases in downstream transcription factor activity. Outside of this general similarity, the mech-

**Figure 3.1:** Two signaling pathways ($S_1$ and $S_2$) can crosstalk in multiple ways with respect to a a response $R$. **a-b**, In the absence of signaling from one or both pathways there can be no crosstalk. **c**, The signals can both affect $R$ in an additive manner, so that the signaling outcomes are independent. Such a case would occur for pathways that use distinct "pools" of the same component. **d**, If both pathways affect the same pool of $R$, or affect one another, then they will show non-additive, interdependent behaviors. This can arise through several distinct topologies that are difficult to distinguish experimentally.

anisms by which these pathways transduce their respective signals are quite distinct. As reviewed in Section 2.2, TGFβ and the related Bone Morphogenic Protein (BMP) ligands cause their serine/threonine kinase receptors to directly phosphorylate the target Smad transcription factors, which subsequently increase in nuclear abundance. Activation of the Wnt pathway, on the other hand, blocks the otherwise constitutive degradation of cytosolic β-catenin, thus leading to a whole-cell increase in the quantity of this transcription factor (reviewed in Section 2.3). The transduction of a Wnt signal requires many protein components, most of which have been implicated in direct interactions with Smad proteins (Section 2.5). Unfortunately, identification of interactions between these pathways has so far led to contradictory and context-dependent outcomes, suggesting that there is not a general mechanism of TGFB$_{sf}$ and Wnt signal integration.

However, as I explain in Section 2.6, it is possible that the methods used to study integration of these pathways are simply incapable of identifying general mechanisms of crosstalk. The majority of the crosstalk studies (and even the studies of each pathway in isolation) rely on overexpression or ablation of pathway components, which may push cells into abnormal states and thus confound interpretation of experimental results. In this chapter, then, I take an endogenous and mechanism-independent approach to directly test the extent of signaling crosstalk between these pathways.

Signaling pathways can interact in multiple ways during transduction. Without any crosstalk, activation of one pathway will by definition have no effect on the canonical output of another pathway (Fig. 3.1b). With **non-additive** crosstalk, two pathways may affect the same response but do so in an additive manner (Fig. 3.1c). This would be the case, for example, with pathways that use different pools of the same component. Finally, pathways can interact in more complex, non-additive ways in which the response cannot be predicted by knowledge of one pathway alone (Fig. 3.1d). Which of these crosstalk categories Wnt and TGFB$_{sf}$ fall into during signal transduction

has not yet been uncovered.

I therefore designed an experimental approach that allows me to distinguish between the three general classes of interaction described above. By measuring the direct output of signaling for each pathway (i.e. transcription factor nuclear concentration) as a consequence of combinatorial $TGFB_{sf}$ and Wnt ligand inputs, I can infer both the class of crosstalk that the interactions fall into and the quantitative extent of that crosstalk. Using this approach, I show in Section 3.2 that Wnt and $TGFB_{sf}$, in opposition to reports in the literature, do not crosstalk at all during signal transduction.

Further, in Section 3.3 I show that TGFβ is insulated from BMP signaling despite sharing the core component Smad4. Intra-$TGFB_{sf}$ inhibition is widely thought to exist and to be due to competition for limiting Smad4 (see Section 2.2.5). I find instead that neither of these claims are correct: BMP4 and TGFβ3 do not inhibit one another even when Smad4 is brought down to limiting levels.

Taken together, my results suggest then that cellular decision-making with respect to TGFB and Wnt occurs primarily at the level of transcription and not at the level of signaling, thus allowing the cell to create a more accurate nuclear model of the complex extracellular microenvironment than would otherwise be possible.

## 3.1   An endogenous system for studying $TGFB_{sf}$/Wnt crosstalk

I reasoned that, because transcription factor activity is the direct endpoint of signal transduction, I can infer the degree of meaningful $TGFB_{sf}$ and Wnt signaling interaction by measuring how stimulation of one pathway affects the immediate transcription factor response of the other pathway (Fig. 3.2). Current studies typically rely on transcriptional readouts to infer such interactions, but these inferences may be confounded by transcriptional feedback (which I consider to be the result of nuclear decision-making, not signal transduction). Therefore, to accurately interpret cross-pathway effects with respect to signaling, I use experimental timepoints that are as close to the initial signaling event as possible. Additionally, it is widely believed that both Wnt and $TGFB_{sf}$ signaling can be highly context-dependent. I therefore repeat the experiments in this chapter using multiple cell types, chosen to represent divergent cellular contexts. By doing so, the hope is that any resulting shared properties of signal integration can be more confidently extrapolated to other cellular systems.

In order to rigorously quantify single-cell responses to the many experimental conditions required for this study, and to make use of the expertise within the Altschuler & Wu lab, I use high-throughput immunofluorescence imaging as my primary experimental platform. Unless otherwise indicated, the presented measurements in this chapter originate from the total-intensity feature values of individual nuclei. I typically report the population-medians of these single-cell values, from replicate experimental setups (as shown schematically in Fig. 3.3; see Chapter 4 for more detail on my approach to image analysis).

**Figure 3.2:** Schematic of the experimental system for TGFB$_{sf}$/Wnt crosstalk. To measure inter-pathway signaling crosstalk, ligands can be applied combinatorially, followed by measurement of nuclear transcription factor levels. As shown, Wnt3A increases global quantities of β-catenin, TGFβ3 causes translocation of bulk Smad2/3 to the nucleus, and levels of nuclear phospho-Smad1/5/8 increase upon BMP4 treatment. HCECs, with ("High") or without ("Low") 2hr treatment by Wnt3A (240ng/mL), TGFβ3 (4ng/mL), or BMP4 (73ng/mL). Outlines are of nuclei, segmented from the Hoechst channel (channel not shown) using the same threshold-segmentation method as in all imaging studies in this chapter. Fields are chosen to demonstrate visually obvious outcomes of signaling, though the diversity of responses is quite high for all pathways (see Figs. 3.5 & 3.7).



**Figure 3.3:** The fluorescence intensities reported in this chapter are population-level, based on single-nuclei measurements. As shown schematically here, cells are treated in 96- or 384-well plates (left) and then fixed, immunostained, and imaged (middle, top). Nuclei are then identified computationally (see Section 4.4) so that the distributions of single-nuclei total fluorescence are obtained (middle, bottom). The medians of these distributions are then calculated for each replicate well. The reported values in this chapter are the means and standard deviations of these median values, nearly always from $n=3$ replicates. P-values are then calculated using an unpaired, two-tailed Student's t-test, and significant differences (p<0.05) are indicated with an asterisk. Figure legends indicate which samples are being compared for statistical significance.

**Figure 3.4:** A screen for Wnt3A (**a**) and TGFβ1-responsiveness (partial-replicate screens **b** and **c**) across ∼20 cell lines revealed consistent responses by human colonic epithelial cells (HCECs) and two melanoma lines (SKMEL2 and MALME3M), highlighted with bright red text. The measured single-cell feature is the nuclear mean of fluorescence intensity by imaging (arbitrary units). Plots show the median of these single-cell values across all cells in a treated well (as in Fig. 3.3). The control means are subtracted from all values, per cell line, to show absolute changes in nuclear intensity. Concentrations: 1ng/mL TGFβ1, 200ng/mL Wnt3A. Timepoints: 1.25hrs (Wnt3A), 1hr (TGFβ1).

### 3.1.1 Choosing cell types

In order to choose useful cell lines, ligands, and readouts for the study of Wnt and TGFB$_{sf}$ signaling, I was confronted with something of a chicken-or-the-egg problem. However, ongoing work within the Altschuler & Wu lab had made use of purified recombinant TGFβ1 for stimulating the TGFβ pathway and a total-Smad2/3 antibody for measuring the response. Additional work had shown the efficacy of a β-catenin antibody for measuring cellular responses to purified Wnt3A. I therefore first made use of these reagents, using literature-supported concentrations, to identify cell lines that show responsiveness to these pathways.

I first selected an immortalized (via telomerase and CDK4 expression) but non-transformed human colonic epthelial cell (HCEC) line. This cell line has stable ploidy and properties consistent with it being a pseudo-differentiated cell type [213]. Given the importance of Wnt and TGFB$_{sf}$ signaling in the gut (Section 2.4.1), this cell line makes for a reasonable model system for studying crosstalk between these pathways. As a control, I also chose the rat small intestine epithelial line (IEC6) [214], reasoning that it should display similar properties to HCECs. I found that both of these cell types respond strongly (in a statistical sense) to TGFB ligands, moderately to BMP4,

and weakly (but measurably) to Wnt3A.

As discussed in Section 1.4.3, signaling pathways generally show context-dependency, especially with respect to differences in cell type. This non-generality is especially true of the TGFB$_{sf}$ and Wnt pathways (Chapter 2). Therefore, in order to discover general properties (if indeed they exist) we need to look for commonalities across cell types; the two intestinal cell lines chosen may not be sufficient to infer generality. One approach to increase generality would be to exhaustively test a large number of cell types, so that with each successive cell type we gain confidence in the generality of a biological phenomenon. Unfortunately, this approach is costly and difficult, and still does not lead to certainty in the generality of a discovered phenomenon. I opted then for a simpler approach, which is to test a small number of divergent cell types. In this way, any behaviors consistent across cell types can still be extrapolated as "general" behaviors, though the resulting confidence in such a generalization may be somewhat lower.

To identify additional cell types for studying TGFB$_{sf}$ and Wnt crosstalk, I screened a panel of cancer cell lines for responsiveness to TGFβ1 and Wnt3A. Additional selection criteria included cell morphology and growth patterns that would allow for accurate image segmentation (see Section 4.4), as well as cellular growth rates and adherence properties that would allow for the throughput needed for the many experimental conditions required in my studies. Two melanoma cell lines, SKMEL2 and MALME3M, satisfied these criteria and were consistently ranked among the most responsive to both TGFβ1 and Wnt3A (Fig. 3.4). Both of these cell lines can form malignant melanomas in nude mice, and have abnormal ploidy [215]. In particular, I found that SKMEL2 cells always form tri-modal cell cycle distributions, implying the presence of diploid, tetraploid, and octoploid cells within an asynchronous cycling population. To my knowledge, there are no Wnt or TGFB$_{sf}$ pathway mutations in SKMEL2 or MALME3M cells.

These cell lines should not be used to make inferences with respect to differences between "normal" and cancer cells, as these cell lines differ also in tissue of origin. Instead, they should be thought of simply as divergent pairs of cell types that provide consistent "contexts" (e.g. within the two intestinal lines or within the two melanoma lines ) as well as divergent contexts (e.g. between the intestinal and melanoma lines). For all analyses in this chapter, I use only those cells within the first peak of the imaging-based cell cycle distributions (see Section 4.5.1). I verified for each individual pathway that the position of a cell within the cell cycle distribution was not predictive of pathway behavior (data not shown) and restriction to the G1/0 cell cycle phase otherwise reduces both experimental noise and unimportant biological variation.

### 3.1.2 Choosing signaling inputs and outputs

Having chosen the cellular contexts in which to measure the signaling crosstalk between TGFB$_{sf}$ and Wnt, I then needed to choose inputs and outputs that could be believably interpreted to represent these pathways. As discussed in Section 1.4, as a rule it is not known what aspects of a signal or a readout are the most relevant carriers of information for a cell. It is generally believed, however, that for morphogenic pathways it is the extracellular concentration of a ligand and the nuclear

**Figure 3.5:** Saturating ligand concentrations of TGFβ1/3 (12ng/mL), Wnt3A (730ng/mL), and BMP4 (8ng/mL) are informative with respect to their canonical transcription factors. MI, mutual information between the ligand and readout (in bits, maximum MI is 1, see Section 4.4.4). Histograms, distributions of single-cell total nuclear intensities for the immunostained transcription factors. $n$, number of cells per histogram, pooled from 3 replicate wells (color-coded). Gray, untreated. SKMEL2 cells, 1hr treatment. Arbitrary fluorescence units, frequencies scaled to have the same maximum value for display purposes.

concentration of a corresponding transcription factor that are the relevant parameters. I therefore tested several ligands and readouts in order to identify experimental inputs and outputs that are both meaningful and practical.

**Choosing prototypical ligands**

I first assayed several pathway inputs for information content and signaling specificity. For the BMP2/4 pathway, I found that all cell lines were non-responsive to purified BMP2 (data not shown) but responsive to purified BMP4. These two ligands are highly homologous and are thought to act through the same receptors (Section 2.2.2) therefore this absence of effective BMP2 signaling does not have a clear interpretation (perhaps a faulty reagent). In any event, the BMP4 signal is informative (using the mutual information metric, as described in Section 4.4.4); the distribution of single-cell responses to saturating BMP4 is wide but significantly different from the control distribution (Fig. 3.5).

For the TGFβ pathway I compared TGFβ1 and TGFβ3, as these ligands are considered to be essentially interchangeable (Section 2.2.2). Indeed, within a single experiment these two ligands generated similarly broad and separated Smad2/3 responses with similar mutual information (Fig. 3.5). Dose-response curves for the two ligands have the same maxima and similar hill coefficients, though TGFβ3 is ~10 time more efficacious than TGFβ1 (data not shown). TGFβ3 generally yielded more reliable responses, and so I use this ligand for the remaining experiments in this chapter.

Finally, I chose to use purified Wnt3A to stimulate the canonical Wnt pathway. Wnt3A is considered the prototypical ligand for this pathway (Section 2.3.2), and its commercially-available purified form has been widely used throughout the literature. Importantly, I discovered that what is likely the most commonly-used form, a low-purity (~75%) version from R&D Biosystems, is sufficient to send Smad2/3 to the nucleus with the same kinetics and dose-response Hill coefficient as seen with TGFβ treatment (Fig. 3.6a). I was unable to measure contaminating TGFβ ligands by Western (Fig. 3.6b), though this could be due to low concentrations of a high-efficacy TGFβ variant.

**Figure 3.6:** The commonly-used low-purity (75%) Wnt3A and Wnt5a ligands have trace amounts of contaminating TGFβ. **a,** Dose-response curves for Wnt3A and Wnt5A against Smad2/3 yield $EC_{50}$ concentrations for the Wnts that are comparable to those used in the literature to stimulate Wnt responses (>100ng/mL is common). Duplicate experiments. Median of the mean nuclear intensity is plotted for each well (as in Fig. 3.3). **b,** Western blot of the purified Wnt3A ligands using a pan-TGFβ antibody. HP/CF is the high-purity/carrier-free ligand used throughout this chapter, +BSA is the low-purity Wnt3A. BSA is shown as a reference. The low-purity Wnt3A lane should contain ~25µg of BSA in this blot, which is enough protein to soak up significant TGFβ antibody. TGFβ1 is ~12 kiloDaltons. **c,** The apparent Wnt→Smad2/3 response is completely blocked by co-treatment with a pan-TGFB antibody that does not block the ability of Wnt3A to stimulate a β-catenin response. Low-purity Wnt5A also stimulates Smad2/3, but a carrier-free (CF) version does not. The high-purity (>90%) carrier-free (HP/CF) Wnt3A shows a minimal Smad2/3 response that is reproducible in other experiments, though not significant in this one. Concentrations: Wnt3A (200ng/mL), Wnt5A (100ng/mL), αTGFB (5µg/mL). Y-axes and p-values as in Fig. 3.8. **a,c,** HCECs, 1hr treatment.

In any event, this pathway crosstalk is likely a consequence of trace amounts of contaminating TGFβ in the purified Wnt, as a pan TGFβ-blocking antibody is sufficient to block this response. Further, high-purity and carrier-free variants of the product do not activate Smad2/3 (Fig. 3.6c).

This artifactual crosstalk also occurs with the prototypical non-canonical Wnt5A (Fig. 3.6a,c), which casts uncertainty onto recently-published work linking Wnt5A and TGFβ signaling (reviewed in Section 2.5.1) [197]. The presence of such contamination may complicate the interpretation of many other published studies, as it suggests that treatment with purified Wnt3A or Wnt5A may generally include treatment with TGFβ, such that some resulting phenotypes could be due to stimulation by Wnt, TGFβ, or both.

I therefore use high-purity/carrier-free Wnt3A for all experiments in this chapter. I note, however, that even this ligand often causes a small increase in Smad2/3 and so crosstalk experiments must be interpreted with this fact in mind.

**Figure 3.7:** Nuclear phospho-Smads (pSmads) contain more information about ligand concentrations than do nuclear total-Smads, though the sum and mean of intensity features are comparable. The untreated (gray) and treated (colored) distributions of single-nuclei measurements show more overlap of total-Smad readouts, resulting in relatively less mutual information (MI, in bits, maximum of 1, see Section 4.4.4). The median (top row) and sum (bottom row) of nuclear intensities for a single readout have nearly identical information content. $n$, number of cells per histogram (color-coded). Ligand concentrations: saturating 10ng/mL TGFβ3 (blue) or 50ng/mL BMP4 (green). SKMEL2 cells, 1.5hr treatment. Arbitrary fluorescence units, frequncies normalized to have same maximum value.

**Choosing prototypical readouts**

Having identified robust pathway inputs, I then needed to validate antibody-based pathway outputs for immunofluorescence imaging. For the TGFβ pathway, which specifically activates Smad2/3, and the BMP4 pathway, which specifically activates Smad1/5/8 (see Section 2.2.4), I find that the nuclear fractions of both active phospho-protein and total-protein levels can respond robustly to pathway activation. However, it is unclear from the literature whether it is the total concentration or the phospho-state concentration that encodes extracellular ligand levels.

I therefore measured the mutual information between these readouts and their ligands, which revealed that the phospho-state is indeed more informative (Fig. 3.7). Care should be taken when interpreting this data, however, as the difference in information content may also be due to antibody specificities, along with a myriad of other causes. I note that I have never observed maximum mutual information values of more than ∼1.2 bits for any type of input/output relationship from imaging data, consistent with published reports [16]. The most fair statement, then, is that these particular approximations of the phospho-state are more informative than these approximations of total protein levels.

The total-Smad2/3 antibody yielded more-consistent results across experiments than did the

**Figure 3.8:** The measured canonical transcription factor responses are a specific consequence of ligand treatment. **Left**, A pan-TGFB antibody (αTGFB) blocks 2hr TGFβ3-induced nuclear Smad2/3 accumulation in HCECs. **Middle**, Dikkopf, a soluble antagonist (Section 2.3.3), blocks β-catenin accumulation due to 2hr Wnt3A in SKMEL2s. **Right**, Noggin, a soluble antagonist (Section 2.2.2), blocks 1.5hr BMP4-induced phospho-Smad1/5/8 in SKMEL2s. Y-axes, median across single-cell nuclear intensities (the total feature) within a well. As in Fig. 3.2, points are the mean and standard deviation of $n = 3$ replicate wells and '*' indicates one-sided p-value $<0.05$ (Student's t-test) compared to control. Intensities normalized by $R_{i,norm} = \frac{R_i - mean(R_{ctrl})}{R_{ligand} - mean(R_{ctrl})}$. Concentrations: TGFβ3 (10ng/mL), αTGFB (5µg/mL), Wnt3A (200ng/mL), Dikkopf (1µg/mL), BMP4 (50ng/mL), Noggin (100ng/mL).

phospho-Smad2/3 antibody, however, and has similarly-high information content. I therefore use the total-Smad2/3 and the phospho-Smad1/5/8 (pSmad1/5/8) antibodies for the crosstalk experiments. For canonical Wnt signaling there is strong evidence that it is the total-protein level of β-catenin, not the phospho-state, that encodes ligand concentration (Section 2.3.4). I therefore use a total-protein β-catenin antibody to measure canonical Wnt responses.

**Validating input/output relationships**

Aside from measuring the mutual information between each input and output it is essential to ensure that the input/output relationships are not artifactual. To do so, I verified that each output could be blocked by a highly specific antagonist (Fig. 3.8). Additionally, while I am primarily interested in the signal transduction process, it is important to ensure that a transduced signal leads to a nuclear decision. In the case of TGFB$_{sf}$ and Wnt signaling, this decision is a change to the transcriptional network. The conserved target for Wnt3A is Axin2, a negative auto-regulator (Section 2.3.4). For TGFB$_{sf}$, the conserved targets are the inhibitory-Smads (iSmads), Smad6/7 (Section 2.2.4). I therefore measured mRNA levels of these transcriptional targets, finding that stimulation does indeed yield transcription in all cases (Fig. 3.9).

### 3.1.3 Interpretation of single-cell immunofluorescence measurements

As for any method of measurement, it is important to take a step back and think carefully about the biological interpretation of the resulting data. For the image-based single-cell data obtained in my work, cellular nuclei are identified using Hoechst-staining and threshold segmentation, and I only analyze those nuclei that likely belong to the G1 cell cycle phase (see Section 4.4.1). Within these nuclei I tally up all pixel values to obtain the total intensity feature, which I use as a proxy for the quantity of the immuno-labeled target protein within the nucleus. This feature and the mean intensity feature, which is a proxy for concentration, are similarly informative in my assays (Fig. 3.7), and so I use the total-intensity feature for the practical reasons described in Section 4.4.3.

**Figure 3.9:** The TGFβ3, Wnt3A, and BMP4 ligands cause expression of canonical target genes in the cell lines used in this dissertation. Fold-change is relative to the control condition. 2hr treatment. Concentrations: 10ng/mL TGFβ3, 25ng/mL BMP4, 200ng/mL Wnt3A. Mean and standard deviation of 3 replicates, '*' indicates two-sided p-value <0.05 (Student's t-test). mRNA levels measured by TaqMan qPCR (see Methods for experimental details).

How do we interpret changes in this feature? Fold-change over a reference (e.g. the untreated state) is a commonly used metric, but the interpretation of fold-change is unclear in cases where either the basal state is near zero (as any fold-change value will move towards infinity) or is non-zero because of "background" signal (which pushes all fold-change values towards 1). Using the fluorescence image model that I present in Section 4.2, we can model the total intensity feature $T$ for any given nucleus $n$ by Equation 3.1. In this model the intensity of each pixel $p$ is the sum of multiple fluorescence sources, such as non-specific fluorescence ($F_{\text{nonspecific},p}$, e.g. due to off-target antibody binding), non-signaling fluorescence ($F_{\text{nonsignaling},p}$, e.g. due to a pool of the target protein that is not involved in the studied signaling process), and the actual signaling fluorescence $F_{\text{signaling},p}$.

$$T_n = \sum_p \left( F_{\text{nonspecific},p} + F_{\text{nonsignaling},p} + F_{\text{signaling},p} \right) \tag{3.1}$$

Take the case of β-catenin as an example. Basal β-catenin levels within Wnt3A-unstimulated cells are thought to be essentially zero (Section 2.3.4), and yet the measured basal levels are quite high. This is due in large part to the presence of a membrane-associated β-catenin pool that does not participate in Wnt signaling. It is impossible to know how the total nuclear intensity breaks down into the various components of Equation 3.1, and therefore a metric like fold-change is uninterpretable in terms of the extent of change for β-catenin, though it can be used as a normalization method to allow for measurements of relative change.

An alternative metric is to simply subtract a reference intensity measurement from the exper-

**Figure 3.10:** The cell lines tested in this dissertation have similar but non-identical timelines (**a**) and dose-response curves (**b**). Shown are SKMEL2 (solid lines) and HCEC (dashed lines) cell lines. Cell lines show measurable responses between 60-120 minutes and are saturated by 10ng/mL TGFβ3, 1000ng/mL Wnt3A, or 100ng BMP4. Dose-responses are at 1hr. Y-axes as in Fig. 3.8, normalized so that the minimum and maximum responses are 0 and 1. Concentrations used in **a**: TGFB3 (4ng/mL), Wnt3A (240ng/mL), BMP4 (70ng/mL).

imental intensity, as the only term remaining will be $F_{\text{signaling,experiment}} - F_{\text{signaling,ctrl}}$. While this metric has a simple interpretation (absolute change in signaling-associated fluorescence intensity), unfortunately it cannot be used to infer the absolute magnitude of response since fluorescence units are arbitrary. Neither division nore subtraction can be used at the single-cell level for fixed-cell assays, as the relative contribution of each fluorescent component may vary from cell to cell. Both metrics can be used at the population level, however.

In summary, it is impossible to infer the absolute magnitude of change for a signaling molecule using image-based immunofluorescence without making assumptions that would be indefensible for the immunofluoresce studies in this chapter. External means of validation are then required to show that the observed change is large enough to have a meaningful effect (as in Fig. 3.9, where I show that the same experimental treatments yield sizeable changes to target gene expression). Relative comparisons between experimental perturbations are possible in any case, and are easy to understand, and so for convenience I use units normalized to a 0/1 scale throughout this chapter.

### 3.1.4   Timepoints and concentrations

Many studies of morphogenic signaling pathways conflate the signal transduction and the transcriptional decision-making processes due to use of long-term transcriptional readouts. This conflation may be acceptable for some experimental questions, but here my aim is to study integration specifically at the level of signal transduction.

To do so, it is then necessary to minimize the effects of transcriptional feedback on the stimulated

pathways. I therefore performed time-course experiments for each cell line and pathway in order to identify the earliest timepoints that could still yield robust pathway responses (shown for HCEC and SKMEL2 in Fig. 3.10a). For all cell lines and pathways the responses were measurable by 1hr and maximal by 2hrs, and so my crosstalk experiments take place within this temporal window.

As I discuss in Section 1.5, it is not necessarily true that the use of a constant ligand concentration across multiple cell lines is the appropriate way to ensure that the treatment is the "same." Fortunately, the dose-response curves for each cell line and pathway are not wildly different (shown for HCEC and SKMEL2 in Fig. 3.10b) and so I was able to determine concentrations that would yield approximately half-maximal responses (EC50) or maximal responses across all cell types.

The reader may notice the use of differing ligand concentrations between the experiments in this chapter. For initial experiments, I did not yet know how responsive the cell types were to each ligand, and so concentrations were based on prior work in the lab or on literature-obtained values. For later experiments, doses were first chosen according to whether a half-maximal or maximal response was experimentally required, and these doses were then kept high over time to maintain saturating levels in the face of slowly-degrading reagents and idiosyncratic sensitivities of cell lines. Importantly, saturating concentrations minimize the consequences of pipetting error, since larger error can be tolerated before a measurable difference in cellular responses will appear.

## 3.2 Signaling integration between Wnt and TGFB$_{\text{sf}}$

Up to this point, the content of this dissertation has all been designed to set up the question: To what extent do Wnt and TGFB$_{\text{sf}}$ interact during signal transduction, prior to nuclear entry? As discussed throughout Chapter 2, these pathways have extensive opportunity for interaction and are widely believed to integrate information at the level of transcription. Further, despite a lack of shared core components, these signaling pathways have been tied together by numerous studies showing cross-pathway protein-protein interactions (Section 2.5). What is still unclear, however, is whether these interactions take place and have functional consequences to signal transduction under endogenous levels of pathway components.

### 3.2.1 Wnt3A and TGFβ3 show complete signaling insulation

To measure the signaling crosstalk between TGFβ and canonical Wnt, I made use of the experimental strategy described in the previous section by treating cells with combinations of each ligand. If these pathways were to interact in a manner that affects signaling (i.e. in a manner that transfers information) then the direct outcome of signaling (that is, nuclear canonical transcription factor levels) would be affected.

To test this, I treated all four cell types with combinatorial inputs of TGFβ3 and Wnt3A, and then measured nuclear accumulation of the transcription factors Smad2/3 and β-catenin by single-cell image analysis. Data for all four cell lines are shown for 1, 2, and 18hr timepoints in (Fig. 3.11). This is a lot of data, and so for simplicity I refer the reader to Fig. 3.12, which shows only the

**Figure 3.11:**  Wnt and TGFβ are insulated during signal transduction.  **a**, Wnt3A causes little or no modulation of Smad2/3 responses at both short (1-2hr) and long (18hr) timepoints for all cell lines tested. **b**, TGFβ3 causes little or no modulation of β-catenin responses, except in the case of long-term treatment in HCECs (note the 18hr TGFβ3 and TGFβ3+Wnt3A responses). Y-axes and p-values as in Fig. 3.8, with normalization per timepoint to set the control to 0 and the canonical-input-only condition to 1. These values fixed by normalization are in gray. $n = 3$ replicates per point. p-values indicate whether each condition differs from control ($\neq$ctrl) or from the canonical-input-only condition (either $\neq$TGFβ3 or $\neq$Wnt3A). Concentrations (1 and 2hr): 0.2ng/mL TGFβ3, 100ng/mL Wnt3A. Concentrations (18hr): 10ng/mL TGFβ3, 200ng/mL Wnt3A.

essential data for HCECs and SKMEL2s.

I performed the initial experiment at 1 and 2hr timepoints, using ligand concentrations that ranged from EC50 to saturating across the cell lines. In all cases, Wnt3A had small or statistically insignificant effects on Smad2/3 levels, even when co-treated with TGFβ3 (Fig. 3.11a, 1 and 2hr timepoints; Fig. 3.12a, top).  The small Wnt3A-induced Smad2/3 increases may be real, but are likely due to the trace contamination discussed earlier (see Fig. 3.6). The same absence of signaling integration occurs in the other direction, from TGFβ3 to β-catenin (Fig. 3.11b, 1 and 2hr timepoints; Fig. 3.12a, middle). In this case, the presence of TGFβ3 had no significant effects on β-catenin levels in any context. Therefore, Wnt3A and TGFβ3 are completely insulated during signal transduction (Fig. 3.12a, bottom).

Because the literature is full of examples of these two pathways interacting over longer timescales, I decided to repeat the experiment with a more distant timepoint. To ensure that the absence of crosstalk was not due to low activity of the pathways. To my surprise, even at 18 hours there was almost no measurable modulation of the transcription factor activity of one pathway by the other (Fig. 3.11, 18hrs). There was one exception, however: HCECs show TGFβ3→β-catenin interaction

**Figure 3.12:** Summary figure for insulation between Wnt3A and TGFβ3, limited to SKMEL2s and HCECs at 2 and 18hrs. **a,** At 2hrs, there is complete signaling insulation between Wnt3A and TGFβ3. **b,** However, context-dependent transcriptional integration already occurs at the same timepoint: HCECs show inhibition of Axin2 mRNA expression by TGFβ3 treatment, while SKMEL2s show complete transcriptional insulation. **c,** This context-dependent effect shows up at the level of signaling hours later, in HCECs. The increase in β-catenin due to TGFβ3 signaling likely stems from the inhibition of Axin2 mRNA (since Axin2 is a negative auto-regulator of Wnt3A). Thus, signaling insulation (**a**) combined with transcriptional integration (**b**) leads to a new, biased signaling state over time (**c**). Daggers indicate significant departure from pathway insulation. Data for **a** and **c** from Fig. 3.11. Data for **b** from Fig. 3.13.

(Fig. 3.12c, middle). Indeed, TGFβ3 treatment alone was sufficient to activate β-catenin, and co-treatment with Wnt3A yielded an approximately additive effect.

The complete absence of cross-pathway transcription factor modulation at early timepoints strongly suggests that the Wnt3A and TGFβ3 signaling cascades are completely insulated from one another, displaying no signaling crosstalk whatsoever (Fig. 3.12a, bottom). Further, the frequent absence of cross-pathway modulation even given significant time for transcriptional feedback shows that pathway insulation can be maintained even after the transcriptional network has been remodeled by morphogenic signals. Finally, the fact that HCECs lose this insulation at later timepoints demonstrates a case of context-dependency (i.e. dependency on cell type) in cellular decision-making despite context-independent insulation of signal transduction (Fig. 3.12c, bottom).

**Figure 3.13:** Wnt and TGFB show context-dependent transcriptional crosstalk. **a**, In the melanoma cell lines, TGFβ3/BMP4 have no effect on Axin2 expression, while Wnt3A treatment has no effect on Smad7 expression. 2hr treatment. **b**, In HCECs, TGFβ3 reduces baseline and blocks Wnt3A-induced Axin2 expression. This effect is already apparent at 2hrs and is maintained at 6hrs. Wnt3A does not affect Smad7 expression in any context. Bold '+' indicates doubled TGFβ3 concentration, demonstrating transcriptional saturation of this pathway. Fold-change is relative to the control condition. Mean and standard deviation of 3 replicates, '*' indicates two-tailed p-value <0.05 (Student's t-test). Concentrations: 10ng/mL TGFβ3, 25ng/mL BMP4, 200ng/mL Wnt3A. See Methods for qPCR details.

### 3.2.2 Wnt and TGFB_sf show context-dependent transcriptional integration

Due to the surprising result that the direct outcomes of signaling (transcription factor levels) are completely insulated between Wnt3A and TGFβ3, I decided to test for insulation at the level of transcription. By measuring mRNA expression 2hrs after treatment, I reasoned that I could identify whether insulation at the level of signaling (as already demonstrated) necessarily implies insulation at the level of transcription. Because Axin2 and Smad6/7 are the only prototypical transcription targets of the TGFB_sf and Wnt pathways, I measured mRNA levels of these genes following combinatorial ligand treatment (Fig. 3.13; simplified in Fig. 3.12b).

By qPCR, neither of the melanoma cell lines show transcriptional crosstalk between the canonical TGFB_sf and Wnt3A outputs Smad7 and Axin2 (Fig. 3.13a), implying that these pathways are insulated both in terms of the quantity of the transcription factors sent to the nucleus (as shown in Fig. 3.11) and in the general activity of those transcription factors. (This is not to say that these transcription factors do not affect one another for any transcriptional targets.) However, as before, HCECs display a different behavior. By 2 hours HCECs show strong modulation of Axin2 expression by TGFβ3 treatment(Fig. 3.12b), and this effect increases over time (Fig. 3.13b).

The modulation of Axin2 by TGFβ3 in HCECs is intriguing for several reasons. First, it provides a simple explanation for the 18hr transcription factor modulation result observed in (Fig. 3.12c),

**Figure 3.14:** The TGFβ and BMP signaling pathways are additive, and do not compete for Smad4. **a**, BMP4 causes no modulation of Smad2/3 responses at both short (1-2hr) and long (18hr) timepoints across all cell lines tested. **b**, TGFβ3 additively modulates pSmad1/5/8 responses to BMP4. **c**, Measurement of TGFβ3/BMP4 crosstalk by Western is consistent with the imaging data. SKMEL2, 2hr treatment. **a,b**, Y-axes, normalization, and p-values as in Fig. 3.11. Concentrations (1 and 2hr): 0.2ng/mL TGFβ3, 5ng/mL BMP4. Concentrations (18hr): 10ng/mL TGFβ3, 25ng/mL BMP4. Western courtesy Curtis A. Thorne (Altschuler & Wu lab, UT Southwestern).

since repression of Axin2, a negative autoregulator of β-catenin, could block the negative feedback otherwise present after Wnt3A stimulation. Second, it clearly shows that HCECs can simultaneously display signaling insulation and transcriptional crosstalk, demonstrating that these two processes can be completely independent. Finally, the general lack of crosstalk across all cell lines, coupled with the single instance of transcriptional crosstalk in HCECs, is suggestive that the oft-cited idiosyncratic outcomes of Wnt/TGFB$_{sf}$ crosstalk are predominantly due to context-dependent transcriptional crosstalk.

## 3.3 Signaling insulation between BMP and TGFβ

The above result is perhaps not so surprising, that two pathways (Wnt3A and TGFβ3) lacking any shared core components do not modulate one another during signal transduction. Under this rationale the BMP2/4 and TGFβ1/3 pathways might then be expected to interact, given their

**Figure 3.15:** Summary figure showing the lack of Smad4 competition between BMP4 and TGFβ3, limited to SKMEL2s and HCECs at 2hrs. **a,** There is additive signaling crosstalk from TGFβ3 to pSmad1/5/8, but not from BMP4 to Smad2/3. **b,** At the same time, the transcriptional output from the combined pathways also appears to be additive. **c,** Smad4 RNAi reduces protein levels of Smad4 in HCECs. Histone H3B serves as a loading control. **d,** Smad4 RNAi in HCECs reduces overall TGFβ3 responsiveness (top) but not pSmad1/5/8 responsiveness (middle), while the signaling crosstalk between TGFβ3 and BMP4 remains approximately additive. Thus, competition for Smad4 does not cause cross-pathway signaling inhibition between TGFβ3 and BMP4. Normalization for **d** uses the control and single-ligand responses from the scramble siRNA treatment to define 0 and 1. Y-axes, normalization, and p-values as in Fig. 3.11. Daggers indicate significant departure from pathway insulation. Data for **a** from Fig. 3.14. Data for **b** from Fig. 3.13. Western courtesy Curtis A. Thorne (Altschuler & Wu lab, UT Southwestern).

shared requirement for Smad4. Indeed, the claim of intra-TGFB$_{sf}$ crosstalk via Smad4 competition has been cited in many reviews and papers, but to my knowledge has not been directly tested (Section 2.2.5). I therefore decided to use the same experimental setup as above to measure the extent of crosstalk between BMP4 and TGFβ3 (Fig. 3.14). The essential data is again summarized in a simpler figure (Fig. 3.15). If these pathways do compete for Smad4, then co-treatment with saturating concentrations of both ligands (to maximize sequestration of Smad4) should cause one or both pathways to be attenuated.

BMP4 treatment had absolutely no effect on Smad2/3 in any cell line (Fig. 3.14a; Fig. 3.15a,top). There is crosstalk in the other direction: treatment by TGFβ3 is sufficient to activate pSmad1/5/8, and the presence of both ligands yields an additive behavior (Fig. 3.14b). Importantly, these two pathways also have an approximately-additive effect on transcription of their shared downstream target, Smad7 (Fig. 3.13; Fig. 3.15b,top). As noted in Section 2.2.4, while the BMP2/4 and TGFβ

pathways are generally considered to be distinct, they have been shown to cross-activate one another in various settings. Western blotting was not quantitative enough to confirm the activation of pSmad1/5/8 by TGFβ3, but is consistent with an absence of negative cross-regulation (Fig. 3.14c). In any case, both the insulation of Smad2/3 from BMP4 and the additive crosstalk between TGFβ3 and pSmad1/5/8 directly argue against competitive inhibition between these pathways during signal transduction.

I then wondered if the absence of inhibitory crosstalk between BMP4 and TGFβ3 was a context-dependent phenomenon. For example, the additive behavior is consistent with the quantity of Smad4 being so high as to be non-limiting, in which case both pathways would effectively have distinct pools of Smad4 (as in Fig. 3.1c). I therefore used siRNA to knock down Smad4 in order to make Smad4 a limiting factor in an effort to modulate the form of crosstalk (e.g. to a non-additive form as in Fig. 3.1d).

As a consequence of Smad4 depletion (Fig. 3.15c), TGFβ3 signaling was strongly reduced overall but was still not inhibited by co-treatment with BMP4 (Fig. 3.15d, top). Levels of pSmad1/5/8, on the other hand, were not affected overall by Smad4 knockdown (Fig. 3.1d, top) and the effect of co-treatment remained roughly additive. Therefore, in contrast to expectations, BMP4 and TGFβ3 do not negatively regulate one another at all at the level of signal transduction, and do not compete for Smad4.

## 3.4    Discussion

In this chapter I measured the extent of crosstalk between key morphogenic pathways, using a mechanism-independent and endogenous approach that kept signaling crosstalk separable from transcriptional crosstalk. By doing so, I discovered that canonical Wnt and TGFβ do not crosstalk at all during signal transduction, and that TGFβ and BMP do not compete for Smad4 during signaling. Both of these results yield important simplifications to the ever-increasing apparent complexity of signal transduction that stems from discoveries of putative nodes of crosstalk. Below, I discuss specific implications in more detail.

### 3.4.1    Morphogenic signaling is insulated

As I explain in Section 2.6, crosstalk during signal transduction is likely to decrease the accuracy of the intra-nuclear model of the extracellular environment. In such a case, the information that eventually passes to the nucleus for transcriptional processing represents a simplified view of the cellular microenvironment, such that the nucleus then has to make a decision with more uncertainty about the state of the outside world. By instead allowing information to pass relatively unfiltered from the extracellular milieu to the nucleus (e.g. by insulating information channels from one another) the nucleus has access to a more accurate model of the microenvironment and, presumably, can then make more useful decisions.

The results presented in this chapter suggest that the morphogenic Wnt and TGFB$_{sf}$ pathways

are indeed isolated from one another, such that they can each pass information to the nucleus without cross-pathway interference. As each pathway is only capable of sending a small amount of information regarding ligand concentration (not shown, but one can infer this from the broad distributions in Figs. 3.7 & 3.5), maintaining this information from distinct pathways may be the only way that a cell can obtain an accurate internal model of its environment. Therefore, I propose that signaling insulation may be a general property of morphogenic pathways whose key outputs are changes to the transcriptional network.

Signaling insulation allows for the internalization of a somewhat unfiltered view of the cellular microenvironment (Fig. 3.12a, cartoon). Context-dependent transcriptional feedback (Fig. 3.12b, cartoon) can then modulate the signaling pathways over time (either by auto-regulation or cross-pathway regulation) (Fig. 3.12c, cartoon). The result of this nuclear decision-making would be to essentially bias the cell's view of what could even be an unchanging environment. Thus, by mixing insulated signaling with long-term feedback, cells can maintain a relatively complete internal model of the environment but interpret that environment in a temporally biased manner.

An obvious argument against general morphogenic insulation in my own data is that TGFβ3 appears to activate the BMP-specific Smads (Smad1/5/8). I note that I was unable to verify that this crosstalk is specific (i.e. not due to a contaminant) as addition of an anti-TGFβ antibody along with TGFβ3 treatment did not significantly block the interaction (data not shown). For this reason, I only interpret the BMP/TGFβ crosstalk data to say that one does not inhibit the other.

An additional counter-argument is that these pathways may in fact interact, but I have measured the wrong thing to be able to identify the interactions. This would be a useful discovery, as what my work shows is that the generally-accepted method of encoding used by these pathways shows complete insulation. It seems to me quite possible that there exists some information channel, at the level of signaling, that does indeed integrate information from both Wnt and TGFβ signals.

### 3.4.2  BMP and TGFβ do not compete for Smad4

While it is widely believed that co-activation of the BMP and TGFβ pathways should result in some sort of cross-pathway inhibition via Smad4 competition, I find no evidence of this phenomenon. Indeed, I am unaware of any studies that explicitly show this form of crosstalk, and some studies have shown that Smad4 is highly abundant [92, 97].

Interestingly, my data show that even when Smad4 is brought down to limiting levels, there is still no cross-pathway inhibition. Smad4 is seen as the factor that allows receptor-Smad (rSmad) access to the nucleus, and it is generally believed that this interaction is stoichiometric and non-transient [96]. If this were the case, how could limitation of Smad4 selectively ablate one arm of Smad responses but not the other?

One possibility is that my observed difference between an ablated Smad2/3 and maintained pSmad1/5/8 responses (Fig. 3.14e) is due to some difference between the behavior of the bulk protein on one hand versus the phosphorylated form on the other. However, such an effect would be difficult to reconcile with the standard model that the phospho-state allows Smads to stay in the

nucleus (i.e. the nuclear phospho- and total-protein levels should vary together).

What is, to me, a more likely explanation is that one of the consequences of ablating Smad4 signaling, wich requires 48 hours of transfection, is a change to the transcriptional network with respect to $TGFB_{sf}$ pathway components. For example, loss of Smad4 may have resulted in a decrease in the quantity of TGFβ receptors, co-receptors, or even an increase in secreted antagonists. I did not test this, and therefore cannot make a conclusion with respect to the differential effect of Smad4 RNAi on Smad2/3 and pSmad1/5/8. Therefore, the most reasonable takeaway from this data is simply that reduction in Smad4 does not force pathway crosstalk.

The lack of competition for even low levels of Smad4 leads to a potential modification of the current model of Smad nucleo-cytoplasmic shuttling (see Section 2.2.4). It may be that interactions of the rSmads with Smad4 are more transient than is currently believed, and that association with Smad4 is not required for nuclear receptor-Smad activity and maintained localization. In such a case, Smad4 could behave more like an enzyme than a stoichiometric scaffold, in that single Smad4 molecules could mediate the nuclear translocation of numerous rSmads. If this process were fast enough, then the two classes of rSmads would effectively have access to different pools of Smad4; activity of one pathway would not "soak up" the co-Smad even at pathway saturation.

### 3.4.3 Future directions

While this work has demonstrated that morphogenic pathway integration may not be as complicated as we think, there are many unanswered questions. Perhaps the most straightforward questions to address are on the generality of morphogenic insulation. Does insulation extend to other classic morphogenic pathways (such as Hedgehog and Notch)? Does it extend to yet more diverse cell lines than those tested here? I suspect that signal transduction insulation is a general principle, and future work using the approaches in this chapter could provide the answers to these questions.

It would also be informative to study the kinetics of crosstalk. For example, in the case of HCECs, which show cross-pathway modulation of transcription factors after 18 hours (Fig. 3.11), one could measure how much time is required after the signal this crosstalk becomes apparent. Such data could be used to determine how long it takes for a cell to build a new, biased model of its environment, and how stable that biased model is in the presence of either a constant or changing signal.

Throughout this dissertation I have been harping on this concept of information transfer, as we must always make assumptions regarding which cellular and protein properties are carrying signaling information. My data for Wnt and $TGFB_{sf}$ suggest that these pathways carry relatively little information about absolute extracellular ligand concentrations. This lack of information implies either that cells are terrible concentration detectors or that we are not looking at the right encoding relationship between these ligands and the internal cellular model of those ligands. A comprehensive study designed to sort out the sources of information, and the precise intracellular properties into which that information is encoded, will be essential to our understanding of how cells make decisions as a result of $TGFB_{sf}$ and Wnt signals.

**Table 3.1:** Recombinant proteins used in this chapter. Vendors: CST, Cell Signaling Technology Inc (Danvers, MA); Life, Life Technologies (Grand Island, NY); R&D, R&D Biosystems (Minneapolis, MN).

| Protein | Vendor | Catalog # | Lot # |
|---------|--------|-----------|-------|
| BMP4 | CST | 4697 | |
| DKK1 | R&D | 5439-DK | SMR2713041 |
| Noggin | R&D | 6057-NG | |
| TGFβ1 | Life | PHG9204 | |
| TGFβ3 | CST | 8425 | |
| Wnt3A-LP | R&D | 5036-WN | RSK31 |
| Wnt3A-HP/CF | R&D | 5036-WNP/CF | SVH0813081 |
| Wnt5A-LP | R&D | 645-WN | |
| Wnt5A-LP/CF | R&D | 645-WN/CF | MCR4513111 |

## 3.5   Methods

The rationale and general methodology for image correction and analysis are provided in Chapter 4. Specific experimental details are provided in the image captions; this section provides information that is broadly applicable or more detailed than is appropriate for a figure legend.

**Cell culture.** I maintained all cells in RPMI1640 (Cellgro #10-040) with 5% FBS (Gemini Bio-Products #100-106) with antibiotics/antimycotics under standard tissue culture conditions (37°C, 5% CO2). In preparation for imaging, I plate $\sim 10^4$ cells per well of a 96-well plastic plate (Corning #353219) and a fifth of that for 384-well glass plates (Nunc #164586). I use DMEM (Gibco #11965-126) at the time of seeding when starvation conditions are needed. In either case, cells are left to adhere overnight before being treated the following day. For treatments, I dilute recombinant protein into the same media that the cells are grown in. The proteins used for treatment are listed in Table 3.1. Most of the plates used in this chapter are 384-well glass plates. Human colonic epithelial cells were a kind gift from Dr. Jerry Shay (University of Texas Southwestern, Dallas TX) [213]. The other cell lines were obtained from American Type Culture Collection (Manassas, VA).

**Gene expression.** I plated cells in 96-well plastic plates as above and left them to adhere overnight before treatment the following day. Each treatment was performed in triplicate wells. I used an Ambion Cells-to-CT kit (Life Technologies) to extract mRNA according to manufacturer instructions, and passed these samples on to the UT Southwestern Medical Center Microarray core facility for cDNA library preparation and TaqMan qPCR. The samples were given obfuscating identifiers and within-plate positions so as to blind the core facility to the treatments and replicates. TaqMan probes (Applied Biosystems) comprised Smad7 (Hs00998193_m1), Axin2 (Hs00610344_m1), and 18S rRNA (Hs99999901_s1). The core facility returned the raw threshold cycle ($C_t$) values for each probe/sample combination, which I first internally normalized by 18S rRNA levels to yield $C_{t(probe,norm)} = C_{t(probe)} - C_{t(18S)}$ and then converted these values to fold-change over control.

**RNA interference.** I used a Dharmacon Smad4 siGENOME SMARTpool (GE Life Sciences, #M-003902-01) to ablate Smad4 in HCEC cells. For transfection I used Lipofectamine RNAiMax (Life Technologies) according to manufacturer instructions. The transfection media was left on cells for 48hrs, after which I replaced the media with fresh media for >2hrs prior to experimental treatment.

**Western blots.** SDS-PAGE and western blotting were performed using standard techniques.

**Table 3.2:** Antibodies used in this dissertation. Vendors: CST, Cell Signaling Technology Inc (Danvers, MA); BD, BD Biosciences (San Jose, CA); Abcam (Cambridge, MA); Life, Life Technologies (Grand Island, NY); R&D, R&D Biosystems (Minneapolis, MN). The indicated dilutions are for immunofluorescence unless appended with a 'w' for Westerns. The appendix 'b' indicates use as a blocking antibody.

| Target | Dilution | Source | Vendor | Catalog # | Lot # |
|---|---|---|---|---|---|
| Smad2/3 | 1/1000 | Rabbit | CST | 8685 | 3 |
| pSmad1/5/8 | 1/100 | Rabbit | CST | 9511 | 8 |
| β-catenin | 1/100 | Mouse | BD | 624084 | |
| Smad4 | 1/100 | Rabbit | CST | 9515 | 4 |
| Smad4 | 1/100 | Mouse | Abcam | 3219 | GR94411-1 |
| Smad1 | 1/100 | Rabbit | CST | 6944 | 2 |
| Smad2 | 1/100 | Rabbit | CST | 5339 | 4 |
| pSmad2/3 | 1/100 | Rabbit | CST | 8828 | |
| H3B | 1/1000w | Rabbit | CST | 9715 | |
| TGFβ1-3 | 1/1000bw | Mouse | R&D | MAB1835 | CCI1512031 |
| AlexaFluor 546 Anti-Mouse IgG (H+L) | 1/1000 | Goat | Life | A11003 | 1256168 |
| AlexaFluor 488 Anti-Rabbit IgG (H+L) | 1/1000 | Goat | Life | A11008 | 1470706 |

I plated cells in 6-well dishes and treated them the following day. After treatment, I washed the wells with ice-cold PBS and then lysed with RIPA buffer (50 mM Tris (pH 8.0), 150 mM NaCl, 1% NP40, 0.5% deoxycholic acid, 0.1% SDS, 0.5 mM EDTA) containing protease and phosphatase inhibitors (Sigma Aldrich). Antibodies used are shown in Table 3.2, generally with 1:1000 dilutions (the CST rabbit Smad4 antibody was used for blotting).

**Immunostaining**. All solutions are made in PBS (Gibco #70013). All wash steps are repeated 3 times with 0.1% Tween20 (Fisher #BP337). Antibodies are diluted into 2.5% BSA (Fisher #NC9871802). Cells are fixed in 4% paraformaldehyde (Electron Microscopy Sciences #15710) for 10min, permeabilized with 0.2% Triton X-100 (Sigma #93443) for 10min, then washed. I then incubate the samples overnight at 4°C with appropriate primary antibodies (Table 3.2). After washing, I secondary-stain for 2hrs at room temperature with added 2.5μg/mL Hoechst. Finally, I wash the samples once more before storing them in PBST for imaging. For all experiments, I prepare 6-18 uniformly-fluorescing wells for measuring image shading, using the same dissolved secondary antibodies and Hoechst.

**Imaging.** I imaged all stained plates on a Nikon Eclipse Ti-E2000 microscopes controlled by NIS Elements version 4, with an Andor Zyla sCMOS 11-bit camera. I wrote custom image coordinate-generating software in Python for increased stage precision. To obtain the signal from the camera alone, I take detector images at low exposure times with no light source.

**Image Correction.** The background and shading correction followed the pipeline described in Section 4.3.3 using custom Matlab software.

**Segmentation and quality control.** For all analyses, I wrote a custom Matlab threshold segmentation algorithm to automate detection of nuclei using the Hoechst fluorescence channel. I manually set filters for size and shape, by cell type, to remove objects that are likely to be artifacts. For each nucleus, I then extract its area (in pixels) and the total of all pixel intensities for all imaged channels. Finally, I follow the DNA-based quality control and single-cell regression-based correction described in Section 4.5 using custom R software. The resulting high-quality G1-phase cells are used for analysis.

**Mutual information measurement.** I implemented the unbiased mutual information algorithm as described in [16] and discussed in Section 4.4.4, using a custom Matlab script.

# Chapter 4

# On quantitative imaging of single cells

## 4.1  Introduction

Fluorescence microscopy is a powerful tool for measuring single-cell properties, especially when sub-cellular spatial resolution is required. With modern computing power and robotics, we can now use this technology to amass large quantities of image data and so study cell biology with increasing breadth and depth. Accurate high-throughput imaging and image analysis of single-cell data across large numbers of conditions is becoming routine in some labs and will continue to become even more commonplace as the technology simultaneously improves and becomes cheaper. The availability of the technology has led, in recent years, to a boom in large-scale microscopy studies of single-cell phenotypes [27, 216–221]. However, our ability to generate microscopy data is outpacing our collective ability to analyze and interpret it.

Fluorescence microscopy has been an experimental staple of biology for many years, though it is still infrequently used for rigorous quantitative analysis. More often, researchers use fluorescence imaging to demonstrate a particular phenotype (as I do in Fig. 3.2), or to manually count the instances of a visually obvious phenotype. Our own visual perception is notoriously prone to error, however, as our brains can choose to perceive patterns that do not exist, thus allowing for latent biases to affect how we manually tally the data. As biologists we are often painfully aware of this problem, and so there seems to be a general unease when it comes to evaluating imaging data. We all know the joke that, in published figures, "representative image" is a euphemism for "the best image I could find."

The distrust is understandable, as proper evaluation of imaging data is a non-trivial task, and researchers rarely publish enough details to even allow for thorough evaluation in the first place. This problem is made worse by a lack of imaging and analysis standards, and by the difficulty inherent to making large image-datasets publicly available.

Biological image quantification is a difficult and generally unsolved problem [222], and its various approximate solutions tend to require a level of expertise in mathematics and computer programming, and at least a cursory understanding of microscopy optics. Importantly, it also requires expertise in the experimental biology under study. Few, if any, biology training programs prepare

students for such broad practical knowledge.

Cell biologists, like myself, are rarely trained in quantitative fields and so we tend to rely on our colleagues in analytical fields to do the quantitative work for us. Those analysts, in turn, typically lean back on the biologists for understanding the "what" and "why" of the analysis. Unfortunately, the large knowledge and language asymmetry between biologists and analysts creates steep communication barriers that are hard to break down. Therefore, the full benefits to cell biology of careful quantitative microscopy may be going untapped.

My own work relies almost entirely on imaging (Chapter 3) and, having spent so much time thiking about and performing image analysis, I have come to prefer imaging over other methods not only for its efficiency in measuring single-cell properties across large numbers of conditions, but also because it provides a literal view into the mysterious and beautiful world of the cell. I therefore hope that careful single-cell image analysis will someday become as commonplace as Western blots.

My overarching goal in writing this chapter is to provide an image analysis resource that is both accessible to classically trained biologists and useful for analysts from non-biological fields that may be preparing to move into the field. In general, I aim to provide careful biological and analytical reasoning for a subset of the many choices that must be made with respect to fluorescence imaging experiments. In specific, I focus on dealing with experimental and imaging artifacts and noise, and how to choose biologically meaningful single-cell measurements. The methods and rationale in this chapter are used extensively in Chapter 3.

## 4.2 Images as layers of fluorescence

The goal of quantitative biological imaging, in overly general terms, is to perform a set of meaningful mathematical operations on fluorescent images. We therefore require a mathematical model of fluorescence images to use as a reference when discussing image correction and analysis. Note that the model I construct below is designed to be intuitive with respect to cell culture microscopy, and so it differs slightly from the more general models that it is founded upon [223–226].

A fluorescence microscopy image $I$ can be thought of as a series of fluorescing layers, each with its own distinct properties, that add together to create the final image (Fig. 4.1). In this simple model, a fluorescence image is composed of multiple **foreground** and **background** layers that come from distinct components of the sample. (As described below, I use specific definitions of these terms "foreground" and "background" that may differ from one's intuition.) It is important to note that additivity is a special property of fluorescence microscopy, as bright-field and other non-fluorescence signals do not necessarily behave in this way [225]. Further, there are cases where the layers of fluorescence will become non-additive, as the fluorophores and the camera will display non-linear behaviors in certain ranges [227]. The model I describe here, and all analysis in this chapter, assumes that all contributing components are behaving within their linear ranges.

**Figure 4.1:** A fluorescence image $I$ of a cell is the sum of distinct fluorescence layers. Here, $F_1$ is the foreground signal of interest, $F_2$ is non-specific staining within the cell, and $B$ is background fluorescence from the imaging surface.

**Experimental image layers $F$ and $B$**

In this dissertation, foreground ($F$) refers to any image layer that emits fluorescence in a spatially non-uniform manner within an image; the layer is not "flat." The most useful foreground layer is due to the specific binding of a fluorescent probe to its target, as this layer is likely the one that is under study. Where my definition of "foreground" may diverge from others is that I include spatially non-uniform fluorescence artifacts as foreground layers. Such artifacts might include non-specific cellular staining, cellular autofluorescence (as for layer $F_2$ in Fig. 4.1), and staining artifacts such as halos, bubbles, or bright puncta. Many biologists refer to these artifacts as "background", but I have a specific definition for this term as well.

Also specific to this dissertation, "background" ($B$) refers to any image layer that emits fluorescence in a spatially uniform manner within an image (e.g. layer $B$ in Fig. 4.1). In other words, the layer is "flat," aside from variation between pixels due to measurement error. Background layers may include reflections from the imaging surface, autofluorescence from unbound fluorophore in the solvent, or fluorophore that has adhered to the imaging surface.

Finally, all fluorescent image layers scale linearly with excitation light intensity and exposure time (again, so long as we are in the linear range for all image components). For convenience we can fold excitation intensity and exposure time into a single term $t$, which I refer to simply as "exposure." Taken together we get the (yet-incomplete) image model in Equation 4.1 that contains $n$ foreground and $m$ background layers.

$$I = t \left( \sum_{i=1}^{n} F_i + \sum_{j=1}^{m} B_j \right) \tag{4.1}$$

**Image modification by the microsope**

The model in Equation 4.1 describes the fluorescing layers of an image, but there are also non-fluorescent properties of images. In digital fluorescence microscopy, the camera typically has some non-zero baseline value that I refer to as the "detector value" $D$. This value is a constant and does not change with exposure [228].

**Figure 4.2:** A fluorescence image can be modeled by $I = D + S(F + B)$, where $D$ is the baseline detector value, $S$ is shading, $F$ is the foreground, and $B$ is the background. Here, synthetic images of two foreground objects demonstrate each of these components. Plots indicate pixel values along the white line drawn across the image as i **a**. **a**, Two foreground objects that vary 2-fold in intensity are shown. **b**, Addition of background reduces the apparent fold-difference between the two foreground objects, as does addition of the detector (**c**). **d** Shading distorts the foreground and background, but leaves the detector contribution unchanged.

Uneven illumination of the image by the light source, which I refer to as "shading" ($S$), is an optical problem inherent to microscopy [227, 229–231]. Shading is a consequence of how light moves through lenses and so it is not just a property of wide-field microscopy, though this is a common misconception: shading is also present in confocal and total internal fluorescence (TIRF) microscopy [231, 232].

Shading can be interpreted as variation in relative exposure as a function of position within the image. In other words, $S$ modifies $t$ in a pixel-coordinate dependent manner. Because fluorescence units are typically arbitrary, in that their absolute values carry no meaning, I simplify the model further by setting $t = 1$ and dropping it from the equation (note that analyses making use of varied exposure times could also collapse $S$ and $t$ into a single term). By including the $S$ and $D$ components, and a noise term $\epsilon$ to absorb experimental error, we get the complete model in Equation 4.2.

I simplify the model further by collapsing all foreground and background layers into single terms. The resulting image model in Equation 4.3 is used throughout this text. This simplification is useful because the basic image correction and analysis methods presented in this chapter do not distinguish between foreground layers or between background layers. It is useful to keep the multi-layer model in mind, however, as it will help when trying to untangle the overall fluorescence behavior of experimental images. See Fig. 4.2 for a visual demonstration of the simple model.

$$I = D + S \left( \sum_{i=1}^{n} F_i + \sum_{j=1}^{m} B_j \right) \epsilon \qquad (4.2)$$

$$I = D + S(F + B)\epsilon \qquad (4.3)$$

**Figure 4.3:** An image $I_{r,c}$ is a matrix of intensity values, with pixel coordinates given by row $r$ and column $c$. Analyses can be performed on single pixels, the entire image, image neighborhoods, or image stacks as shown. In the case of an image stack, $z$ refers to the image position within the stack. Note that $z$ need not refer to a height, as in the common z-stacks of confocal microscopy, but can also indicate different color channels or even entirely unrelated images. In effect, $z$ is a "height" within the image stack that may not correspond to the physical height of an imaged within a sample.

### 4.2.1 Properties of the image components

With a model in place, it is useful to obtain some intuition for how to think about images in this framework. First, note that the image itself, and each layer that makes it up, is a matrix of intensity values (Fig. 4.3). Many analytical operations can be performed per pixel coordinate, in effect ignoring the presence of neighboring pixels. Other operations take into account those neighbors (this is especially true of segmentation, discussed later). Finally, we can think of a "stack" of images, in the same way one would stack a deck of cards. We can perform operations "down the stack" at particular pixel coordinate. Such operations include "per-pixel" means and medians.

Each component of the image model in Equation 4.3 is a matrix of the same dimensions, and each has distinct properties. The foreground components of $F$ have completely idiosyncratic values, both within and between images, as these values depend on where cells are, what is being stained, and what kinds of randomly-positioned artifacts are present. Indeed, the unknown behaviors of the foreground layers are what we typically aim to understand via image analysis.

The background layers of $B$ are defined to be more predictable, in that they are unchanging by position within an image. The background should be constant between images as well, when identical experimental conditions are used to obtain those images. However, the total background may change as a consequence of some experimental perturbation. For example, some small molecules used as drugs may fluoresce and so add an additional background layer. Finally, note that due to shading these "flat" background layers will appear distorted in uncorrected images. To clarify, then, I define background layers as those that are constant across an image in the absence of shading.

The detector layer $D$ has the simplest behavior, as it can be considered constant regardless of the imaging and experimental conditions. "Constant" in this case means that the value at any given pixel position $D_{r,c}$ does not change over time or between images. The values within an image, on the other hand, may vary (see Fig. 4.4a).

The shading value $S$ is problematic for analysis in that it distorts the $B$ and $F$ layers in non-

**Figure 4.4:** Example within-image variation of the detector $D$ and shading $S$ image components. **a**, Per-pixel averages of 10 detector images from two cameras are shown. Camera 1, Andor Zyla sCMOS 11-bit. Camera 2, Roper Scientific CoolSnap HQ2 CCD 14-bit. Intensities and images are scaled independently. **b**, Shading patterns from two optical channels. Uniformly-fluorescent background images were made with dissolved Hoechst or wheat germ agglutinin (WGA)-TRITC in the DAPI and TRITC optical channels. Note that each channel has a dramatically different shading pattern and overall degree of shading. Histograms show the all-pixel distributions of shading values, which are defined to have a median of one (a robust variant of $E[S] \equiv 1$, defined in Section 4.3).

trivial ways. Like $D$, this layer also shows variation within an image. The shading pattern caused by the objective lens tends to show brighter fluorescence at image centers and weaker fluorescence at the edges. However, the presence of other components in the light path, such as filters, can modify this shading pattern (Fig. 4.4b) [233]. Unlike $D$, the shading pattern can vary between images as well, though in Section 4.3 I show that this variation can be both predictable and correctable. $S$ distorts the $F$ and $B$ layers multiplicatively, and can cause as large as 1.5- to 2-fold intensity differences across an image [234].

Finally, we are left with the noise term $\epsilon$. I use this term specifically to capture measurement noise, not true biological variability. The largest source of measurement noise for modern fluorescence imaging is probably due to the combination of two factors: the error in converting photon counts to electrons, and the error introduced during transmission of those electrons from the camera. Combined, this measurement error introduces an intensity-dependent uncertainty in the total measured intensity of the pixel $I_{r,c}$. In general, the relative error decreases as a function of the square root of intensity [228].

## 4.3   On image correction

Fluorescence microscopy images contain non-biological data, distortions, and artifacts that may confound analytical results if not addressed. In the terms of the image model from the previous section (4.3), the subject of study in fluorescence imaging is contained within a subset of the foreground layers. All non-foreground components should thus be removed in order to obtain meaningful quantitative results. Ideally, the foreground layers that are not of interest (especially artifacts) should also be removed, though this is a much more difficult problem that I do not address in this dissertation. In essence, from the original image $I = D + S(F + B)\epsilon$, we need to obtain $F$: obtaining $F$ is the goal of image correction.

It is important to be aware that image correction is not a solved problem. There are many ways in which it can be done, several of which I review in Section 4.3.2, but the most commonly-used methods produce incomplete correction and/or are prone to generating artifacts. Further, methods sections of papers that rely on imaging often do not explain their correction methodology at all, making it difficult to evaluate some published results. I therefore make the importance of image correction a focus of this chapter, and describe an image correction approach that I developed for accurate determination of $F$ in the difficult context of high-throughput microscopy Section 4.3.3 [235].

Before going into the details of image correction, I should first address whether this step is even necessary. After all, any processing step could introduce artifacts and thus potentially do more harm than good. Normal cell-to-cell variability is already relatively high, with standard deviations of ~15-30% for protein concentrations [28]. One might then wonder whether error introduced by non-foreground image components would make much of a difference. Here I show that it can indeed make a difference. Unsurprisingly, the size of the difference (and thus the importance) is highly dependent on the properties of each particular image dataset, the measurement methods used, and the experimental goals.

### 4.3.1   Non-foreground components distort single-cell phenotypes

A particular cellular phenotype can be represented as a set of measurements. For example, a cell phenotype might be composed of its measured size, shape, texture, and average fluorescence intensities in multiple color channels. Each measurement is generally referred to as a "feature." The diversity of features that can be measured for a single cell is only limited by the imagination of the investigator. It would therefore be impossible to develop a single comprehensive argument for the importance of image correction that covers all possible ways of measuring cellular phenotypes. Instead, I make a case study of several commonly-used and biologically-interpretable single-cell features: the average, total, and ratios of fluorescence intensity.

How to quantify the effects of image correction on data is not obvious, since that data can be used in many ways depending on experimental goals. There are, however, aspects of single-cell distributions that are meaningful across a broad array of experiments, and so I use these as metrics when measuring the consequences of image correction. These are the mean ($\mu$), standard deviation

**Table 4.1:** Symbols used in the mathematical derivations of the consequences of $S$ and $B$ on single-cell feature distributions.

| Symbol | Meaning |
| --- | --- |
| $E[X]$ | Expected value (i.e. the mean) of the random variable $X$ |
| $F_c$ | Average foreground intensity within cell $c$ |
| $S_c$ | Average shading value within cell $c$ |
| $B$ | Background intensity per pixel, a constant |
| $\alpha_c$ | Area of cell $c$ |
| $A$ | Average intensity feature |
| $T$ | Total intensity feature |
| $R$ | Ratio of intensities feature |
| $Z_c$ | The value of feature $Z$ for a single cell, $c$. |
| $\sigma_Z$ | The standard deviation of feature $Z$ across all cells. |
| $\mu_Z$ | The mean of feature $Z$ across all cells. |
| $cv_Z$ | The coefficient of variation of feature $Z$ across all cells ($\sigma_Z/\mu_Z$). |

($\sigma$), and coefficient of variation ($cv = \sigma/\mu$) of a feature across a population of cells.

The standard deviation and $cv$ are measures of distribution widths, which are used to determine the statistical separability of distributions. Inaccurate measurements of the true variability may consequently reduce statistical power. The mean of a feature distribution, on the other hand, is typically used to determine how large of an effect an experimental perturbation has had. Inaccurate measurements of the true mean may cause strong results to appear weak, or vice versa, leading to false negatives or false positives. These distribution metrics are therefore useful as readouts for the utility of image correction.

The mathematics in the following discussion were worked out in conjunction with my co-authors on the relevant publication [235]: Satwik Rajaram, Chonlarat Wichaidit, Steven Altschuler, and Lani Wu. The text and figures draw heavily from the same publication. Those readers who do not need convincing that image correction is important may skip to page 85 without a loss in coherence of this chapter.

**Mathematical definitions of commonly used single-cell features**

Commonly used single-cell measurements include pixel intensity averages, totals, and ratios within some cellular compartment $c$ (such as the nucleus, cytosol, or whole cell). By defining these features mathematically we can determine their general behaviors as a consequence of the presence of image background or shading. Refer to Table 4.1 for the list of mathematical symbols, to Table 4.2 for a summary of the statistical properties used in the mathematical derivations, and Fig. 4.5 for a case study of these behaviors.

For each cellular object, I define $F_c$ and $B_c$ as the average foreground and background intensities within $c$, while $S_c$ is the average shading. $B_c$ is a constant for all cells, as the background is assumed to be the same for all pixels in the absence of noise, and so I drop the subscript from this term. I refer to the area of each cell, measured in pixels, as $\alpha_c$. Finally, for the derivations I assume that the detector value has been subtracted from all pixels and that the image contains no measurement

**Figure 4.5:** Image background and shading cause feature-dependent changes in estimates of phenotypic variability. Human colonic epithelial cells ($n > 3700$), stained for DNA (using Hoechst) and Smad (using a Smad2/3 antibody), imaged at 10X. The distributions of nuclear feature values were then compared before (black histograms) and after (dotted blue histograms) artificial background (top row, with background $\sim$16% of Hoechst or $\sim$50% of Smad foreground) or shading (bottom row, linear gradient with maximum 1.5 fold intensity difference) were added to each image. $\mu$, mean; $\sigma$, standard deviation; $cv = \sigma/\mu$. Inset, top left, relative size of change to the shown distributions. Inset, top right, arrows indicate the direction of change in the general case (question marks indicate uncertainty due to dependency on other variables). $x$-axes in arbitrary fluorescence units. A version of this figure is published as Fig. 1a in [235].

noise ($D = 0$ and $\epsilon = 1$).

For each cell I can then define the three simple intensity features: total intensity $T$, average intensity $A$, and the ratio of intensities $R$ between two independent foreground signals $F_{c1}$ and $F_{c2}$ (e.g. the ratio of nuclear Hoechst and Smad intensities, as in Fig. 4.5). Because the ratio takes two signals into account, each may come from a distinct fluorophore and optical setup and so have distinct shading and background values. Additionally, those features could be defined within distinct cellular compartments, such that the compartment sizes may also differ. The three features are thus defined for single cells by Equations 4.4-4.6.

$$A_c = S_c(F_c + B) \tag{4.4}$$

$$T_c = \alpha_c A_c = \alpha_c S_c(F_c + B) \tag{4.5}$$

$$R_c = \frac{\alpha_{c1} S_{c1}(F_{c1} + B_1)}{\alpha_{c2} S_{c2}(F_{c2} + B_2)} \tag{4.6}$$

For simplicity of the following analysis, I take the special case where $F_c$, $S_c$, $\alpha_c$, and $B$ are all statistically independent (i.e. cells do not spatially arrange themselves within an image by phenotype, and foreground intensity is independent of cell size). Further, I assume that cells are small relative to the spatial rate of change of $S$ across an image. These assumptions allow me to use the properties in Table 4.2. Note that these assumptions will be valid for some experimental

**Table 4.2:** Statistical properties used in the derivations of the effects of shading $S$ and background on the distributions of single-cell average $A$, total $T$, and ratio $R$ intensity features. $X$ and $Y$ are independent random variables, and $k$ is a constant.

| Index | Property |
|:-----:|:--------:|
| 1 | $\mathrm{E}[S_c] \equiv 1$ |
| 2 | $\mu_{X+k} \equiv \mathrm{E}[X + k] = \mathrm{E}[X] + k$ |
| 3 | $\mu_{XY} = \mathrm{E}[X]\mathrm{E}[Y]$ |
| 4 | $\sigma_X^2 \equiv \mathrm{Var}[X] = \mathrm{E}[X^2] - \mathrm{E}[X]^2 \implies \mathrm{E}[X^2] \geq \mathrm{E}[X]^2$ |
| 5 | $\sigma_{[X+k]}^2 = \sigma_X^2$ |
| 6 | $\sigma_{XY}^2 = \mathrm{E}[X^2]\mathrm{E}[Y^2] - \mathrm{E}[X]^2\mathrm{E}[Y]^2$ |

cases, but certainly not all. For the case study shown in Fig. 4.5 they are appropriate: I verified that nuclear size and staining intensity were uncorrelated with each other and with position, and that the two foreground signals (Hoechst and Smad2/3) are independent (data not shown).

Finally, I noted earlier that units of fluorescence are typically arbitrary and that $S$ causes a multiplicative change in relative intensity across an image. This means that we are free to choose how to define the expected value of $S$: a value other than 1 would cause a scaling of the foreground and background values, but this scaling would retain the relative relationship between all measured intensities and so would be of no consequence. For convenience, then, I define shading so that its expectation value across all images and pixels is $\mathrm{E}[S] = 1$ so that, for a large number of cells, $\mathrm{E}[S_c] \approx 1$. This definition simplifies the mathematical derivations below, as the $\mathrm{E}[S_c]$ term can be dropped from several formulae.

**Effects of background $B$ on the average intensity feature $A$**

For this case, we ignore the effects of $S$ and focus on $B$. We therefore set $B > 0$ and $S_c = 1$, so that the average feature from Equation 4.4 simplifies to $A_c = F_c + B$. This results in the distribution properties for this feature in Equations 4.7-4.9.

$$\mu_A \equiv \mathrm{E}[A_c] = \mathrm{E}[F_c + B] = \mathrm{E}[F_c] + B \tag{4.7}$$

$$\sigma_A \equiv \sqrt{\mathrm{Var}[A_c]} = \sqrt{\sigma_{F_c+B}^2} = \sigma_{F_c} \tag{4.8}$$

$$cv_A \equiv \frac{\sigma_A}{\mu_A} = \frac{\sigma_{F_C}}{\mathrm{E}[F_c] + B} \tag{4.9}$$

It is clear that $B$ will always cause an increase in the mean of average intensities, $\mu_A$ (Equation 4.7). The standard deviation, $\sigma_A$, is unaffected by background (Equation 4.7, using Property 5 from Table 4.2). As a consequence of the constant $\sigma_A$ and increased $\mu_A$, the coefficient of variation will decrease with increasing background. In summary, background will cause an overestimation of $\mu_A$, an underestimation of $cv_A$, and will not affect $\sigma_A$ (ee the case study in Fig. 4.5, top left panel).

**Effects of background $B$ on the total intensity feature $T$**

The situation is the same as the previous case, except with the inclusion of cell size. The total intensity feature Equation 4.5 therefore simplifies to $T_c = \alpha_c(F_c + B)$, resulting in the distribution properties shown in Equations 4.10-4.12.

$$\mu_T \equiv \mathrm{E}[T_c] = \mathrm{E}[\alpha_c(F_c + B)] = \mathrm{E}[\alpha_c](\mathrm{E}[F_c] + B) = \mathrm{E}[\alpha_c]\mathrm{E}[F_c] + \mathrm{E}[\alpha_c]B \tag{4.10}$$

$$\sigma_T \equiv \sqrt{\mathrm{Var}[T_c]} = \sqrt{\sigma^2_{\alpha_c(F_c+B)}} = \sqrt{\sigma^2_{\alpha_c F_c} + 2\sigma^2_{\alpha_c}\mathrm{E}[F_c]B + \sigma^2_{\alpha_c}B^2} \tag{4.11}$$

$$cv_T \equiv \frac{\sigma_T}{\mu_T} = \frac{\sqrt{\sigma^2_{\alpha_c F_c} + 2\sigma^2\mathrm{E}[F_c]B + \sigma^2_{\alpha_c}B^2}}{\mathrm{E}[\alpha_c]\mathrm{E}[F_c] + \mathrm{E}[\alpha_c]B} \tag{4.12}$$

Though somewhat less obvious than for the average feature, it should be clear that increasing $B$ will cause an increase in the mean of total intensities, $\mu_T$ (Equation 4.10, using Property 3 from Table 4.2). Importantly, each cell will be affected differently by background, depending on its size. Unlike the average feature, the standard deviation $\sigma_T$ increases with background (Equation 4.11, using Properties 3 & 6 from Table 4.2).

Because both $\sigma_T$ and increased $\mu_T$ increase with increasing $B$, it is not immediately obvious what the effect should be on the coefficient of variation, $cv_T$ (as the numerator and denominator in Equation 4.12 both are proportional to $B$). However, if we take the derivative of the $cv$ with respect to a changing background, we get Equation 4.13. Because this derivative is always $\leq 0$, the $cv_T$ will decrease with increasing background. In summary, background will cause cell size-dependent overestimation of $\mu_T$ and $\sigma_T$, and underestimation of $cv_T$ (see the case study in Fig. 4.5, top middle panel).

$$\frac{d}{dB}(cv_T^2) = -2\frac{\sigma^2_{F_c}\mathrm{E}[\alpha_c^2]}{\mathrm{E}[\alpha_c]^2(\mathrm{E}[F_c] + B)^3} \leq 0 \tag{4.13}$$

**Effects of shading $S$ on the average intensity feature $A$**

We now move on to the effects of shading on the average and total features, and therefore set $B = 0$. Note that increasing the magnitude of $\mathrm{E}[S_c]$ will have no effect on any of these features, as it is the same as a change in units. Because shading is a variation in intensity across an image, we can therefore modulate its strength by changing the variance of this image component. The larger the variance, the more shading. For this case, then, we set $\mathrm{Var}[S_c] > 0$. The average intensity feature from Equation 4.4 therefore simplifies to $A_c = S_c F_c$. This results in the distribution properties shown in Equations 4.14-4.16.

$$\mu_A \equiv \mathrm{E}[A_c] = \mathrm{E}[S_c F_c] = \mathrm{E}[S_c]\mathrm{E}[F_c] = \mathrm{E}[F_c] \tag{4.14}$$

$$\sigma_A = \sqrt{\sigma^2_{S_c F_c}} = \sqrt{\mathrm{E}[S_c^2]\mathrm{E}[F_c^2] - \mathrm{E}[S_c]^2\mathrm{E}[F_c]^2} = \sqrt{\mathrm{E}[S_c^2]\mathrm{E}[F_c^2] - \mathrm{E}[F_c]^2} \tag{4.15}$$

$$cv_A \equiv \frac{\sigma_A}{\mu_A} = \frac{\sqrt{\mathrm{E}[S_c^2]\mathrm{E}[F_c^2] - \mathrm{E}[F_c]^2}}{\mathrm{E}[F_c]} \tag{4.16}$$

From Equation 4.14 it is obvious that $S$ has no effect on the mean of average intensities, $\mu_A$, as it falls out of the formula entirely (using Properties 1 & 3 from Table 4.2). The standard deviation, $\sigma_A$, however will increase with background (Equation 4.15, using Properties 1, 4, & 6 from Table 4.2). As a consequence of the increasing $\sigma_A$ and constant $\mu_A$, the coefficient of variation increases with increasing background. In summary, shading will cause overestimation of variation for $A$ (see the case study in Fig. 4.5, bottom left panel). Shading does not affect $\mu_A$, but this is not surprising since I defined shading to have a mean of 1 specifically so that it would not affect the mean.

**Effects of shading $S$ on the total intensity feature $T$**

The situation is the same as the previous case, except with the inclusion of cell size. The total intensity feature Equation 4.5 therefore simplifies to $T_c = \alpha_c S_c F_c$. This results in the total intensity distribution properties shown in Equations 4.17-4.19.

$$\mu_T \equiv \mathrm{E}[T_c] = \mathrm{E}[\alpha_c S_c F_c] = \mathrm{E}[S_c]\mathrm{E}[\alpha_c F_c] = \mathrm{E}[\alpha_c]\mathrm{E}[F_c] \tag{4.17}$$

$$\sigma_T = \sqrt{\sigma_{S_c F_c}^2} = \sqrt{\mathrm{E}[S_c^2]\mathrm{E}[\alpha_c^2 F_c^2] - \mathrm{E}[S_c]^2\mathrm{E}[\alpha_c F_c]^2} = \sqrt{\mathrm{E}[S_c^2]\mathrm{E}[\alpha_c^2 F_c^2] - \mathrm{E}[\alpha_c F_c]^2} \tag{4.18}$$

$$cv_T \equiv \frac{\sigma_A}{\mu_A} = \frac{\sqrt{\mathrm{E}[S_c^2]\mathrm{E}[\alpha_c^2 F_c^2] - \mathrm{E}[\alpha_c F_c]^2}}{\mathrm{E}[\alpha_c F_c]} \tag{4.19}$$

The derivation is nearly the same as that for the previous case for the average feature. As before, from Equation 4.17 we see that $S$ has no effect on the mean of total intensities, $\mu_T$. Also as before, the standard deviation, $\sigma_T$, will increase with shading, as will $cv_T$ (see the case study in Fig. 4.5, bottom middle panel). Thus, increasing the shading always artificially increases the variation of the average and total intensity features.

**Effects of $B$ and $S$ on $R$**

I now turn to the ratiometric feature (Equation 4.6), which turns out to be the least generalizable even within the narrow constraints defined at the start of this section. And so to simplify, I further constrain the discussion to the ratio of average intensities. This allows us to set $\alpha = 1$ for both signals (note that these terms also cancel when taking the ratio within a single compartment). Thus, we have Equation 4.20.

$$R_c = \frac{S_{c1}(F_{c1} + B_1)}{S_{c2}(F_{c2} + B_2)} \tag{4.20}$$

There are a few distinct cases we can examine to get an idea of how this ratio behaves. In the first case, we can have the signals originating from the same channel (for example, the ratio of nuclear to cytosolic Smad2/3). Here, $B_1 = B_2$ and, under my earlier assumption that cells are small relative to the rate of change of $S$, $S_{c1} = S_{c2}$. The ratio thus becomes $R_c = \frac{F_{c1}+B}{F_{c2}+B}$ and is unaffected by shading. The effects of $B$ are less obvious because it is in both the numerator and denominator. In the limiting case, however, $\lim_{B \to \infty} R_c = 1$. Since all values are forced to 1 with large $B$, the variation across the population will be forced to 0. In less extreme cases, however, background can

**Figure 4.6:** Ratiometric features are idiosyncratically affected by background. The same initial data as in the right panels of Fig. 4.5, but where each single-nucleus ratio after addition of background is plotted against its true value (1000/3700 cells shown). To each average was added 0 (black), 1000 (red), or $10^{10}$ (blue) background units before taking the ratio. The first two artificial cases are the same as those in Fig. 4.5. The impact on the apparent ratio is dependent on the relative sizes of the foregrounds and the background, such that the data becomes scrambled in the red curve. Note that the distribution shape is also affected, such that adding background decreased the skewness. $F_1$, average single-nuclear Hoechst; $F_2$, average nuclear Smad2/3.

either cause an increase or decrease in the apparent variation, depending on its size relative to both foreground terms. Because of this, the error in the measured ratio can vary from cell to cell within the same population. For example, if two cells have the same ratio but different absolute foreground values, $B$ will cause their ratios to diverge (e.g. $\frac{1+B}{2+B} \neq \frac{10+B}{20+B}$). This effect causes the scrambling seen in Fig. 4.6.

In the second case, we could allow $F_{c1}$ and $F_{c2}$ to come from distinct channels, which would also allow them to have different shading and background. As a consequence, the effects on the distribution properties are extremely difficult to predict due to the presence of many independent variables (as indicated in the case study in Fig. 4.5, right panels). I therefore leave the discussion here, with the conclusion that the mean and variation of the ratio feature can both increase or decrease with different ranges of component values. This unpredictability has important ramifications for applications such as fluorescence resonance energy transfer (FRET) where interpretation relies on accurate cross-channel ratios [236].

**On the importance of image correction**

Unfortunately, the mathematical discourse above leaves us with the dissatisfying result that the importance of image correction is highly dataset dependent. In Fig. 4.5 I show the size of the distortion for an example dataset with experimentally-reasonable amounts of background and shading,

which shows effects that are measurable and sometimes relatively large. However, different analytical requirements can tolerate quite different amounts of feature distribution distortion before data interpretation is affected. The best I can do then, besides the easy blanket statement "always correct your images!" is to provide a few rules of thumb for deciding on the importance of image correction.

First, I completely ignored the detector contribution ($D$) in the above discussion because the matrix $D$ is unchanging and is therefore trivial to subtract from images. Removal of $D$ should be a part of all image analysis pipelines. This step is rarely explicitly performed in the literature but can have a large impact on measurements of weakly-fluorescent samples. In particular, shading will be underestimated without subtraction of the detector fluorescence contribution.

Second, shading tends to increase feature distribution widths and so should be corrected whenever measurement of true biological variability is important. Additionally, if the foreground to background ratio is low, then shading may cause foreground values in one part of an image to fall below background values in another part. This can have important consequences to image segmentation (Section 4.4.1) [232]. However, in the case that biological variability is high relative to the variation in shading it would not be necessary to correct for $S$. This is because the observed average and total distributions are the convolution of the shading and foreground distributions. If two normal distributions with variances $\sigma_1^2$ and $\sigma_2^2$ are convolved, they yield a new distribution with $\sigma_3^2 = \sigma_1^2 + \sigma_2^2$. Because shading is a relative term, we have to scale it to the size of $F_c$ to make use of this property: $\sigma_{\text{observed}}^2 = (F_c \sigma_{S_c})^2 + \sigma_{F_c}^2$. Thus, if $F_c \sigma_{S_c}^2 << \sigma_{F_c}^2$ the observed distribution will be close to the true distribution. While this is a useful rule of thumb, I note that I have never observed normally-distributed shading values (see examples in Fig. 4.4b).

Third, background can dramatically shift feature distribution means. Background should therefore be corrected either when accurate means or accurate ratios between two means are needed. When background is low compared to even the lowest foreground values, however, it will have little impact on the feature distributions discussed above.

Finally, when using cross-channel ratios extra care should be taken to ensure that both background and shading are corrected. These ratios should always be interpreted with an eye towards the possible effects of imaging artifacts, as artifacts can distort rations in unpredictable ways.

### 4.3.2   Review of image correction methods

Now that I have given some motivation for the importance of image correction, I turn to the available methods for this process. Here, I briefly review the correction methods that are commonly employed in the literature. In general, when deciding on a method the investigator should test its theoretical performance given the image model described in this chapter, $I = D + S(F + B)$. By taking this approach, deficiencies or important assumptions of the methods should become clear.

There are many published methods for fluorescence image correction, perhaps as many methods as there are labs, due to the idiosyncrasies of imaging data and the lack of standardized approaches. I group the most common methods into two broad, non-exhaustive categories, which I refer to as

"reference-image" and "per-image" correction. Reference-image methods obtain correction parameters from one image and then use those parameters to correct another image. Per-image methods find such parameters in the very image that is to be corrected.

Reference-image methods are straightforward and so are commonly used and recommended throughout the literature [226, 233, 237–239]. These methods make a key assumption: that the reference image has approximately the same shading and background as do the images to be corrected. A good example of this approach uses two reference images to flatten shading, subtract background, and normalize the fluorescence intensity to a standard [226]. One reference image, $I_{\text{uniform}}$ contains a dissolved fluorophore; because this image would be flat without shading, it can be used to determine how much shading is present in the real images. The other reference, $I_{\text{background}}$ is the same as the sample images but contains no sample. This image can thus be used to estimate background. Equations 4.21-4.24 demonstrate this method.

$$I_{\text{sample}} = D + S(F_{\text{sample}} + B) \tag{4.21}$$

$$I_{\text{background}} = D + SB \tag{4.22}$$

$$I_{\text{uniform}} = D + S(B_{\text{uniform}} + B) \tag{4.23}$$

$$I_{\text{corrected}} = \frac{I_{\text{sample}} - I_{\text{background}}}{I_{\text{uniform}} - I_{\text{background}}} = \frac{[D + S(F_{\text{sample}} + B)] - [D + SB]}{[D + S(B_{\text{uniform}} + B)] - [D + SB]} = \frac{F_{\text{sample}}}{B_{\text{uniform}}} \tag{4.24}$$

While reference-image methods are straightforward and often easy to perform, some of those recommended in the literature are only partially corrective. This is especially true with respect to the detector value $D$, which I have not seen explicitly accounted for in these methods. To find out if a given method performed complete correction, investigators can plug the image model into the method and check that, algebraically, the output is either $F$ or some normalized form of it (as in Equation 4.24). However, note that partial correction may be sufficient in some cases, particularly when background intensity is low compared to foreground intensity and when within-image shading variation is low compared to foreground variation.

Per-image methods, on the other hand, have the challenging task of measuring all image components ($D$, $S$, $F$, and $B$) within the image that is to be corrected. These methods typically work by trying to fit a model to the combination of non-foreground layers, $D + SB$. Therefore the primary difficulty is that images typically contain varying fractions of foreground pixels (e.g. due to variation in cell density), so that determination of which pixels consist of background is non-trivial. Further, even if identification of background pixels is straightforward within an image, the size of the background signal "underneath" a foreground object is necessarily unknown and must be predicted. The method for this prediction is what separates the different per-image correction algorithms.

Some per-image methods fit a polynomial [225] or spline surface [240–242] to predicted non-foreground pixels, therefore assuming a particular structure to the shading patterns. These methods may also assume properties of the foreground objects (e.g. fluorescing cells), such as a maximum size in the case of the rolling ball algorithm employed by ImageJ [243] and Fiji [244], and may

perform non-linear transformations of the underlying images. Additionally, the accuracy of these approaches necessarily decreases with increasing cell density as there are fewer background pixels from which to estimate $I_{\text{background}}$. High confluency or cell clumping can thus cause per-image methods to introduce artifacts.

When they work properly, per-image based methods generate $I_{\text{background}}$ (Equation 4.22) from each sample image, $I_{\text{sample}}$ (Equation 4.21). From this point, then, the reference-based and per-image based correction methods are the same; the only difference is in how $I_{\text{background}}$ is obtained. The question that then remains in both cases is which mathematical operation to perform in order to correct the sample images. In the literature, subtraction ($I_{\text{sample}} - I_{\text{background}}$) is frequently used. However, subtraction does not remove shading, as the algebraic result is $SF$ instead of $F$. The standard rolling ball algorithm employed by ImageJ and FIJI uses this subtractive approach.

The better approach is to use division to remove shading. This can be done in combination with subtraction to remove both background and shading, as in the example above (Equation 4.24). However, many studies use simple division ($I_{\text{sample}}/I_{\text{background}}$), which results in $\frac{D+S(B+F)}{D+SB} = 1 + \frac{SF}{D+SB}$. This result is a background-normalized foreground with incompletely-corrected shading. The correction can be completed by prior subtraction of $D$ from both images, followed by subtraction of 1, which would yield the normalized image $F/B$.

I use a different approach from those listed above, which is to independently estimate all non-foreground parameters so that I can perform the image algebra that expicitly returns $F$. I discuss this method next.

### 4.3.3   An improved correction method for high-throughput imaging

Having determined the importance of image correction, and observed the diversity and frequent inaccuracy of the available correction methods, I sought to develop a simple and robust method for my own imaging applications. All of the work in this dissertation took advantage of the throughput of microtiter plates (e.g. 96- and 384-well plates), and so in particular I needed to be able to accurately correct large numbers of images with minimal computational overhead. Further, I needed the correction to be robust to cell density, so that I could correctly measure single-cell biological variability under a wide variety of experimental conditions.

The approach that I developed (published in [235]), is a reference-based method that relies on a key observation: shading patterns are highly predictable in microtiter plates. This allows for a correction method that forgoes the need to estimate parameters for every image, yet is more accurate than using only a single set of correction parameters. Below I show the evidence that shading is indeed predictable, and then outline a correction method for taking advantage of this fact.

**The shading pattern is a function of within-well position**

I was initially surprised at the high variability in shading patterns that I observed within imaging datasets from microtiter plates. In the literature there appears to be an implicit assumption that such variability is common and unpredictable, as the most common correction methods used in big

**Figure 4.7:** For a within-well position, shading patterns are consistent throughout the entirety of a microtiter plate. I imaged 3x3 grids within every well of a 384-well plate containing dissolved fluorescein. Optical setup: black plastic 384-well plate (Corning #3712), FITC filter, 20X objective, Andor sCMOS camera. I cropped all images identically to remove pixels that contained well edges. Within-well images were montaged into a single image, and then normalized to the median pixel value of that montaged image (e.g. in **a**, top). This converts the images to an estimate of shading as defined in this chapter. Reference 3x3 image montages were made by taking the per-pixel average across all 3x3 montages (from all wells). This was done either with all 9 images in the 3x3 grid (Position-dependent correction) or with only the central image from the grid repeated 9 times (Position-independent correction). Finally, the resulting grids were montaged to show the image properties across the entire 384-well plate (**b**). Panel **a** shows larger thumbnails of the 3x3 image grids in well A1 before and after position-dependent correction. The histograms in **b** shows the distributions of all relative pixel intensities in the montaged images. Tighter distributions indicate more accurate correction. White objects in the montaged images are auto-fluorescing debris; small black corners in remaining in corrected images are due to inclusion of a small portion of the black well edge.

datasets are per-image. However, shading is an artifact that is generated by the light path, which is a static aspect of the microscope optical setup, and so I would have expected it to be unchanging within a dataset. Indeed, images taken at different positions along a glass slide seem to show an unchanging shading pattern. I therefore reasoned that it was the microtiter plates themselves that modified the shading pattern. Further, since microtiter plates are essentially arrays of identical wells, I predicted that the shading modification must be a function of the image location within a single well (as opposed to the position within the plate). The source of the within-well shading modulation could be due to reflections from well edges, local distortions of the imaging surface, or lensing by the solvent meniscus.

This prediction of within-well positional dependence of shading patterns bore out, as can be seen in Fig. 4.7. Further, this turned out to be consistent for both 384- and 96-well plates from various manufacturers and with different physical specifications and materials. In 96-well plates this positional effect is less dramatic, which is consistent with shading patterns being modified by

**Figure 4.8:** Shading ($S$) can be estimated using sample-containing images. Media-only wells are controls (0% confluency). The spatially uniform fluorescence in these control wells was used to estimate the "reference" shading by per-pixel averaging across 42 media-only control wells. I define "confluency" as the fraction of pixels in an image with intensities $> 3\sigma$ above background. I used A549 cells expressing two differently-colored fluorescent proteins (left, nuclear CFP; right, cytosolic mCherry) at three different seeding densities (each seeding density had 84 replicates), effectively yielding six different confluency levels. I additionally created mixed-confluency image sets by randomly selecting across all confluency levels within each color channel. For each fixed number of images, $n$, selected from the same within-well position, I computed the "estimated shading" as the per-pixel median of $n$ randomly selected images. The error in the shading estimates are computed as the $cv$ of the per-pixel ratio of (estimated shading)/(reference shading). For each confluency level, the inaccuracy of the sample-based shading estimate generally decrease with increasing numbers of images.

well edges (as these edges are much further apart than are those in 384-well plates). I also note that a position-independent correction method (that uses a single reference image instead of one per within-well position, Fig. 4.7, blue) dramatically improves the images as well. Therefore this even simpler method may be suitable for some datasets (though the resulting multi-modality in Fig. 4.7 could, in principle, generate artificial cellular subpopulations (Section 1.4.2)).

**Pipeline for within-well position-based image correction**

The within-well positional shading constancy allowed me to implement a simple and robust image correction pipeline for images from microtiter plates. This pipeline is based on the existing reference-based approaches discussed above and in other sources [234, 245].

1. **Calibrate the microscope stage for the microtiter plate.** Ideally, images should be acquired near well centers, and relative within-well image positions should have minimal drift between wells (e.g. images in well A1 should not be closer to the left well edge than those in H12). It is important to be aware that this method will become inaccurate with increasing positional drift. Also note that microscope stage-driving software can vary in its accuracy, and so custom solutions might be required for plates with small wells.

2. **Measure the dark current component, $D$.** The simplest approach is to capture images

without a light source, or with the light path diverted from the camera, and per-pixel average the images. In practice, a small number of images is sufficient ($n$=6-20 in my analyses).

3. **Subtract dark current, $D$, from all images in the dataset.** The resulting images $(I-D)$ are then modeled by $S(B+F)$. This important step should be performed regardless of the correction method subsequently used, for the reasons explained above. Without subtracting $D$, the estimated shading patterns can become increasingly inaccurate with large foreground values or small background values.

4. **Estimate shading $S$ for each distinct within-well position.** This step is the most complicated in the pipeline and should be performed with care. The goal is to obtain a reference shading pattern for each within-well position. For example, if an investigator has 9 images/well she will need to estimate 9 shading patterns. There are at least two possible approaches to this estimate.

   - **Uniform reference images.** This approach is the most robust, but is only possible if there are extra wells that can be reserved solely for acquiring reference images. In these extra wells, add dissolved fluorophore at the appropriate concentration for the intended exposure time. These wells should then be imaged along with the other sample wells. In practice, I have found that a small number (e.g. 6) of such uniform images is often sufficient so long as the wells are free from fluorescing debris. Across all wells $w$, for a given within-well position $p$, the images $I_{w,p}$ should be per-pixel averaged to obtain a reference image $R_p = \mathrm{mean}_w(I_{w,p}-D)$. Note that the per-pixel median or other quantile may be more robust.

   - **Sample-based reference images.** In high-throughput studies extra wells may be unavailable for acquiring reference images. In this case, the investigator can take advantage of the large number of available sample images and otherwise use the same method as for uniform reference images. Because these sample images contain foreground, many more values per coordinate are needed to get an accurate estimate of shading. In Fig. 4.8 I show how accuracy is dramatically affected by the number of images used. The same figure suggests that 20-40 cell-containing images may often be sufficient, but this is dependent on the cellular confluency of those images.

   The results of this step will be one reference image per within-well position. As defined earlier, the shading values should be centered on 1 to maintain the original intensity range. Therefore once all reference images are collected each image $R_p$ should be divided by the median or mean of all reference image pixel values (with coordinates $(r,c)$. The resulting shading patterns are described by $S_p = \frac{R_p}{\mathrm{mean}_{r,c,p}R_{r,c,p}}$

5. **Correct for the shading $S$ at each within-well position.** To correct for shading, every image $I_p$ should be per-pixel divided by the corresponding shading pattern $S_p$ obtained in the previous step: $\frac{I_p}{S_p}$. The resulting images then contain only $B+F$.

6. **Subtract the background from each image.** There are multiple approaches to this task. The proper choice depends on the particulars of the dataset, and so I illustrate two cases here. In both cases, after subtracting background we end up with the approximate foreground signal, $F$.

   - **Global background subtraction.** In some datasets, the background may vary little between images. In this case, a global background value can be estimated by averaging background pixels from a representative image. Subtract this value from all pixels in all images.

   - **Per-image background subtraction.** In other datasets, there may be significant variability in background from well to well. This could be due to errors in staining or variation in exposure time. Background values can be estimated per image by e.g. Otsu thresholding to automatically identify background pixels [246]. Alternatively, a low quantile pixel value can be taken as background (the quantile choice is dependent on cell density). Then, for each image, subtract its estimated background value.

     I note that, in the case of background differences being due to something that would imply foreground differences, such as due to exposure time or staining variation, these per-image background estimates can be used to normalize intensities between images. For example, to normalize all images to some reference $B_{ref}$, determine the normalization factor by taking its ratio with the per-image background $B_{ref}/B_{sample}$. The image can then be multiplied by this factor at every pixel, bringing it to the same intensity level as the reference image. Care should be taken with this approach, however, as it is not generally obvious when background variation is predictive of foreground variation.

Note that this pipeline is only meant to remove shading $S$, background $B$, and detector $D$ from images. What is left will be all foreground layers, which may include various artifacts. Further, the foreground values may vary for reasons independent from imaging. For example, it is well established that assays in microtiter plates can show batch, edge, row, and column effects. These effects may non-biologically change the true values of $F$ and should thus be normalized after image correction [247–251].

### 4.3.4   Image correction quality control

As with any data manipulation, the image correction method outlined above should be checked to ensure that no errors are introduced. An obvious approach is to simply visually inspect a subset of images, as in Fig. 4.9a, though automated approaches are also feasible.

For background correction, images should be checked to ensure that only the background has been subtracted. This can be done by examining background pixels. Since a perfect subtraction will have set the centroid of the background pixel distribution to zero, roughly half of all post-correction background pixels should be zero while the other background pixels take on small values. Unfortunately, this task can be difficult to automate for the same reason that errors may be introduced:

**Figure 4.9:** Images can be over- or under-corrected, and therefore require quality control. **a**, Visual inspection should reveal a flat background and a foreground that does not vary in a spatially predictable manner after correction. Histone H2B-Cyan Fluorescent Protein labeled A549 cells. Image courtesy Jungseog Kang (Altschuler & Wu lab, UT Southwestern). **b**, After segmentation and feature-extraction, single-cell features can be tested for within-image spatial dependence. Top left, a sample image of Smad2/3-stained human colonic epithelial cells, showing the radial metric of "distance from image center." Plots show nuclear area and total nuclear Smad2/3 or Hoechst as a function of within-image radial position ($n$=1000/~20000 randomly-chosen cells). Intensities were median-normalized by well to prevent true experimental variation from affecting this test of positional phenotype dependence. Inset numeric value is the Pearson correlation coefficient for the two variables plotted.

cell density may vary dramatically within a dataset, complicating the automated identification of background pixels.

For shading correction, the likely artifacts will be spatial. In other words, if there is over-correction or under-correction, this will cause certain regions of every image (and the cells within those regions) to have systematically higher or lower intensities than the global mean. This can be checked in an automated way after cell segmentation by testing for dependence of single-cell features on within-image position (see an example of this approach in Fig. 4.9b).

## 4.4 On segmentation and single-cell features

Once a dataset has been collected, and corrected as in the previous section, the analysis begins. In general, we want to obtain interesting properties of the foreground layers in our images. In specific, it is nearly always the foreground properties within individual cells that we care about. "Segmentation" is the process of identifying those cells (or intracellular objects) and extracting them from the rest of the image. Once cells have been segmented, properties of their pixel values and spatial arrangement can be measured and stored as sets of features.

Cellular segmentation is still an unsolved computer vision problem [222], in large part because

experimental needs are too specific to allow for a good general solution. There is, however, an array of partial solutions available. These range from simple to highly complex, and vary tremendously in their accuracy and utility. In this section I briefly discuss segmentation algorithms and how one can go about choosing single-cell features that are biologically meaningful. My goal in part is to provide the rationale for my own choices for the experimental work in Chapter 3.

### 4.4.1 Segmentation approaches

There are many approaches to cellular segmentation. Software solutions such as CellProfiler [252] implement many of the most-used segmentation algorithms so that the user can choose one appropriate to the experiment. Unfortunately, which algorithm should be chosen is not a trivial matter. Many of these algorithms are complex and therefore difficult to understand and use. Further, the simpler algorithms may fail to provide sufficiently accurate segmentation for some image properties. As a consequence the path of least resistance is often manual segmentation using tools such as ImageJ [243], but this labor-intensive approach limits the resulting sample size and may yield biased investigator-specific outcomes.

The value of automated segmentation approaches should be obvious: they allow for the rapid and reproducible identification of huge numbers of cells. Automated approaches are also biased, due to the choice of parameters for the algorithm, but the bias is systematic and does not change between images. There are plenty of difficulties with automation, however. With huge datasets comes the inability to perform thorough quality control. Additonally, there may be no single set of segmentation parameters, or even a single algorithm, that will successfully segment all cellular phenotypes in a diverse experimental setup (e.g. a drug screen). Automation thus requires extensive testing and a wary mindset.

As a consequence of these issues, I strongly advocate for use of the simplest segmentation method that is capable of answer a given experimental question. Because so many aspects of images can be measured, it is easy to get carried away with trying to obtain every single piece of data that the images contain. As already noted, however, the number of potential measurements is large; obtaining all data from images is not only impossible, it probably is not wise since each extracted feature should be understood at a biological level before it is used (and, as I discuss below, biological interpretation of features is not a trivial task). The use of simpler methods makes the resulting segmentation more understandable, and so conditions that will cause the algorithm to produce garbage are more predictable and intuitive. A few of the more straight-forward and commonly-used methods are threshold, watershed, and voronoi segmentation, discussed next (these are depicted graphically in Fig. 4.10).

Threshold segmentation (the approach that I use in this dissertation) is probably the simplest method after manual segmentation. It works by assuming that the pixel intensities within cells are generally higher then those in the image background Fig. 4.10a. Therefore a threshold can be chosen either manually or via some mathematical or algorithmic approach (e.g. Otsu thresholding [246]) that best separates background and foreground pixels. Neighboring pixels classified as foreground

**Figure 4.10:** Cartoons of three basic segmentation methods. **a**, Threshold segmentation separates background $B$ from foreground $F$ by fluorescence intensity. Neighboring pixels are then considered part of the same object. The histogram represents the distribution of all pixel intensities within an image. **b**, Watershed starts with a known intracellular point (e.g. the nucleus, blue) and moves outwards (gray, arrows) until it reaches cell boundaries. **c** Voronoi segmentation assigns to a set of already-known objects all of the space closer to that object than any other. Red lines, boundaries of the Voronoi cells, using nuclei centroids as the known points. Gray area, segmented region for the middle cell.

can be grouped into objects, such as nuclei. This method tends to fail for whole-cell segmentation when cellular density is high, as neighboring cells can be segmented as single objects. Additionally, it is not necessarily true that the nuclear or cytosolic compartments contain higher pixel intensities than do the background. This is dependent on the probes used, and is of particular difficulty for live-cell imaging (see Appendix A for an experimental solution). Cellular nuclei are frequently segmented using this method, as they tend to be spatially distinct even with high cell density, and DNA stains such as Hoechst and DAPI create bright foreground signals. Threshold-segmented nuclei often form the basis for more complex algorithms.

There are many specific algorithms for watershed segmentation [46, 253, 254] but the general approach is roughly the same. The algorithm starts with a "seed," which is a pixel coordinate in the image. This seed can be randomly generated or chosen from centroids of threshold-segmented nuclei. The algorithm then searches away from that seed point, following paths of increasing intensity (Fig. 4.10b). Sharp decreases in intensity, as frequently occurs between cells or between a cell edge and the background, will cause the outward movement to stop. These algorithms are thus useful for segmenting irregularly-shaped and slightly crowded cytosolic regions. However, they may require many parameters, and can yield over-segmentation (the breaking up of single cells into many objects).

The final method that I want to make note of is Voronoi segmentation. I have not frequently seen this method used in the context of segmenting cells in tissue culture, but it can be useful in the case that cells are at an extremely high density (so that there is no background) and when those cells are similar in size and relatively round or cuboidal. With this approach, a set of seeds are again needed. These will typically be the centroids of threshold-segmented nuclei. For each centroid, then, the algorithm assigns to it all space closer to that centroid than to any other (Fig. 4.10c). This is a simple and efficient method for roughly segmenting cellular cytosolic compartments.

With this brief overview of few segmentation methods, a few points should be clear. First,

many segmentation approaches are fluorescence-based (though some use brightfield) and so require staining of the cellular compartments to be segmented (see Appendix A for a live-cell solution to this). Second, nuclei are much simpler to segment than cytosolic regions, because nuclei typically have narrower ranges of size, shape, and staining intensity. Therefore threshold segmentation of nuclei is generally considered to be accurate, and is a frequent first-step for more complicated segmentation pipelines. Third, different foreground objects may require different algorithms for accurate segmentation. Taken together, the above points help to explain why cellular segmentation does not have a general solution.

### 4.4.2   Understanding and choosing single-cell features

Given the nearly unlimited set of features to choose from it can be difficult to determine the subset that is most appropriate to the study at hand. There are a few broad approaches to this problem. The obvious, but non-trivial, approach is to first choose the biological property of interest and then identify or create features that approximate that property. A less obvious approach is to obtain a large number of features and use computational methods to choose those that are the most informative [27]. In the latter case, the features need not be biologically interpretable at all. For this discussion I focus on the first case.

Biologically-motivated features are necessarily approximations of the underlying property of interest. It is therefore important to be aware of how the features are defined so that the data can be interpreted properly. As an example, a project in the Altschuler & Wu lab required measurement of neutrophil polarity, but needed that measurement to be performed on fluorescently-labeled cytoskeletal components. There is no obvious mathematical feature of fluorescently-labeled actin or microtubules that would indicate the degree of polarization of a neutrophil. An approximate solution was therefore developed, which measures how close together in space are the brightest pixels [7, 255].

This polarity feature will decrease in value with, for example, increased collection of actin to one side of the cell. In the case that the actin intensity is diffuse throughout the cell, the brightest pixels will also be diffuse and so the feature value increases. This feature therefore serves as a useful proxy for polarity. However, the feature is sensitive to bright punctate artifacts that are common to immunostained images and is distorted by differences in cell size. By being aware of these pitfalls they can be addressed. In this case, visual quality control was performed on every cell to ensure the absence of artifacts, and the feature calculation was modified to compensate for distortions due to cell size [7, 255].

Investigators will more frequently use combinations of simpler features, such as those describing aspects of cell size and shape ("morphological features") and of fluorescence intensity ("intensity features"). It is important to be aware that these two classes are not completely independennt. Intensity features in particular can be highly dependent on morphological features.

As a simple example, the average intensity is dependent on the area. It is therefore important to interpret changes to the average with care, as the change could result from a change in cell size,

a change in concentration of the labeled protein, or a combination of the two. The ratio of nuclear to cytosolic average intensities additionally suffers from this potential confusion. It has an added difficulty, however, in that a change in the ratio can be due to movement of the labeled protein from one compartment to the other, or to independent changes in one or both compartments.

Ideally, then, a feature should be carefully chosen to have an unambiguous biological interpretation and be independent of other features whose changes are not important to the study. I argue that interpretability is more important than having a feature that more closely approximates the biological property of interest. For example, in my own work (Chapter 3) I am interested in the concentrations of nuclear transcription factors. Concentration is intuitively approximated by the average intensity feature, yet I instead chose the total intensity feature for all of my analysis. I did so because certain properties of the total intensity feature, discussed next, allow for less ambiguity in how to interpret changes in its values.

### 4.4.3   Benefits of the total nuclear intensity feature

The total intensity feature is a proxy for the absolute count of a fluorescently-labeled target molecule. This feature has the advantage of being independent of cell size, such that changes in cell morphology may change the average but not total intensity. Unfortunately, intensity from widefield microscopy images does not just come from the focal plane, but also from the space both above and below that plane. As a consequence, in image represents a messy volumetric cross section through the $z$-axis. This fact adds some complication to the interpretation of the total intensity feature: is it measuring the total number of molecules in the entire cell volume, from a thin cross-section, or something in between?

Fortunately, my focus on transcription factors in Chapter 3 allows me to mitigate this concern by measuring only the nuclear intensities. The nucleus tends to maintain a taller stature in the $z$-axis than does the rest of the cell (resulting in the famous "fried egg" appearance of cells in tissue culture). Also, nuclei tend to be of more similar size and shape than do cytosolic compartments. For these reasons, I can reasonably assume that whatever the thickness of the imaged section, I will be imaging a similar thickness for all nuclei.

Finally, and perhaps most importantly, the nucleus has a built-in "ground truth" for this feature that allows for both quality control and removal of measurement error (Section 4.5). This ground truth comes from the DNA content of cells, as each cell within e.g. the G1 phase of the cell cycle in reality has a near-identical total DNA content. The total intensity feature is a proxy for this content. Importantly, total nuclear intensity is the only feature with such a "ground truth" reference. The small-molecule stain Hoechst provides a robust and DNA-specific signal that I use throughout this dissertation to measure total DNA content. I therefore use the term "DNA" interchangeably to refer to the actual molecule and as a short-hand for "Hoechst-stained DNA."

**Figure 4.11:** The mutual information metric can be interpreted as yielding $\log_2$ of the number of distinct signal-response relationships. **a**, A toy cases with only two signals, $S_1$ and $S_2$, and two responses, $R_1$ and $R_2$, with uniform joint probabilities. Darkened boxes show occurring signal-response relationships and their joint probabilities. Left, two completely distinct signal-response pairs yield $\log_2(2) = 1$ bit of mutual information. Middle, both signals cause response $R_1$, meaning that observation of $R_1$ provides insufficient information to know which signal caused that response. The mutual information thus decreases to $\log_2(81/16) \approx 0.78$ bits (using Equation 4.25). Right, with both signals causing both responses with equal probability, there is no mutual information ($\log_2(1) = 0$ bits). **b**, Mutual information measurements of TGFB1-Smad2/3 responses. Distributions show the wide variability of nuclear Smad2/3 accumulation even in response to saturating (10ng/mL) ligand concentrations. Observe that the control and saturating doses yield distributions that are nearly non-overlapping, and as a consequence yield ∼1 bit of mutual information (as there are 2 distinct signal-response relatinships). It should be clear that the mutual information between all concentrations and outputs can be calculated at once, but there is significant overlap between all but the outermost distributions. Thus, addition of each subsequent intervening signal will in this case yield diminishing improvements to the mutual information between TGFB1 and nuclear Smad2/3. $n{>}3000$ human colonic epithelial cells per condition. I used an implementation of the mutual information algorithm as described in [16].

### 4.4.4 Measuring information content of a feature

One of the major hurdles in experimental cell biology is our lack of initial knowledge about which environmental signals ($S$) and cellular responses ($R$) cells care about (discussed in Chapter 1). As a consequence, it is generally unclear how much information about the environment a cell can accurately process and store, though estimates suggest that the features we believe cells care about contain somewhat unimpressive information content [16].

Measuring the information content of a feature with respect to the signal under study can therefore be useful when trying to choose between a set of potential $(S, R)$ combinations. This can be done using the "mutual information" metric between the signal and response, $M(S; R)$ [16]. Mutual information has advantages over other statistical metrics, such as Z-scores and the like, in that it is completely non-parametric (i.e. does not assume a distribution shape) and uses units that can be directly interpreted as "information content," measured in bits. I use this metric in Chapter 3 to compare the information content of ligand concentrations for TGFB$_{sf}$ and Wnt, and so here I dig into mutual information a bit deeper[1] with the goal of providing an intuition to the reader regarding its interpretation.

The mutual information between a set of signals $S$ (e.g. ligand concentrations) and responses $R$ (e.g. nuclear transcription factor concentrations) is defined by Equation 4.25. In the formula, $P(R, S)$ is the joint probability of each signal-response pair and $P(R)$ and $P(S)$ are the marginal

---

[1]Did you catch the joke?

probabilities. The mutual information value returned by this formula is in units of "bits," and can be interpreted as the $\log_2$ of the number of distinct signal-response relationships (see Fig. 4.11a for a toy example). In essence the quantity describes how accurately we would be able to guess the signal if we were told the response (and vice versa). The goal of the mutual information metric then is to measure the degree of overlap between signal-response distributions, it is not to measure how far apart those distributions are. For example, two completely separated response distributions will always have 1 bit of mutual information even if they are infinitely far apart. Therefore the maximum possible information content is $\log_2$ of the number of distinct signals (Fig. 4.11a, left). The minimum mutual information is zero, which occurs when the distributions completely overlap (Fig. 4.11a, right).

$$M(R;S) = \sum_S \sum_R P(R,S) \log_2 \left( \frac{P(R,S)}{P(R)P(S)} \right) \tag{4.25}$$

Importantly, single-cell variability in nuclear transcription factor concentrations is high, such that even the untreated and saturating ligand doses yield a small overlap in cellular responses (Fig. 4.11b) [16]. In other words, if we were given a randomly drawn cellular Smad2/3 response from a dose-response curve for TGFB1 treatment, we would have high uncertainty as to the precise TGFB1 concentration that caused the drawn response.

In summary, mutual information is a metric that measures the overlap of distributions, and so can be interpreted to indicate the number of distinct signal-response pairs for a given feature. This metric can thus be used to directly measure the relative information content of different features $R$ or signals $S$ (as I do in Section 3.1).

## 4.5   Dealing with variation

Experimental error is always present in our measurements. Additionally, any given feature may show extensive biological variability even within apparently homogeneous populations. The biological variation can be informative, as it gives us insight into the limits to accuracy of cellular processing and can reveal phenotypically distinct subpopulations (see Chapter 3). However, if experimental error is mis-interpreted as biological variability, we lose statistical resolution or may assign unwarranted meaning to non-biological variation. It is then important to be able to separate experimental from biological variation.

Experimental variation in fluorescence imaging comes from many sources. The imaging plane itself may cause variation by intersecting cells at different relative heights. This focal plane effect can cause cell-to-cell differences in the degree of focus and in how much off-plane fluorescence is captured. As discussed in Section 4.3.1, image shading and background can also contribute to artificial variation due to microscopy.

Aside from microscopy artifacts, the process of preparing cells for imaging may also generate distortions of true biological variability. For example, variation in cellular surface area or volume may lead to differences in how well an antibody or non-permeable dye can access an intracellular

**Figure 4.12:** An asynchronous cell population can be fit to a simple model of the cell cycle with reasonable accuracy. Left, the theoretical asynchronous cell-cycle distribution consists of delta functions for G1 and G2 cells (i.e. all cells in these populations have identical DNA content) and a uniform distribution for S-phase cells that are moving from the G1 to G2 states at a constant rate. Middle, the total DNA feature of G1 and G2 nuclei shows log-normal variation that can thus be fit to normal distributions after log-transformation. Right, the cell subpopulation models add up to a reasonably accurate estimate of the cell cycle distribution. Gray, filled histograms are of $\log_2$(total DNA), with DNA in arbitrary units, of $\sim 2 \times 10^4$ Hoechst-stained human colonic epithelial cells. Dashed lines are the actual fits to this data using the method described in the text.

target. Such differences in cell morphology may be enhanced or dampened by fixatives, which can cause cells to shrink in the $z$-axis [256].

Finally, non-biological variation can be introduced during segmentation. This can be due to outright errors (e.g. a cell being split into two objects) or to the more subtle fact that the accuracy of a set of segmentation parameters will vary from cell to cell. For example, a chosen threshold that perfectly separates background from foreground for one cell may end up discarding the outer edges of a dimmer cell.

## 4.5.1   Using DNA features for quality control

Due to the error sources discussed above (among others) there may be many outlier cells to discard. Manual or pseudo-manual approaches are often used for this task. Visual inspection of a random subset of segmented cells is a common approach, though automated solutions are needed for large datasets. An example is the identification of out-of-focus images using image-level features from tools like PhenoRipper [257]. I use an automated statistical approach that takes advantage of "ground truth" aspects of total DNA content in cells. This approach uses population-level statistics of DNA features to determine which cells are likely to be properly segmented and in focus. In effect, I make the assumption that cells with biologically-unlikely DNA feature values will have non-biological values in other features as well.

The first step in my quality control pipeline is to identify cell cycle subpopulations by fitting a cell cycle model to the total DNA histograms. This allows for later isolation of these subpopulations and removal of outliers. Note that in tissue culture microscopy mitotic (M-phase) cells are often lost during sample washes and so are already excluded from analysis.

An asynchronous cell population can be accurately fit to a simple model of the cell cycle. The theoretical, variation-free model of this cell cycle consists of of delta functions for the G1 and G2

**Figure 4.13:** DNA features can be used for single-cell quality control. Bottom left, the Total DNA feature is fit to a cell cycle plot (fit not shown) and the G1 population is gated as all cells within $\mu_{G1} \pm 2\sigma_{G1}$ (blue). The population is also statistically gated using the nuclear area (bottom, middle) and the intra-nuclear intensity $cv$ (bottom, right). The gating is $\pm 3 \times$ MAD from the median of each feature, where MAD is the median absolute deviation (median($|X - \text{median}(X)|$). ($3 \times$ MAD $\approx 2\sigma$ for normal distributions.) The gated points are considered to be in-focus G1 cells that are likely segmented properly. In the top plot, these quality-controlled cells are found as red points within the blue box. Data from >4000 Hoechst-stained human colonic epithelial cells.

peaks with a uniform S-phase distribution in between (Fig. 4.12, left). In other words, all cells within G1 or G2 have the exact same DNA content, and cells moving through S-phase increase their DNA content at a constant rate. In reality, the cell cycle distribution arising from measured total DNA consists of log-normally distributed G1 and G2 peaks, with a variable-shaped S-phase distribution in between. I note that, on a log-scale, the G1 and G2 distributions have near-identical standard deviations ($\sigma_{G1} \approx \sigma_{G2}$).

I implemented a simple variant of the Dean-Jett-Fox cell cycle model [258,259] that is sufficient to accurately identify G1 and G2 cells from microscopy data. The formulae that comprise this model are in Equations 4.26-4.28, where: $T_c$ is the log$_2$-transformed total DNA of a single cell; $\mu_{G1}$ is the average of this value across all G1-phase cells; $\sigma$ is the standard deviation of the G1 and G2 distributions; $v$ is the height of the S-phase uniform distribution, $f(T_c)$ is the fraction of cells with the same $T_c$ DNA content (in practice, it is the fraction of cells falling into the same histogram bin),

and $w$ values are weights. Fig. 4.12 shows this model graphically, fit to experimental data.

$$f_{G1}(T_c) = \frac{w_1}{\sigma\sqrt{2\pi}}e^{-\frac{(T_c-\mu_{G1})^2}{2\sigma^2}} \tag{4.26}$$

$$f_S(T_c) = \begin{cases} 0 & : T_c \notin [\mu_{G1}, \mu_{G2}] \\ v & : T_c \in [\mu_{G1}, \mu_{G2}] \end{cases} \tag{4.27}$$

$$f_{G2}(T_c) = \frac{w_2}{\sigma\sqrt{2\pi}}e^{-\frac{(T_c-\mu_{G2})^2}{2\sigma^2}} \tag{4.28}$$

While fitting to a cell cycle distribution model may be sufficient to identify biological outlier cells, I use two additional DNA features to further exclude low-quality images of nuclei. These features are the nuclear area and the coefficient of variation ($cv$) of intra-nuclear DNA intensity. The $cv$ is a rough proxy for the DNA texture, and so can be used to identify out-of-focus cells (lower $cv$) or those with chromatin condensation or punctate artifacts (higher $cv$). I therefore statistically gate the population by rejecting those cells that are too far from the median of either of these features (see Fig. 4.13).

Finally, I restrict my analyses to cells in the G1 phase (Fig. 4.13). The rationale for this is that I do not expect G1 and G2 cells to have different qualitative behaviors for the signaling pathways that I study in Chapter 3, though I do expect them to have somewhat different quantitative behaviors. The effect of combining these subpopulations would then be a meaningless increase in apparent signaling variability. I therefore chose the G1 population as it is typically more populated and is less prone to double-segmentation errors.

### 4.5.2   Using DNA features to correct measurement error

By using DNA features for quality control, we can thus collect all cells within the G1 and/or G2 populations that are high-quality (from an imaging standpoint) and accurately segmented. The quality control described above may be a sufficient level of data clean-up for some experimental goals, but it does not deal with the problem identified at the top of this section: that is, the separation of true biological variation from measurement error. In other words, quality control only discards cells that are too far from the "typical" cell, it does nothing to determine or correct the measurement error in those cells that are kept.

To identify biological variability, then, I again take advantage of the cell cycle "ground truth." As shown in Fig. 4.12, the theoretical cell cycle distribution has no variation in the G1 or G2 populations, as all cells have exactly the same diploid or tetraploid DNA content. The observed variation around the theoretical values then do not carry any biological meaning. (Note that this approximation becomes less accurate with chromosomally-unstable cell populations.)

**Figure 4.14:** The variation in total DNA content (within a single cell cycle peak) is non-biological and therefore should not be predictive of intensity values for other probes. **a**, Total Smad as a function of total DNA. The raw data show a low Pearson correlation coefficient (inset $r$ value) and linear regression slope (black line), which is over-corrected by multiplicative normalization (middle) and corrected by regression-based normalization (right). **b**, Comparison of single-cell values before (x-axis) and after (y-axis) regresson-based correction. This dataset has low correlation to total DNA, and so the change is small. $n = 689$ human colonic epithelial cells (G1 only, quality-controlled) immunostained with anti-Smad2/3 (Smad) and Hoechst (DNA). All $y$-axes on the same scale.

**Single-cell measurement error correction**

If the variation in total DNA carries no biological information, then it should not be predictive of of other feature values: any feature dependence on the DNA content must then be due to a global source of error. In principle, then, we can then reduce this error by removing the meaningless correlations of other features to total DNA content.

I take two single-cell level approaches to correcting total intensity feature errors (Fig. 4.14). The intuitive method is to estimate a normalization factor for e.g. all G1 cells using total DNA ($T_{\mathrm{DNA},c}$), and then use this factor to normalize the total intensity of other fluorescent probes ($T_{\mathrm{probe},c}$) in those same cells (Equation 4.29). Indeed, I have seen this approach used in the literature even without first restricting analysis to one cell cycle subpopulation. This method assumes a simple multiplicative relationship between the DNA and other channels such that, for example, a 10% increase in total DNA content would predict a 10% increase in different total probe intensity. Note that this assumption may not hold true, such that this method can cause over- or under-correction (as in 4.14a, middle).

$$f_{\mathrm{norm}}(T_{\mathrm{probe},c}) = \frac{\mathrm{median}_c(T_{\mathrm{DNA},c})}{T_{\mathrm{DNA},c}} T_{\mathrm{probe},c} \tag{4.29}$$

My preferred method is regression-based, as it guarantees removal of correlation between total Hoechst and the total intensity of another probe (Fig. 4.14a, right). For this method, linear regression is performed to get the function in Equation 4.30 with slope $m$ and intercept $b$. This results in a residual value ($r_{\mathrm{probe},c}$) for every cell. The value of each $T_{\mathrm{probe},c}$ can then be corrected by setting it equal to the median across all values plus the residual value for the same cell (Equation 4.31).

$$T_{\mathrm{probe},c} = f_{\mathrm{regression}}(T_{\mathrm{DNA},c}) = m T_{\mathrm{DNA},c} + r_{\mathrm{probe},c} + b \tag{4.30}$$

$$f_{\mathrm{norm}}(T_{\mathrm{probe},c}) = \mathrm{median}_c(T_{\mathrm{probe},c}) + r_{\mathrm{probe},c} \tag{4.31}$$

**Figure 4.15:** An observed total intensity distribution $f_{\text{observed}}$ is the result of convolution of the true distribution $f_{\text{true}}$ and the measurement error $f_{\text{error}}$. The error function for total nuclear intensity features can be estimated as having $\sigma_{\text{error}} \approx \sigma_{\text{DNA}}$. Cartoon, using distributions from Fig. 4.12.

For the sample data in Fig. 4.14 it is clear that there is low basal correlation between total DNA and total Smad, though I have observed much higher correlations in some datasets. I further note that this same rationale could be extended to non-DNA references and other features for cases in which cross-probe correlations are expected to be meaningless. For all datasets in this dissertation I apply the total DNA regression-based correction when accurate single-cell values are needed (such as for the calculation of mutual information between single-cell distributions).

**Population-level error correction**

While the single-cell correction above can be used to remove those measurement errors that are directly shared by each probe, it is reasonable to expect that much of measurement error is more random and so affects probes independently. Though it is not possible to correct single-cell values for such unpredictable error, we can apply correction at the population level.

An observed feature distribution can be modeled as the convolution of a true biological distribution with a measurement error distribution centered on zero, $f_{true} * f_{error}$. In the case of log-total DNA in G1 cells, $f_{true}$ is a delta distribution, $\delta_{true}$, positioned at $\mu_{G1}$. I can then take advantage of the property that $\delta_{true} * f_{error} = f_{error} + \mu_{G1}$. In other words, the G1 and G2 distributions are themselves estimates of the measurement error distribution, if their means are set to 0 (Fig. 4.15).

For distributions that are log-normal, as are the total intensity features for all probes used in this dissertation, I can also use the property that two convolved normal distributions yield a third normal distribution with mean $\mu_3 = \mu_1 + \mu_2$ and variance $\sigma_3^2 = \sigma_1^2 + \sigma_2^2$. Thus, I can estimate the "true" cell-to-cell total nuclear intensity variation for any probe by Equation 4.32. There of course may be other sources of error not compensated for in this way, and such an approach is only defensible for the total intensity feature.

$$\sigma_{\text{probe,true}} = \sqrt{\sigma_{\text{probe,observed}}^2 - \sigma_{\text{DNA,error}}^2} \tag{4.32}$$

What utility does this deconvolution have? Published measurements of cell-to-cell variability range from 15-30% [28], and my own raw data show values within this same range. However, these values include measurement error that is not being compensated for. Thus, cell-to-cell variability, as measured by microscopy, is necessarily overestimated. The above reasoning shows that this

overestimation is simple to measure, as all that is needed are the log-scale standard deviations of the G1/2 total DNA distributions and the standard deviations of the total-probe values in question.

I have not performed a comprehensive study of the size of this effect, though I have measured it in several independent datasets for various markers. I find that deconvolved total intensity distributions yield $\sim 10\%$ lower standard deviations and $\sim 10\%$ higher information content (measured by mutual information [16]). These inaccuracies are small enough that I feel comfortable stating that measurement errors in high-quality microscopy datasets are much smaller than true biological variation.

It is important to note that this DNA-based deconvolution makes the assumption that the sources of error are the same between Hoechst and other probes. It is possible that nuclear antibody-based stains have a partially non-overlapping set of error sources with small molecules like Hoechst. One possibility to address this, then, would be to use a non-specific secondary antibody in a free channel. That non-specific antibody should not be correlated with the specific antibody staining in other channels, and so any measured correlations could be removed using the same rationale as for DNA content-based correction. Unfortunately, I have had limited success with this approach, though non-specific secondary antibodies do indeed show high single-cell correlation. The problem has been that they also show unexpected properties, like differences in intracellular localization, for which adequate controls are not clear.

## 4.6 Discussion

Unlike other quantitative single-cell methods, such as flow cytometry, the subcellular resolution of image data allows for a stupefyingly large number of feature measurements, and there are no standardized practices for choosing or implementing these features. Further, identification of individual cells takes place at the level of software, not hardware. Quantitative imaging is thus an exceedingly difficult task outside of the labs that specialize in it, and no two of these labs are likely to converge on the exact same solutions.

With imaging we can directly see the beauty of the biology we are studying, and the high information content of images makes this type of data boundless in its potential utility. We are currently not meeting this potential, however, and I firmly believe that this is due to an absence of established standards and approaches that would make quantitative imaging more broadly accessible and interpretable. To that end, I hope that this chapter provides some intuition to those scientists who have not had the opportunity to work and think extensively about image data. Further, I hope that my approach to single-cell image analysis, described in this chapter and demonstrated in Chapter 3 for a specific biological study, provide useful demonstrations of the utility of quantitative single-cell analysis.

## 4.7   Methods

This chapter provides the imaging and image analysis details also used for the experiments in Chapter 3. The methods for cell culture, immunostaining, and imaging are left to that chapter (see Section 3.5). Specific methodological details for the figures in this chapter are primarily provided within the figure legends; this section provides the remaining information.

**Cell culture.** I followed the methods in Section 3.5, with the following additions. For Fig. 4.8, I used a fluorescently-labeled clone of the cell line A549 (ATCC #CCL-185). This clone contains a pSeg vector that co-expresses Cyan Fluorescent Protein fused to Drosophila Histone H2B to label nuclei and the Red Fluorescent Protein variant mCherry to label the whole cell (see Appendix A). The clone was generated by Jungseog Kang and Qi Wu (Altschuler & Wu lab, UT Southwestern).

**Imaging.** Plates were imaged on one of two Nikon Eclipse Ti-E2000 microscopes controlled by NIS Elements version 4, using several optical setups. Image coordinates were generated using NIS Elements or custom software for higher precision. The cameras used were either a Roper Scientific CoolSnap HQ2 CCD 14-bit Fig. 4.5 or an Andor Zyla sCMOS 11-bit.

# Chapter 5

# Conclusion

The study of cellular signaling is a difficult one, in large part because we know so little, *a priori*, about what aspects of a signal a cell cares about and into what intracellular properties it encodes those signals. It may be, however, that the complexity of our understanding of signaling is hiding a reality of underlying simplicity (Chapter 1). In particular, the static network models that we rely on include links derived from different timescales (such that activity along one link should be considered constant relative to activity along another) as well as links that may only be present in a subset of true signaling network instances (e.g. in certain cell types or experimental conditions).

Wnt/β-catenin and TGFB$_{sf}$ together provide a case study for this problem of signaling complexity (Chapter 2). For each pathway alone there is a lack of clarity due to idiosyncratic outcomes between labs and experiments. There is even less clarity when considering how the two pathways are integrated by cells in order for them to make decisions. Much of the confusion regarding crosstalk between these pathways may stem from the reliance of published work on overexpression assays, which can push cells into abnormal states. Further, the same studies rely on transcription-based readouts that are measured long after initial pathway activation, allowing for transcriptional feedback to confound the results (Section 2.6).

Therefore, I tested Wnt/β-catenin and TGFB$_{sf}$ for crosstalk using an assay that relies on endogenous levels of proteins (to prevent the signaling networks from being pushed into abnormal states), short timecourses (to prevent transcriptional feedback from confounding interpretation), and direct readouts of signal transduction (Chapter 3). This study design allowed me to directly test the affect of signaling through one pathway on signaling through the other. In this way, I discovered that Wnt/β-catenin and TGFB$_{sf}$ are in fact completely insulated from one another during signal transduction. Further, I found no evidence for the expected intra-TGFB$_{sf}$ inhibition nor the competition for the shared component, Smad4, thought to cause it.

Though I observed that signaling insulation is a general phenomenon, I did uncover an instance of context-dependent transcriptional crosstalk. Importantly, the resulting transcriptional feedback led to biasing of signaling activity over longer timecourses. In effect, this is a direct example of the fact that some network links may exist only under certain conditions, and that experimental separation of signaling and transcription is essential to understanding the nature of morphogenic

106

pathway integration. From these results, I conclude that there exists a surprising simplicity of interactions between morphogenic pathways: there is a high degree of insulation between Wnt/β-catenin and $TGFB_{sf}$, allowing nuclei to make transcriptional decisions based on more-complete models of their environments (Section 3.4).

For the experimental approach that led to my discoveries, I relied heavily on quantitative, single-cell fluorescence microscopy (Chapter 4). Single-cell imaging is an increasingly available and invaluable approach to studying cell biology (and cell signaling in particular). Unfortunately, single-cell image analysis is a difficult and unsolved problem. An oft-neglected issue is that of image correction, prior to analysis. Indeed, the most commonly-used methods throughout the literature tend to provide incomplete or error-prone image correction. I discovered a solution to this problem, in the particular case of high-throughput microscopy, which takes advantage of predictable optical properties within microtiter plates (Section 4.3).

**On negative results**

When I first found that Wnt3A and TGFβ3 are insulated from one another during signaling, I was hugely disappointed. This disappointment came, in large part, from having spent a year under the impression that there actually was crosstalk. I had performed many experiments and narrowed down the point of inter-pathway crosstalk to the receptor or ligand level, before I finally discovered contamination of the purified recombinant Wnt3A ligand (Fig. 3.6).

The artifactual result itself was not the most disappointing part; it was instead the feeling that I now had a "negative result" and would therefore have to drop the project completely and start over. This attitude had been ingrained within me throughout my science education: I was taught that scientific discovery is about finding new and exciting things. Or, in any event, that my discoveries would need to be new and exciting in order to get published and move my career forward.

While scrambling to make the best out of the situation, I slowly realized that there was not truly a problem. The fact that the crosstalk I observed was an artifact was a real result: numerous studies show signaling interactions between these pathways, and the contaminated reagent is among the most commonly used for studying Wnt/β-catenin signaling. As long as I could unambiguously demonstrate the absence of crosstalk, then my result was not a negative one at all and could potentially reduce some confusion in the field. One person's negative result is another person's positive result, as it were.

In any case, there is value to negative results. In recent years there has been increasing hubbub in the news and in opinion pieces regarding the state of biomedical research. Of particular concern has been the irreproducibility of published results and the toxic effects of using "impact" as the primary metric for determining the merit of both scientific discoveries and of the scientists who make those discoveries [260,261]. Why does the emphasis on impact lead to irreproducibility? A needed statistic when evaluating a claim is the *a priori* likelihood of that claim being true; when the literature is biased towards exciting (i.e. unlikely) results because scientists do not publish "negative" results, we cannot accurately estimate the likelihood of truth for new claims. This problem is (in)famously

addressed in a paper by John Ioannidis that should be required reading for all biologists [262].

**Parting thoughts**

The overarching field of cellular signaling is an exciting one; understanding how cells communicate and process information is absolutely fundamental to our understanding of all of cell biology. Importantly, a deep understanding of this topic will allow us to more ably manipulate biological systems, both in medicine and in bioengineering. To develop that deep understanding, however, it is essential that we revisit the complex static maps of established cellular signaling pathways, and then pay careful attention to how time and context reshape these networks. By identifying underlying simplicity in cell signaling, we will be better equipped to both understand and control cellular information processing and decision-making.

# Appendix A

# pSeg: plasmids for live-cell segmentation

This dissertation uses only fixed-cell microscopy assays, but future studies would ideally include live-cell studies of TGFB$_{sf}$ and Wnt signaling. Live-imaging and fixed-cell imaging can use the same general analytical methods, including the use of nuclear and cytosolic fluorescence to identify these two cellular compartments. However, finding appropriate labels for live-cell assays is much more difficult, as the cells are not permeabilized and toxic stains must be avoided. In order to effectively deploy segmentation algorithms, the nucleus and cytosol should ideally be labeled brightly and homogeneously.



**Figure A.1:** The pSeg plasmids generally consist of two differentially-localized fluorophores. Left, sample images from two human colonic epithelial cell (HCEC) pSeg clones. The top image is from the LYiHnCL pSeg construct shown at right, and the bottom image is from an LRiHnCL construct. See Table A.1 for symbols. Right, structure of pSeg construct. The 5' and 3' Murine Leukemia Virus LTRs (long terminal repeats) flank the construct, which also has loxP sites oriented in the same direction so that the construct can be excised by Cre recombinase, thus reverting a clone to a "parental" phenotype. AmpR, ampicillin resistance gene. psi+, viral packaging sequences. Puro, puromycin. IRES, internal ribosome entry site.

**Table A.1:** List symbols for the pSeg library in Table A.2.

| Symbol | Abbreviation | Full |
|--------|--------------|------|
| C | CFP | mCerulean |
| G | GFP | enhanced Green Fluorescent Protein |
| H | H2B | Histone H2B (*Drosophila*) |
| i | IRES | Internal ribosome entry site (encephalomyocarditis virus) |
| L | loxP | loxP site |
| n | NLS | Nuclear Localization Signal (SV40) |
| m | | membrane- localized Gap43 N-terminus (*Danio rerio*) |
| R | RFP | mCherry |
| Y | YFP | enhanced Yellow Fluorescent Protein |

One approach to address this problem is via "central dogma" (CD) tagging, in which a fluorescent-protein-encoding gene surrounded by splice acceptor/donor sites is randomly integrated into the genome. Any integration that lands within an intron of a gene has a chance to become an exon for that same gene. Thus,the genetic protein product will then contain a fluorescent protein as one of its domains, allowing the localization and dynamics of that protein to be visualized [263,264].

High-throughput studies have used this technique to obtain nucleus- or cytosol-localized labels for the purpose of segmentation, with the purported benefit that it requires no exogenous expression system and thus minimizes perturbations to the cell [28,265–267]. Unfortunately, the CD-tagging process is slow, since exon-generating integration events are rare and must be subsequently screened to identify the rarer-still events that yield the desired localization patterns. Additionally, the tagging of a protein is itself a perturbation that may disrupt localization or function. Chosen CD-tagged clones must then be carefully studied to ensure that the cells (and the tagged proteins) behave properly.

A much faster alternative is to express fluorescent proteins from exogenous promoters, though there is always the concern that the presence of an exogenous promoter may affect the cell in some unpredictable way. To me, this is just as problematic as the concern that CD-tagging will disrupt some unknown function of the tagged protein. I therefore prefer the efficiency of the exogenous approach. To take advantage of this efficiency for future live-cell projects, I created a library of viral constructs, each generating distinct localization patterns of red, yellow, and cyan fluorescent protein expression in mammalian cells. I refer to these as "segmentation plasmids" (pSegs).

**pSeg structure**

Each pSeg plasmid has three variable components: a fluorescent protein (FP) with a whole-cell expression pattern, a fluorescent protein with a nuclear expression pattern, and a selection cassette. I chose the mCherry, eYFP, and Cerulean fluorescent proteins because they have wide spectral separation and so can all be used simultaneously in single cells. The general structure of the insert is shown in Fig. A.1. Whole-cell expression patterns are generated by either free FP or a FP fused to the Gap43 membrane localization signal [268]. This signal results in faint whole-cell and strong membrane fluorescence. Nuclear expression is mediated by addition of either a nuclear localization sequence (NLS) or fusion to *Drosophila* histone H2B. The H2B-fused FP is useful because it can be used to monitor chromatin condensation [269] and seems to be benign even when generally expressed

**Table A.2:** Full list of plasmids in the pSeg library. See Table A.1 for symbol meanings. "Floxed" refers to flanking by loxP. PuroR and NeoR refer to puromycin or neomycin resistance cassettes.

| floxed & PuroR | floxed & NeoR | NeoR |
|---|---|---|
| LmYiHnCL | LYiHnCL | YiHnC |
| LmCiHnCL | LGiHnCL | YiHnY |
| LmRiHnCL | LRiHnCL | CiHnC |
| LmYiCnnCL | LYiCnnCL | CiHnY |
| LmCiCnnCL | LGiCnnCL | YiCnnC |
| LmRiCnnCL | LRiCnnCL | YiYnnY |
| LYiHnCL | | CiCnnC |
| LGiHnCL | | |
| LRiHnCL | | |
| LYiCnnCL | | |
| LGiCnnCL | | |
| LRiCnnCL | | |

in mice [270].

Each construct is expressed as a single mRNA, using the retroviral promoter and an internal ribosome entry site (IRES) for differential translation of the the protein products. Finally, the constructs are also flanked by loxP sites and so can be removed after genomic integration, thus generating unlabeled "parental" cell lines. The pSeg library consists of 25 color/localization/selection combinations in Puromycin or Neomycin-resistant backgrounds, with or without flanking loxP sites. Table A.2 shows all of these constructs (refer to Table A.1 for symbol definitions).

I built the library into the pMYs retroviral backbone, which contains a modified version of the Murine Leukemia Virus promoter that improves integration efficiency in some stem cell systems [271]. The pMYs backbone contains the viral long terminal repeats (LTRs) and a viral packaging signal (Table A.1), allowing for the creation of virus using MLV packaging cell lines.

**Generating cellular pSeg clones**

Integration of the pSeg constructs into target cells requires the generation of and infection by live virus (see Methods). A few days after infection, I FACS (Flow-assisted cell sorting)-sort single cells into each well of 384-well plates. I find that survival rates for singly sorted cells are quite low (∼20-30%), though I could double these rates by pre-seeding the wells with un-labeled cells. The feeder cells can be selected against later according to the chosen pSeg cassette. With this appraoch, a clone from the A549 cell line was created in the Altschuler & Wu lab by Jungseog Kang and Qi Wu, which forms the foundation of a CD-tagged library (unpublished).

The slowest part of the process is the establishment of clonal cell lines from single cells, which can take 6 weeks. However, this step is essential. Cell-to-cell variability can be high, even within a clonal population [27]. Further, I found that all pSegs yield populations with diverse localization patterns. For example, the LRiHnCL construct should express mCherry and Cerulean in the nucleus, and yet all possible patterns are observed (Fig. A.2). The reason for this is not clear, though retroviruses (such as MLV) have diploid genomes within viral particles and are known to have recombination events [272]. The high degree of repetition within these constructs (the different fluorescent proteins have high homology) may then lead to truncations or expansions of the integrated sequence, such

**Figure A.2:** A population of cells infected with the pSeg.puro-LRiHnCL construct show all possible nuclear localization patterns. Percents show the relative population size for each group. The bottom-left (black) cells are unlabeled, the top-right (dark blue) cells are correctly labeled with both fluorescent proteins, and the other two populations express only one of the two proteins. $n=1770$ human colonic epithelial cells, stained with Hoechst for segmentation. Arbitrary total fluorescence units.

that it re-arranges the localization signals on the fluorescent proteins.

Finally, I warn that the Murine Leukemia Virus has a strong preference for integration into the 5' end of highly-expressed genes [273–275]. Despite this preference, the virus does seem to avoid house-keeping genes, though I note that this may simply be due to the death of cells that do have integrations at such sites. As a consequence of integration, then, functionally important genes may be disrupted. A simple test for this is to select multiple labeled clones and compete each against the unlabeled parental population. The clones with the most similar-to-parental growth rates can then be kept for further experiments.

## Methods

**Molecular cloning.** I constructed the pSeg library using a combination of standard moleculuar cloning techniques and Gibson assembly (New England Biolabs #E5510) [276]. I verified each construct by Sanger sequencing, and confirmed each phenotype by checking the localization patterns in transiently-transfected HEK293T cells.

**Tissue culture.** The tissue culture methods and imaging followed those in Section 3.5 with the exception of viral production and infection. I used the Platinum-A HEK293T packaging cell line (Cell Biolabs #RV-102) and closely followed the protocol used by Uri Alon's group for CD-tagging [266]. In brief, virus is generated by transfecting packaging cells and collecting the supernatant 2-3 days later. This supernatant is then applied to the target cells, which are left for 48-72 hours to allow for genomic integration. The multi-day timeline is essential, as MLV can only integrate into immediately-post-mitotic genomes [277].

**Analysis.** For the analysis in Fig. A.2 I used the ImageJ implementation of the rolling ball algorithm for correction and a custom Matlab nuclear threshold segmentation algorithm. I used

the Hoechst channel to manually gate the G1 population and to subsequently correct the mCherry and Cerulean channels using the regression method described in Section 4.5.2. I used the k-means algorithm implemented in R to automatically identify each of the subpopulations shown in the figure.

# Glossary

**canonical pathway** A well-established series of biochemical signaling steps, which effectively pass information from one molecule to another. Often defined by genetic means with epistasis mapping. Also used to refer to pathways that are more easily studied than alternatives (e.g. compare canonical and non-canonical Wnt signaling).

**canonical Wnt signaling** The branch of Wnt signaling that results in increased intracellular β-catenin levels (therefore also called Wnt/β-catenin signaling. This form of Wnt signaling is much easier to study than the others, because β-catenin is easy to measure and is relatively insulated from other signaling pathways.

**channel** In the imaging sense, used to refer to a fluorescence color channel (e.g. Hoechst and fluorescein fluoresce in different channels). In the information-carrying sense, a channel is a distinct path of information flow.

**crosstalk** A transfer of information between signaling components of canonically distinct pathways.

**decision-making** The mapping of an internal model of a signaling event to some response. An example internal model might be the nuclear concentration of a transcription factor, while the decision is then how much of some transcriptional target to produce.

**detector** The camera used to acquire fluorescence images. It will typically have a constant baseline value added to images that must be subtracted during image correction.

**DNA** Deoxyribonucleic acid (only a massochist would use this acronym for something else). Used also when referring to the total fluorescent Hoechst signal from stained cells, as this molecule intercalates into the DNA backbone and thus serves as a proxy for DNA content.

***Drosophila melanogaster*** The classic fruit fly model system. If you are a biologist, you cannot be forgiven for having to look up this term.

**edge** A link between nodes (borrowed from graph theory). For a signaling network, indicates some interaction (e.g. phosphorylation) or transfer of information between nodes.

**encoding** The manner in which information is converted from one type to another. For example a piece of text can be encoded into a binary format, or the extracellular concentration of a ligand can be encoded into the nuclear concentration of a transcription factor.

**endogenous** Usually used to refer to the normal cellular concentrations of some factor (contrast to overexpression and exogenous).

**exogenous** Usually used to refer to an addition to the normal concentrations of some factor (contrast to endogenous). For example, an added purified ligand is an exogenous source of that ligand, and an overexpressed protein creates an exogenous pool in addition to the endogenous pool.

**feature** A single type of measurement in image analysis. Example features include nuclear area or total cytosolic intensity.

***Homo sapiens*** If you are reading this, you are probably one of these.

**homology** Having a shared evolutionary ancestor. Thus genes with high homology have similar sequences and recent ancestry (or high selective pressure). It is important to note that sequence similarity does not necessarily imply homology (i.e. "homologous" does not mean "similar").

114

The noun form of this term is "homolog," which is a more general term than are paralog and ortholog.

**image correction** The removal of detector, background, and shading components from an image.

**immunostain** The use of an antibody to attach a fluorophore, or other measurable item, to a target molecule.

**knockout** The removal of a gene from the genome. Typically used to refer to the removal of both alleles in a diploid organism, though the term "homozygous knockout" means this more specifically. Thus a "Wnt5A knockout mouse" likely lacks both endogenous alleles of Wnt5A.

**LiCl** Lithium Chloride. Can be used to inhibit GSK3β.

**ligand** A protein or other molecule that is recognized by a cellular receptor, thus leading to some internal representation of properties that molecule.

**microenvironment** A term that is thrown around in the literature extensively but frequently (and perhaps purposely) left undefined. Here, it is used to refer to the collection of environmental parameters that a single cell is exposed to. A subset of such parameters include juxtacrine and paracrine signals from neighboring cells, as well as any more global properties (e.g. temperature). In the context of an experiment, this also includes any perturbations that a cell should be able to sense.

**node** A conceptual unit that may interact with another unit (borrowed from graph theory). For a signaling network, may indicate a protein or protein state.

**non-canonical Wnt signaling** The branches of Wnt signaling that do not result in increased intracellular β-catenin levels (in the longer-term, these pathways may inhibit canonical Wnt signaling). This form of Wnt signaling has been difficult to study, because its readouts (including transient $Ca^{2+}$ signaling) are hard to measure and are integrated with many other signaling pathways. For these reasons, it is not clear how many non-canonical pathways there are, how distinct one is from another, and what the impacts of this

form of signaling are on canonical signaling (besides general long-term inhibition).

**ortholog** Homologous sequences, in two different organisms, that resulted from a speciation event (i.e. they are the "same" sequence).

**overexpression** The exogenous expression of some protein, for example by transient expression from a plasmid, constitutive expression by a genomically-integrated viral construct, or controllable expression from an inducible promoter. Generally implies that there is more than a normal quantity of the expressed protein.

**paneth cell** A long-lived cell type living in the base of crypts. It is thought to provide the stem cell niche in the small intestine.

**paralog** Homologous sequences, within the same genome, that resulted from a gene duplication event.

**probe** Short for "fluorescent probe," used to refer to a fluorescent small molecule or antibody that binds to a specific molecular target.

**proteosome** A large protein complex that degrades proteins in a highly regulated manner, typically after those proteins have been modified by the covalent addition of ubiquitin.

**RNA** Ribonucleic acid. Used with various prefixes to indicate the specific type of this molecule. Types include messenger RNA, small interfering RNA, and ribosomal RNA.

**signaling** Also referred to throughout the text as "cellular signaling" and "signal transduction." I use this term specifically to refer to the process by which an external stimulus is encoded into an internal representation of that stimulus.

**transcription factor** A molecule (typically a protein) that directly binds to DNA, or to other molecules that bind DNA, and thus can cause a change in the transcription rate of a gene.

***Trichoplax adhaerens*** The most basal known metazoan, with a simple multicellular structure.

***Xenopus laevis*** A frog used as a model organism. Wnt and BMP have been heavily studied in this organism, particularly with respect to development.

[1] M. R. Bennett, W. L. Pang, N. A. Ostroff, B. L. Baumgartner, S. Nayak, L. S. Tsimring, and J. Hasty, "Metabolic gene regulation in a dynamically changing environment.," *Nature*, vol. 454, pp. 1119–22, Aug. 2008.

[2] M. Acar, J. T. Mettetal, and A. van Oudenaarden, "Stochastic switching as a survival strategy in fluctuating environments.," *Nature genetics*, vol. 40, pp. 471–5, Apr. 2008.

[3] M. Natarajan, K.-M. Lin, R. C. Hsueh, P. C. Sternweis, and R. Ranganathan, "A global analysis of cross-talk in a mammalian cellular signalling network.," *Nature cell biology*, vol. 8, pp. 571–80, June 2006.

[4] G. Balázsi, A. van Oudenaarden, and J. J. Collins, "Cellular decision making and biological noise: from microbes to mammals.," *Cell*, vol. 144, pp. 910–25, Mar. 2011.

[5] B. N. Kholodenko, "Cell-signalling dynamics in time and space.," *Nature reviews. Molecular cell biology*, vol. 7, pp. 165–76, Mar. 2006.

[6] T. Ideker and N. J. Krogan, "Differential network biology.," *Molecular systems biology*, vol. 8, p. 565, Jan. 2012.

[7] C.-J. Ku, Y. Wang, O. D. Weiner, S. J. Altschuler, and L. F. Wu, "Network crosstalk dynamically changes during neutrophil polarization.," *Cell*, vol. 149, pp. 1073–83, May 2012.

[8] K. A. Janes and D. A. Lauffenburger, "Models of signalling networks - what cell biologists can gain from them and give to them," *Journal of Cell Science*, vol. 126, pp. 1913–1921, May 2013.

[9] E. R. Zhang, L. F. Wu, and S. J. Altschuler, "Envisioning migration: mathematics in both experimental analysis and modeling of cell behavior.," *Current opinion in cell biology*, vol. 25, pp. 538–42, Oct. 2013.

[10] D. Pe'er and N. Hacohen, "Principles and strategies for developing network models in cancer.," *Cell*, vol. 144, pp. 864–73, Mar. 2011.

[11] O. Shoval and U. Alon, "SnapShot: network motifs.," *Cell*, vol. 143, pp. 326–e1, Oct. 2010.

[12] Y. Wang, C.-J. Ku, E. R. Zhang, A. B. Artyukhin, O. D. Weiner, L. F. Wu, and S. J. Altschuler, "Identifying network motifs that buffer front-to-back signaling in polarized neutrophils.," *Cell reports*, vol. 3, pp. 1607–16, May 2013.

[13] L. Goentoro, O. Shoval, M. Kirschner, and U. Alon, "The incoherent feedforward loop can provide fold-change detection in gene regulation," *Molecular cell*, vol. 36, no. 5, pp. 894–899, 2009.

[14] J. Dutkowski and T. Ideker, "Protein Networks as Logic Functions in Development and Cancer," *PLoS Computational Biology*, vol. 7, p. e1002180, Sept. 2011.

[15] V. Becker, M. Schilling, J. Bachmann, U. Baumann, A. Raue, T. Maiwald, J. Timmer, and U. Klingmüller, "Covering a broad dynamic range: information processing at the erythropoietin receptor.," *Science (New York, N.Y.)*, vol. 328, pp. 1404–8, June 2010.

[16] R. Cheong, A. Rhee, C. J. Wang, I. Nemenman, and A. Levchenko, "Information transduction capacity of noisy biochemical signaling networks.," *Science (New York, N.Y.)*, vol. 334, pp. 354–8, Oct. 2011.

[17] O. Shoval, L. Goentoro, Y. Hart, A. Mayo, E. Sontag, and U. Alon, "Fold-change detection and scalar symmetry of sensory input fields.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 107, pp. 15995–6000, Sept. 2010.

[18] R. E. C. Lee, S. R. Walker, K. Savery, D. A. Frank, and S. Gaudet, "Fold Change of Nuclear NF-$\kappa$B Determines TNF-Induced Transcription in Single Cells," *Molecular cell*, Feb. 2014.

[19] S. Tay, J. J. Hughey, T. K. Lee, T. Lipniacki, S. R. Quake, and M. W. Covert, "Single-cell NF-kappaB dynamics reveal digital activation and analogue information processing.," *Nature*, vol. 466, pp. 267–71, July 2010.

[20] J. Kang, B. Xu, Y. Yao, W. Lin, C. Hennessy, P. Fraser, and J. Feng, "A Dynamical Model Reveals Gene Co-Localizations in Nucleus," *Gene*, vol. 7, no. 7, 2011.

[21] B. Snijder and L. Pelkmans, "Origins of regulated cell-to-cell variability.," *Nature reviews. Molecular cell biology*, vol. 12, pp. 119–25, Feb. 2011.

[22] S. J. Altschuler and L. F. Wu, "Cellular heterogeneity: do differences make a difference?," *Cell*, vol. 141, pp. 559–63, May 2010.

[23] S. Huang, "Non-genetic heterogeneity of cells in development: more than just noise.," *Development*, vol. 136, pp. 3853–62, Dec. 2009.

[24] J. E. Ferrell and E. M. Machleder, "The biochemical basis of an all-or-none cell fate switch in Xenopus oocytes.," *Science (New York, N.Y.)*, vol. 280, pp. 895–8, May 1998.

[25] L.-H. Loo, H.-J. Lin, D. K. Singh, K. M. Lyons, S. J. Altschuler, and L. F. Wu, "Heterogeneity in the physiological states and pharmacological responses of differentiating 3T3-L1 preadipocytes.," *The Journal of Cell Biology*, vol. 187, pp. 375–84, Nov. 2009.

[26] P. Cluzel, M. Surette, and S. Leibler, "An ultrasensitive bacterial motor revealed by monitoring signaling proteins in single cells.," *Science (New York, N.Y.)*, vol. 287, pp. 1652–5, Mar. 2000.

[27] D. K. Singh, C.-J. Ku, C. Wichaidit, R. J. Steininger, L. F. Wu, and S. J. Altschuler, "Patterns of basal signaling heterogeneity can distinguish cellular populations with different drug sensitivities.," *Molecular Systems Biology*, vol. 6, p. 369, May 2010.

[28] A. Sigal, R. Milo, A. Cohen, N. Geva-Zatorsky, Y. Klein, Y. Liron, N. Rosenfeld, T. Danon, N. Perzov, and U. Alon, "Variability and memory of protein levels in human cells.," *Nature*, vol. 444, pp. 643–6, Nov. 2006.

[29] 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*, vol. 23, pp. 1870–5, Aug. 2009.

[30] T. Kalmar, C. Lim, P. Hayward, S. Muñoz Descalzo, J. Nichols, J. Garcia-Ojalvo, and A. Martinez Arias, "Regulated fluctuations in nanog expression mediate cell fate decisions in embryonic stem cells.," *PLoS Biology*, vol. 7, p. e1000149, July 2009.

[31] S. L. Spencer, S. Gaudet, J. G. Albeck, J. M. Burke, and P. K. Sorger, "Non-genetic origins of cell-to-cell variability in TRAIL-induced apoptosis.," *Nature*, vol. 459, pp. 428–32, May 2009.

[32] P. Navarro, N. Festuccia, D. Colby, A. Gagliardi, N. P. Mullin, W. Zhang, V. Karwacki-Neisius, R. Osorno, D. Kelly, M. Robertson, and I. Chambers, "OCT4/SOX2-independent Nanog autorepression modulates heterogeneous Nanog gene expression in mouse ES cells," *The EMBO Journal*, vol. advance on, Nov. 2012.

[33] H. H. Chang, M. Hemberg, M. Barahona, D. E. Ingber, and S. Huang, "Transcriptome-wide noise controls lineage choice in mammalian progenitor cells.," *Nature*, vol. 453, pp. 544–7, May 2008.

[34] M. B. Elowitz, A. J. Levine, E. D. Siggia, and P. S. Swain, "Stochastic gene expression in a single cell.," *Science (New York, N.Y.)*, vol. 297, pp. 1183–6, Aug. 2002.

[35] B. Munsky, G. Neuert, and A. van Oudenaarden, "Using gene expression noise to understand gene regulation.," *Science (New York, N.Y.)*, vol. 336, pp. 183–7, Apr. 2012.

[36] A. Sanchez and I. Golding, "Genetic determinants and cellular constraints in noisy gene expression.," *Science (New York, N.Y.)*, vol. 342, pp. 1188–93, Dec. 2013.

[37] A. Eldar and M. B. Elowitz, "Functional roles for noise in genetic circuits.," *Nature*, vol. 467, pp. 167–73, Sept. 2010.

[38] J.-W. Veening, W. K. Smits, and O. P. Kuipers, "Bistability, epigenetics, and bet-hedging in bacteria.," *Annual review of microbiology*, vol. 62, pp. 193–210, Jan. 2008.

[39] M. D. Slack, E. D. Martinez, L. F. Wu, and S. J. Altschuler, "Characterizing heterogeneous cellular responses to perturbations.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 105, pp. 19306–11, Dec. 2008.

[40] K. a. Janes, S. Gaudet, J. G. Albeck, U. B. Nielsen, D. a. Lauffenburger, and P. K. Sorger, "The response of human epithelial cells to TNF involves an inducible autocrine cascade.," *Cell*, vol. 124, pp. 1225–39, Mar. 2006.

[41] N. Domedel-Puig, P. Rué, A. J. Pons, and J. García-Ojalvo, "Information Routing Driven by Background Chatter in a Signaling Network," *PLoS Computational Biology*, vol. 7, p. e1002297, Dec. 2011.

[42] R. C. Hsueh, M. Natarajan, I. Fraser, B. Pond, J. Liu, S. Mumby, H. Han, L. I. Jiang, M. I. Simon, R. Taussig, and P. C. Sternweis, "Deciphering signaling outcomes from a system of complex networks," *Science signaling*, vol. 2, p. ra22, Jan. 2009.

[43] T. Gregor, D. W. Tank, E. F. Wieschaus, and W. Bialek, "Probing the limits to positional information.," *Cell*, vol. 130, pp. 153–64, July 2007.

[44] J. O. Dubuis, G. Tkacik, E. F. Wieschaus, T. Gregor, and W. Bialek, "Positional information, in bits.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 110, pp. 16301–8, Oct. 2013.

[45] J. O. Dubuis, R. Samanta, and T. Gregor, "Accurate measurements of dynamics and reproducibility in small genetic networks.," *Molecular systems biology*, vol. 9, p. 639, Jan. 2013.

[46] L.-h. Loo, H.-j. Lin, R. J. S. Iii, Y. Wang, L. F. Wu, and J. Steven, "On an Approach for Extensibly Profiling the Molecular States of Cellular Subpopulations," *Systems Biology*, vol. 6, no. 10, pp. 1–17, 2010.

[47] V. Li, S. Ng, P. Boersema, T. Low, W. Karthaus, J. Gerlach, S. Mohammed, A. Heck, M. Maurice, T. Mahmoudi, and H. Clevers, "Wnt Signaling through Inhibition of $\beta$-Catenin Degradation in an Intact Axin1 Complex," *Cell*, vol. 149, pp. 1245–1256, June 2012.

[48] A. R. Hernández, A. M. Klein, and M. W. Kirschner, "Kinetic responses of $\beta$-catenin specify the sites of Wnt control.," *Science (New York, N.Y.)*, vol. 338, pp. 1337–40, Dec. 2012.

[49] J. Massagué, "TGF$\beta$ signalling in context.," *Nature reviews. Molecular cell biology*, vol. 13, pp. 616–630, Sept. 2012.

[50] B. Schmierer and C. S. Hill, "TGFbeta-SMAD signal transduction: molecular specificity and functional flexibility.," *Nature reviews. Molecular cell biology*, vol. 8, pp. 970–82, Dec. 2007.

[51] A. von Bubnoff and K. W. Cho, "Intracellular BMP signaling regulation in vertebrates: pathway or network?," *Developmental biology*, vol. 239, pp. 1–14, Nov. 2001.

[52] M. Srivastava, E. Begovic, J. Chapman, N. H. Putnam, U. Hellsten, T. Kawashima, A. Kuo, T. Mitros, A. Salamov, M. L. Carpenter, A. Y. Signorovitch, M. A. Moreno, K. Kamm, J. Grimwood, J. Schmutz, H. Shapiro, I. V. Grigoriev, L. W. Buss, B. Schierwater, S. L. Dellaporta, and D. S. Rokhsar, "The Trichoplax genome and the nature of placozoans.," *Nature*, vol. 454, pp. 955–60, Aug. 2008.

[53] L. Huminiecki, L. Goldovsky, S. Freilich, A. Moustakas, C. Ouzounis, and C.-H. Heldin, "Emergence, development and diversification of the TGF-beta signalling pathway within the animal kingdom.," *BMC evolutionary biology*, vol. 9, p. 28, Jan. 2009.

[54] J. Massague, "The transforming growth factor-beta family," *Annual Review of Cell Biology*, vol. 6, pp. 597–641, 1990.

[55] J. Massagué, "How cells read TGF-beta signals.," *Nature reviews. Molecular cell biology*, vol. 1, pp. 169–78, Dec. 2000.

[56] T. U. Wagner, "Bone morphogenetic protein signaling in stem cells–one signal, many consequences.," *The FEBS journal*, vol. 274, pp. 2968–76, June 2007.

[57] D. Umulis, M. B. O'Connor, and S. S. Blair, "The extracellular regulation of bone morphogenetic protein signaling.," *Development (Cambridge, England)*, vol. 136, pp. 3715–28, Nov. 2009.

[58] J. Massague, "Smad transcription factors," *Genes & Development*, vol. 19, pp. 2783–2810, Dec. 2005.

[59] M. Ehrlich, D. Horbelt, B. Marom, P. Knaus, and Y. I. Henis, "Homomeric and heteromeric complexes among TGF-$\beta$ and BMP receptors and their roles in signaling.," *Cellular signalling*, vol. 23, pp. 1424–32, Sept. 2011.

[60] L. M. Wakefield and C. S. Hill, "Beyond TGF$\beta$: roles of other TGF$\beta$ superfamily members in cancer.," *Nature reviews. Cancer*, vol. 13, pp. 328–41, May 2013.

[61] K. Miyazono, S. Maeda, and T. Imamura, "BMP receptor signaling: transcriptional targets, regulation of signals, and signaling crosstalk.," *Cytokine & growth factor reviews*, vol. 16, pp. 251–63, June 2005.

[62] N. Khalil, "TGF-$\beta$: from latent to active," *Microbes and Infection*, vol. 1, no. 15, pp. 1255–1263, 1999.

[63] A. Nohe, "Signal transduction of bone morphogenetic protein receptors," *Cellular Signalling*, vol. 16, pp. 291–299, Mar. 2004.

[64] M. Shi, J. Zhu, R. Wang, X. Chen, L. Mi, T. Walz, and T. A. Springer, "Latent TGF-$\beta$ structure and activation.," *Nature*, vol. 474, pp. 343–9, June 2011.

[65] I. B. Robertson and D. B. Rifkin, "Unchaining the beast; insights from structural and evolutionary studies on TGF$\beta$ secretion, sequestration, and activation.," *Cytokine & growth factor reviews*, vol. null, July 2013.

[66] S. Radaev, Z. Zou, T. Huang, E. M. Lafer, A. P. Hinck, and P. D. Sun, "Ternary complex of transforming growth factor-beta1 reveals isoform-specific ligand recognition and receptor recruitment in the superfamily.," *The Journal of biological chemistry*, vol. 285, pp. 14806–14, May 2010.

[67] L. Zakin and E. M. De Robertis, "Extracellular regulation of BMP signaling.," *Current biology : CB*, vol. 20, pp. R89–92, Feb. 2010.

[68] C. Sieber, J. Kopf, C. Hiepen, and P. Knaus, "Recent advances in BMP receptor signaling.," *Cytokine & growth factor reviews*, vol. 20, pp. 343–55, Jan. 2009.

[69] A. Bandyopadhyay, P. S. Yadav, and P. Prashar, "BMP signaling in development and diseases: a pharmacological perspective.," *Biochemical pharmacology*, vol. 85, pp. 857–64, Apr. 2013.

[70] P. J. Hart, S. Deep, A. B. Taylor, Z. Shu, C. S. Hinck, and A. P. Hinck, "Crystal structure of the human TbetaR2 ectodomain–TGF-beta3 complex.," *Nature structural biology*, vol. 9, pp. 203–8, Mar. 2002.

[71] J. Groppe, C. S. Hinck, P. Samavarchi-Tehrani, C. Zubieta, J. P. Schuermann, A. B. Taylor, P. M. Schwarz, J. L. Wrana, and A. P. Hinck, "Cooperative assembly of TGF-beta superfamily signaling complexes is mediated by two disparate mechanisms and distinct modes of receptor binding.," *Molecular cell*, vol. 29, pp. 157–68, Feb. 2008.

[72] R. Derynck and Y. Zhang, "Smad-dependent and Smad-independent pathways in TGF-beta family signalling," *Nature*, vol. 4, pp. 577–84, 2003.

[73] P. ten Dijke and C. S. Hill, "New insights into TGF-$\beta$âĂŞSmad signalling," *Trends in Biochemical Sciences*, vol. 29, no. 5, pp. 265–273, 2004.

[74] G. J. Inman, "SB-431542 Is a Potent and Specific Inhibitor of Transforming Growth Factor-beta Superfamily Type I Activin Receptor-Like Kinase (ALK) Receptors ALK4, ALK5, and ALK7," *Molecular Pharmacology*, vol. 62, pp. 65–74, July 2002.

[75] A. C. Daly, R. A. Randall, and C. S. Hill, "Transforming growth factor beta-induced Smad1/5 phosphorylation in epithelial cells is mediated by novel receptor complexes and is essential for anchorage-independent growth.," *Molecular and cellular biology*, vol. 28, pp. 6889–902, Nov. 2008.

[76] M.-J. Goumans, G. Valdimarsdottir, S. Itoh, F. Lebrin, J. Larsson, C. Mummery, S. Karlsson, and P. ten Dijke, "Activin Receptor-like Kinase (ALK)1 Is an Antagonistic Mediator of Lateral TGF$\beta$/ALK5 Signaling," *Molecular Cell*, vol. 12, pp. 817–828, Oct. 2003.

[77] I. M. Liu, S. H. Schilling, K. A. Knouse, L. Choy, R. Derynck, and X.-F. Wang, "TGFbeta-stimulated Smad1/5 phosphorylation requires the ALK5 L45 loop and mediates the pro-migratory TGFbeta switch.," *The EMBO journal*, vol. 28, pp. 88–98, Jan. 2009.

[78] K. H. Wrighton, X. Lin, P. B. Yu, and X.-H. Feng, "Transforming Growth Factor {beta} Can Stimulate Smad1 Phosphorylation Independently of Bone Morphogenic Protein Receptors.," *The Journal of biological chemistry*, vol. 284, pp. 9755–63, Apr. 2009.

[79] M. Tojo, Y. Hamashima, A. Hanyu, T. Kajimoto, M. Saitoh, K. Miyazono, M. Node, and T. Imamura, "The ALK-5 inhibitor A-83-01 inhibits Smad signaling and epithelial-to-mesenchymal transition by transforming growth factor-beta.," *Cancer science*, vol. 96, pp. 791–800, Nov. 2005.

[80] J. Vogt, R. Traynor, and G. P. Sapkota, "The specificities of small molecule inhibitors of the TGFßand BMP pathways," *Cellular Signalling*, vol. 23, no. 11, pp. 1831–1842, 2011.

[81] A. A. Ogunjimi, E. Zeqiraj, D. F. Ceccarelli, F. Sicheri, J. L. Wrana, and L. David, "Structural basis for specificity of TGF$\beta$ family receptor small molecule inhibitors.," *Cellular signalling*, vol. 24, pp. 476–83, Feb. 2012.

[82] K. C. Kirkbride, T. A. Townsend, M. W. Bruinsma, J. V. Barnett, and G. C. Blobe, "Bone morphogenetic proteins signal through the transforming growth factor-beta type III receptor.," *The Journal of biological chemistry*, vol. 283, pp. 7628–37, Mar. 2008.

[83] M. Bilandzic and K. L. Stenvers, "Betaglycan: a multifunctional accessory.," *Molecular and cellular endocrinology*, vol. 339, pp. 180–9, June 2011.

[84] S. Kaname and E. Ruoslahti, "Betaglycan has multiple binding sites for transforming growth factor-beta 1.," *The Biochemical journal*, vol. 315 ( Pt 3, pp. 815–20, May 1996.

[85] V. Mendoza, M. M. Vilchis-Landeros, G. Mendoza-Hernández, T. Huang, M. M. Villarreal, A. P. Hinck, F. López-Casillas, and J.-L. Montiel, "Betaglycan has two independent domains required for high affinity TGF-beta binding: proteolytic cleavage separates the domains and inactivates the neutralizing activity of the soluble receptor.," *Biochemistry*, vol. 48, pp. 11755–65, Dec. 2009.

[86] P. Lastres, A. Letamendía, H. Zhang, C. Rius, N. Almendro, U. Raab, L. A. López, C. Langa, A. Fabra, M. Letarte, and C. Bernabéu, "Endoglin modulates cellular responses to TGF-beta 1.," *The Journal of cell biology*, vol. 133, pp. 1109–21, June 1996.

[87] A. Letamendía, P. Lastres, L. M. Botella, U. Raab, C. Langa, B. Velasco, L. Attisano, and C. Bernabeu, "Role of endoglin in cellular responses to transforming growth factor-beta. A comparative study with betaglycan.," *The Journal of biological chemistry*, vol. 273, pp. 33011–9, Dec. 1998.

[88] C. Bernabeu, J. M. Lopez-Novoa, and M. Quintanilla, "The emerging role of TGF-beta superfamily coreceptors in cancer.," *Biochimica et biophysica acta*, vol. 1792, pp. 954–73, Oct. 2009.

[89] L. C. Fuentealba, E. Eivers, A. Ikeda, C. Hurtado, H. Kuroda, E. M. Pera, and E. M. De Robertis, "Integrating patterning signals: Wnt/GSK3 regulates the duration of the BMP/Smad1 signal.," *Cell*, vol. 131, pp. 980–93, Nov. 2007.

[90] X. Guo, A. Ramirez, D. S. Waddell, Z. Li, X. Liu, and X.-F. Wang, "Axin and GSK3-control Smad3 protein stability and modulate TGF- signaling.," *Genes & development*, vol. 22, pp. 106–20, Jan. 2008.

[91] A. Zieba, K. Pardali, O. Söderberg, L. Lindbom, E. Nyström, A. Moustakas, C.-H. Heldin, and U. Landegren, "Intercellular variation in signaling through the TGF-$\beta$ pathway and its relation to cell density and cell cycle phase.," *Molecular & cellular proteomics*, vol. 11, p. M111.013482, July 2012.

[92] F. J. Nicolás, K. De Bosscher, B. Schmierer, and C. S. Hill, "Analysis of Smad nucleocytoplasmic shuttling in living cells.," *Journal of Cell Science*, vol. 117, pp. 4113–25, Aug. 2004.

[93] C. E. Pierreux, F. J. Nicolás, and C. S. Hill, "Transforming growth factor beta-independent shuttling of Smad4 between the cytoplasm and nucleus.," *Molecular and cellular biology*, vol. 20, pp. 9041–54, Dec. 2000.

[94] M. Watanabe, N. Masuyama, M. Fukuda, and E. Nishida, "Regulation of intracellular dynamics of Smad4 by its leucine-rich nuclear export signal.," *EMBO reports*, vol. 1, pp. 176–82, Aug. 2000.

[95] L. Xu, Y. Kang, S. Çöl, and J. Massagué, "Smad2 Nucleocytoplasmic Shuttling by Nucleoporins CAN/Nup214 and Nup153 Feeds TGF$\beta$ Signaling Complexes in the Cytoplasm and Nucleus," *Molecular Cell*, vol. 10, pp. 271–282, Aug. 2002.

[96] A. Warmflash, Q. Zhang, B. Sorre, A. Vonica, E. D. Siggia, and A. H. Brivanlou, "Dynamics of TGF-$\beta$ signaling reveal adaptive and pulsatile behaviors reflected in the nuclear localization of transcription factor Smad4.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 109, pp. E1947–56, July 2012.

[97] D. C. Clarke, M. D. Betterton, and X. Liu, "Systems theory of Smad signalling.," *Systems biology*, vol. 153, pp. 412–24, Nov. 2006.

[98] G. J. Inman, F. J. Nicolás, and C. S. Hill, "Nucleocytoplasmic Shuttling of Smads 2, 3, and 4 Permits Sensing of TGF-$\beta$ Receptor Activity," *Molecular Cell*, vol. 10, no. 2, pp. 283–294, 2002.

[99] P. Lönn, A. Morén, E. Raja, M. Dahl, and A. Moustakas, "Regulating the stability of TGFbeta receptors and Smads.," *Cell research*, vol. 19, pp. 21–35, Jan. 2009.

[100] A. F. Candia, T. Watabe, S. H. Hawley, D. Onichtchouk, Y. Zhang, R. Derynck, C. Niehrs, and K. W. Cho, "Cellular interpretation of multiple TGF-beta signals: intracellular antagonism between activin/BVg1 and BMP-2/4 signaling mediated by Smads.," *Development*, vol. 124, pp. 4467–80, Nov. 1997.

[101] S. Malhotra and P. Kincade, "Wnt-related molecules and signaling pathway equilibrium in hematopoiesis," *Cell Stem Cell*, vol. 4, no. 1, pp. 27–36, 2009.

[102] A. J. Mikels and R. Nusse, "Wnts as ligands: processing, secretion and reception.," *Oncogene*, vol. 25, pp. 7461–8, Dec. 2006.

[103] H. Clevers, "Wnt/beta-catenin signaling in development and disease.," *Cell*, vol. 127, pp. 469–80, Nov. 2006.

[104] J. O. Mason, J. Kitajewski, and H. E. Varmus, "Mutational analysis of mouse Wnt-1 identifies two temperature-sensitive alleles and attributes of Wnt-1 protein essential for transformation of a mammary cell line.," *Molecular biology of the cell*, vol. 3, pp. 521–33, May 1992.

[105] E. Verheyen, "Wnts as Self-Renewal Factors: Mammary Stem Cells and Beyond," *Cell stem cell*, vol. 6, pp. 494–495, 2010.

[106] F. Reichsman, L. Smith, and S. Cumberledge, "Glycosaminoglycans can modulate extracellular localization of the wingless protein and promote signal transduction.," *The Journal of cell biology*, vol. 135, pp. 819–27, Nov. 1996.

[107] K. Willert, J. D. Brown, E. Danenberg, A. W. Duncan, I. L. Weissman, T. Reya, J. R. Yates, and R. Nusse, "Wnt proteins are lipid-modified and can act as stem cell growth factors.," *Nature*, vol. 423, pp. 448–52, May 2003.

[108] C. Y. Janda, D. Waghray, A. M. Levin, C. Thomas, and K. C. Garcia, "Structural basis of Wnt recognition by Frizzled.," *Science (New York, N.Y.)*, vol. 337, pp. 59–64, July 2012.

[109] S. Angers and R. T. Moon, "Proximal events in Wnt signal transduction.," *Nature Reviews Molecular Cell Biology*, vol. 10, pp. 468–77, July 2009.

[110] B. T. MacDonald, K. Tamai, and X. He, "Wnt/beta-catenin signaling: components, mechanisms, and diseases.," *Developmental Cell*, vol. 17, pp. 9–26, July 2009.

[111] V. F. Taelman, R. Dobrowolski, J.-L. Plouhinec, L. C. Fuentealba, P. P. Vorwald, I. Gumper, D. D. Sabatini, and E. M. De Robertis, "Wnt signaling requires sequestration of glycogen synthase kinase 3 inside multivesicular endosomes.," *Cell*, vol. 143, pp. 1136–48, Dec. 2010.

[112] S. Chen, D. Bubeck, B. T. Macdonald, W.-X. Liang, J.-H. Mao, T. Malinauskas, O. Llorca, a. R. Aricescu, C. Siebold, X. He, and E. Y. Jones, "Structural and Functional Studies of LRP6 Ectodomain Reveal a Platform for Wnt Signaling," *Developmental cell*, vol. 3, pp. 848–861, Oct. 2011.

[113] V. E. Ahn, M. L.-H. Chu, H.-J. Choi, D. Tran, A. Abo, and W. I. Weis, "Structural Basis of Wnt Signaling Inhibition by Dickkopf Binding to LRP5/6.," *Developmental cell*, vol. 21, pp. 862–873, Oct. 2011.

[114] Z. Cheng, T. Biechele, Z. Wei, S. Morrone, R. T. Moon, L. Wang, and W. Xu, "Crystal structures of the extracellular domain of LRP6 and its complex with DKK1.," *Nature structural & molecular biology*, vol. 18, pp. 1204–10, Nov. 2011.

[115] K.-A. Kim, M. Kakitani, J. Zhao, T. Oshima, T. Tang, M. Binnerts, Y. Liu, B. Boyle, E. Park, P. Emtage, W. D. Funk, and K. Tomizuka, "Mitogenic influence of human R-spondin1 on the intestinal epithelium.," *Science*, vol. 309, pp. 1256–9, Aug. 2005.

[116] A. Glinka, C. Dolde, N. Kirsch, Y.-L. Huang, O. Kazanskaya, D. Ingelfinger, M. Boutros, C.-M. Cruciat, and C. Niehrs, "LGR4 and LGR5 are R-spondin receptors mediating Wnt/$\beta$-catenin and Wnt/PCP signalling.," *EMBO Reports*, vol. 12, pp. 1055–1061, Sept. 2011.

[117] W. de Lau, N. Barker, T. Y. Low, B.-K. Koo, V. S. W. Li, H. Teunissen, P. Kujala, A. Haegebarth, P. J. Peters, M. van de Wetering, D. E. Stange, J. van Es, D. Guardavaccaro, R. B. M. Schasfoort, Y. Mohri, K. Nishimori, S. Mohammed, A. J. R. Heck, and H. Clevers, "Lgr5 homologues associate with Wnt receptors and mediate R-spondin signalling.," *Nature*, vol. 476, pp. 293–297, July 2011.

[118] R. van Amerongen, "Alternative Wnt pathways and receptors.," *Cold Spring Harbor perspectives in biology*, vol. 4, Jan. 2012.

[119] H. Huang and P. Klein, "Interactions between BMP and Wnt signaling pathways in mammalian cancers," *Cancer Biology & Therapy*, vol. 3, no. 7, pp. 676–678, 2004.

[120] C. Jamieson, M. Sharma, and B. R. Henderson, "Wnt signaling from membrane to nucleus: $\beta$-Catenin caught in a loop.," *The international journal of biochemistry & cell biology*, Mar. 2012.

[121] M. van de Wetering, M. Oosterwegel, D. Dooijes, and H. Clevers, "Identification and cloning of TCF-1, a T lymphocyte-specific transcription factor containing a sequence-specific HMG box.," *The EMBO journal*, vol. 10, pp. 123–32, Jan. 1991.

[122] M. Molenaar, M. van de Wetering, M. Ooster-wegel, J. Peterson-Maduro, S. Godsave, V. Korinek, J. Roose, O. Destree, and H. Clevers, "XTcf-3 Transcription Factor Mediates Beta-Catenin-Induced Axis Formation in Xenopus Embryos," *Cell*, vol. 86, pp. 391–399, Aug. 1996.

[123] A. J. Hanson, H. a. Wallace, T. J. Freeman, R. D. Beauchamp, L. a. Lee, and E. Lee, "XIAP Monoubiquitylates Groucho/TLE to Promote Canonical Wnt Signaling.," *Molecular cell*, pp. 1–10, Feb. 2012.

[124] P. Hatzis, L. G. van der Flier, M. a. van Driel, V. Guryev, F. Nielsen, S. Denissov, I. J. Nijman, J. Koster, E. E. Santo, W. Welboren, R. Versteeg, E. Cuppen, M. van de Wetering, H. Clevers, and H. G. Stunnenberg, "Genome-wide pattern of TCF7L2/TCF4 chromatin occupancy in colorectal cancer cells.," *Molecular and Cellular Biology*, vol. 28, pp. 2732–44, Apr. 2008.

[125] S. Frame and P. Cohen, "GSK3 takes centre stage more than 20 years after its discovery.," *The Biochemical Journal*, vol. 359, pp. 1–16, Oct. 2001.

[126] V. Stambolic, L. Ruel, and J. R. Woodgett, "Lithium inhibits glycogen synthase kinase-3 activity and mimics wingless signalling in intact cells.," *Current biology : CB*, vol. 6, pp. 1664–8, Dec. 1996.

[127] C. M. Hedgepeth, L. J. Conrad, J. Zhang, H. C. Huang, V. M. Lee, and P. S. Klein, "Activation of the Wnt signaling pathway: a molecular mechanism for lithium action.," *Developmental biology*, vol. 185, pp. 82–91, May 1997.

[128] L. Meijer, A.-l. Skaltsounis, P. Magiatis, P. Polychronopoulos, M. Knockaert, M. Leost, X. Ryan, C. Vonica, A. Brivanlou, R. Dajani, and Others, "GSK-3-selective inhibitors derived from Tyrian purple indirubins," *Chemistry & Biology*, vol. 10, no. 12, pp. 1255–1266, 2003.

[129] D. M. Roberts, M. I. Pronobis, J. S. Poulton, E. G. Kane, and M. Peifer, "Regulation of Wnt signaling by the tumor suppressor APC does not require ability to enter the nucleus nor a particular cytoplasmic localization.," *Molecular biology of the cell*, pp. mbc.E11–11–0965–, Apr. 2012.

[130] K. E. Spink, P. Polakis, and W. I. Weis, "Structural basis of the Axin Âś adenomatous polyposis coli interaction," *EMBO Journal*, vol. 19, no. 10, pp. 2270–2279, 2000.

[131] E. Lee, A. Salic, R. Krüger, R. Heinrich, and M. W. Kirschner, "The roles of APC and Axin derived from experimental and theoretical analysis of the Wnt pathway.," *PLoS biology*, vol. 1, p. E10, Oct. 2003.

[132] C. Gao and Y.-G. Chen, "Dishevelled: The hub of Wnt signaling.," *Cellular signalling*, vol. 22, pp. 717–27, May 2010.

[133] S. Acebron, E. Karaulanov, B. Berger, Y.-L. Huang, and C. Niehrs, "Mitotic Wnt Signaling Promotes Protein Stabilization and Regulates Cell Size," *Molecular Cell*, May 2014.

[134] A. De, "Wnt/Ca2+ signaling pathway: a brief overview.," *Acta biochimica et biophysica Sinica*, vol. 43, pp. 745–56, Oct. 2011.

[135] N. Barker, S. Bartfeld, and H. Clevers, "Tissue-resident adult stem cell populations of rapidly self-renewing organs.," *Cell Stem Cell*, vol. 7, pp. 656–70, Dec. 2010.

[136] K. Takahashi and S. Yamanaka, "Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors.," *Cell*, vol. 126, pp. 663–76, Aug. 2006.

[137] L. Fischer, G. Boland, and R. S. Tuan, "Wnt-3A enhances bone morphogenetic protein-2-mediated chondrogenesis of murine C3H10T1/2 mesenchymal cells.," *The Journal of biological chemistry*, vol. 277, pp. 30870–8, Aug. 2002.

[138] A. Nakashima, T. Katagiri, and M. Tamura, "Cross-talk between Wnt and bone morphogenetic protein 2 (BMP-2) signaling in differentiation pathway of C2C12 myoblasts.," *The Journal of biological chemistry*, vol. 280, pp. 37660–8, Nov. 2005.

[139] N. Itasaki and S. Hoppler, "Crosstalk between Wnt and bone morphogenic protein signaling: a turbulent relationship.," *Developmental dynamics : an official publication of the American Association of Anatomists*, vol. 239, pp. 16–33, Jan. 2010.

[140] K. Feigenson, M. Reid, J. See, E. B. Crenshaw III, and J. B. Grinspan, "Canonical Wnt signalling requires the BMP pathway to inhibit oligodendrocyte maturation.," *ASN neuro*, vol. 3, pp. 147–158, Jan. 2011.

[141] A. D. Lander, J. Kimble, H. Clevers, E. Fuchs, D. Montarras, M. Buckingham, A. L. Calof, A. Trumpp, and T. Oskarsson, "What does the concept of the stem cell niche really mean today?," *BMC Biology*, vol. 10, p. 19, Jan. 2012.

[142] K. C. Davidson, a. M. Adams, J. M. Goodson, C. E. McDonald, J. C. Potter, J. D. Berndt, T. L. Biechele, R. J. Taylor, and R. T. Moon, "Wnt/Âǎ-catenin signaling promotes differentiation, not self-renewal, of human embryonic stem cells and is repressed by Oct4," *Proceedings of the National Academy of Sciences*, pp. 1–6, Mar. 2012.

[143] L. Libusova, M. P. Stemmler, A. Hierholzer, H. Schwarz, and R. Kemler, "N-cadherin can structurally substitute for E-cadherin during intestinal development but leads to polyp formation.," *Development (Cambridge, England)*, vol. 137, pp. 2297–305, July 2010.

[144] A. Gregorieff, D. Pinto, H. Begthel, O. Destrée, M. Kielman, and H. Clevers, "Expression Pattern of Wnt Signaling Components in the Adult Intestine," *Gastroenterology*, vol. 129, pp. 626–638, Aug. 2005.

[145] M. van de Wetering, E. Sancho, C. Verweij, W. de Lau, I. Oving, A. Hurlstone, K. van der Horn, E. Batlle, D. Coudreuse, A. P. Haramis, M. Tjon-Pon-Fong, P. Moerer, M. van den Born, G. Soete, S. Pals, M. Eilers, R. Medema, and H. Clevers, "The beta-catenin/TCF-4 complex imposes a crypt progenitor phenotype on colorectal cancer cells.," *Cell*, vol. 111, pp. 241–50, Oct. 2002.

[146] A. Schepers and H. Clevers, "Wnt signaling, stem cells, and cancer of the gastrointestinal tract.," *Cold Spring Harbor Perspectives in Biology*, vol. 4, Jan. 2012.

[147] R. Teo, F. Möhrlen, G. Plickert, W. a. Müller, and U. Frank, "An evolutionary conserved role of Wnt signaling in stem cell fate decision.," *Developmental Biology*, vol. 289, pp. 91–9, Jan. 2006.

[148] 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, and H. Clevers, "Identification of stem cells in small intestine and colon by marker gene Lgr5.," *Nature*, vol. 449, pp. 1003–7, Oct. 2007.

[149] T. Sato, D. E. Stange, M. Ferrante, R. G. J. Vries, J. H. van Es, S. van den Brink, W. J. van Houdt, A. Pronk, J. van Gorp, P. D. Siersema, and H. Clevers, "Long-Term Expansion of Epithelial Organoids from Human Colon, Adenoma, Adenocarcinoma, and Barrett's Epithelium.," *Gastroenterology*, vol. 141, pp. 1762–1772, Aug. 2011.

[150] T. Sato, J. H. van Es, H. J. Snippert, D. E. Stange, R. G. Vries, M. van den Born, N. Barker, N. F. Shroyer, M. van de Wetering, and H. Clevers, "Paneth cells constitute the niche for Lgr5 stem cells in intestinal crypts.," *Nature*, vol. 469, pp. 415–8, Jan. 2011.

[151] S. Yui, T. Nakamura, T. Sato, Y. Nemoto, T. Mizutani, X. Zheng, S. Ichinose, T. Nagaishi, R. Okamoto, K. Tsuchiya, H. Clevers, and M. Watanabe, "Functional engraftment of colon epithelium expanded in vitro from a single adult Lgr5+ stem cell," *Nature Medicine*, vol. epub ahead, Mar. 2012.

[152] S. S. Zeki, T. A. Graham, and N. A. Wright, "Stem cells and their implications for colorectal cancer.," *Nature Reviews. Gastroenterology & Hepatology*, vol. 8, pp. 90–100, Oct. 2011.

[153] N. Barker, R. a. Ridgway, J. H. van Es, M. van de Wetering, H. Begthel, M. van den Born, E. Danenberg, A. R. Clarke, O. J. Sansom, and H. Clevers, "Crypt stem cells as the cells-of-origin of intestinal cancer.," *Nature*, vol. 457, pp. 608–11, Jan. 2009.

[154] L. Attisano and E. Labbé, "TGFbeta and Wnt pathway cross-talk.," *Cancer metastasis reviews*, vol. 23, no. 1-2, pp. 53–61, 2004.

[155] E. Labbé, L. Lock, A. Letamendia, A. E. Gorska, R. Gryfe, S. Gallinger, H. L. Moses, and L. Attisano, "Transcriptional cooperation between the transforming growth factor-beta and Wnt pathways in mammary and intestinal tumorigenesis.," *Cancer research*, vol. 67, pp. 75–84, Jan. 2007.

[156] C. C. W. Chong, R. J. W. Stump, F. J. Lovicu, and J. W. McAvoy, "TGFbeta promotes Wnt expression during cataract development.," *Experimental eye research*, vol. 88, pp. 307–13, Feb. 2009.

[157] P. Minoo and C. Li, "Cross-talk between transforming growth factor-beta and Wingless/Int pathways in lung development and disease.," *The international journal of biochemistry & cell biology*, vol. 42, pp. 809–12, June 2010.

[158] R. Serra, S. L. Easter, W. Jiang, and S. E. Baxley, "Wnt5a as an effector of TGFβ in mammary development and cancer.," *Journal of mammary gland biology and neoplasia*, vol. 16, pp. 157–67, June 2011.

[159] E. Rodríguez-Carballo, A. Ulsamer, A. R. G. Susperregui, C. Manzanares-Céspedes, E. Sánchez-García, R. Bartrons, J. L. Rosa, and F. Ventura, "Conserved regulatory motifs in osteogenic gene promoters integrate cooperative effects of canonical Wnt and BMP pathways.," *Journal of bone and mineral research*, vol. 26, pp. 718–29, Apr. 2011.

[160] G. Chen, C. Deng, and Y.-P. Li, "TGF-β and BMP signaling in osteoblast differentiation and bone formation.," *International journal of biological sciences*, vol. 8, pp. 272–88, Jan. 2012.

[161] A. Akhmetshina, K. Palumbo, C. Dees, C. Bergmann, P. Venalis, P. Zerr, A. Horn, T. Kireva, C. Beyer, J. Zwerina, H. Schneider, A. Sadowski, M.-O. Riener, O. a. MacDougald, O. Distler, G. Schett, and J. H. W. Distler, "Activation of canonical Wnt signalling is required for TGF-β-mediated fibrosis.," *Nature communications*, vol. 3, p. 735, Jan. 2012.

[162] J. Song, J. McColl, E. Camp, N. Kennerley, G. F. Mok, D. McCormick, T. Grocott, G. N. Wheeler, and A. E. Münsterberg, "Smad1 transcription factor integrates BMP2 and Wnt3a signals in migrating cardiac progenitor cells.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 111, pp. 7337–42, May 2014.

[163] B. Iacopetta, "Are there two sides to colorectal cancer?," *International Journal of Cancer*, vol. 101, pp. 403–8, Oct. 2002.

[164] K. Yang, N. V. Popova, W. C. Yang, I. Lozonschi, S. Tadesse, S. Kent, L. Bancroft, I. Matise, R. T. Cormier, S. J. Scherer, W. Edelmann, M. Lipkin, L. Augenlicht, and A. Velcich, "Interaction of Muc2 and Apc on Wnt signaling and in intestinal tumorigenesis: potential role of chronic inflammation.," *Cancer Research*, vol. 68, pp. 7313–22, Sept. 2008.

[165] A. Jemal, R. Siegel, J. Xu, and E. Ward, "Cancer statistics, 2010.," *CA: A Cancer Journal for Clinicians*, vol. 60, no. 5, pp. 277–300, 2010.

[166] P. Simon-Assmann, C. Spenle, O. Lefebvre, and M. Kedinger, "The role of the basement membrane as a modulator of intestinal epithelial-mesenchymal interactions.," *Progress in Molecular Biology and Translational Science*, vol. 96, pp. 175–206, Jan. 2010.

[167] C. Kosinski, D. Stange, C. Xu, A. Chan, C. Ho, S. Yuen, R. Mifflin, D. Powell, H. Clevers, S. Leung, and Others, "Indian hedgehog regulates intestinal stem cell fate through epithelial-mesenchymal interactions during development," *Gastroenterology*, vol. 139, no. 3, pp. 893–903, 2010.

[168] P. W. Ingham, Y. Nakano, and C. Seger, "Mechanisms and functions of Hedgehog signalling across the metazoa.," *Nature reviews. Genetics*, vol. 12, pp. 393–406, June 2011.

[169] N. Lahar, N. Y. Lei, J. Wang, Z. Jabaji, S. C. Tung, V. Joshi, M. Lewis, M. Stelzner, M. G. Martín, and J. C. Y. Dunn, "Intestinal subepithelial myofibroblasts support in vitro and in vivo growth of human small intestinal epithelium.," *PloS One*, vol. 6, p. e26898, Jan. 2011.

[170] H. Snippert, L. Van Der Flier, T. Sato, J. Van Es, M. Van Den Born, C. Kroon-Veenboer, N. Barker, A. Klein, J. Van Rheenen, B. Simons, and H. Clevers, "Intestinal crypt homeostasis results from neutral competition between symmetrically dividing Lgr5 stem cells," *Cell*, vol. 143, pp. 134–144, 2010.

[171] G. T. Eisenhoffer, P. D. Loftus, M. Yoshigi, H. Otsuna, C.-B. Chien, P. A. Morcos, and J. Rosenblatt, "Crowding induces live cell extrusion to maintain homeostatic cell numbers in epithelia.," *Nature*, Apr. 2012.

[172] P. Buske, J. Galle, N. Barker, G. Aust, H. Clevers, and M. Loeffler, "A comprehensive model of the spatio-temporal stem cell and tissue organisation in the intestinal crypt.," *PLoS Computational Biology*, vol. 7, p. e1001045, Jan. 2011.

[173] T. Hinoi, A. Akyol, B. K. Theisen, D. O. Ferguson, J. K. Greenson, B. O. Williams, K. R. Cho, and E. R. Fearon, "Mouse model of colonic adenoma-carcinoma progression based on somatic Apc inactivation.," *Cancer Research*, vol. 67, pp. 9721–30, Oct. 2007.

[174] P. Rizk and N. Barker, "Gut stem cells in tissue renewal and disease: methods, markers, and myths.," *Wiley interdisciplinary reviews. Systems biology and medicine*, May 2012.

[175] J. C. Hardwick, G. R. Van Den Brink, S. A. Bleuming, I. Ballester, J. Van Den Brande,

J. J. Keller, G. A. Offerhaus, S. J. Van Deventer, and M. P. Peppelenbosch, "Bone morphogenetic protein 2 is expressed by, and acts upon, mature epithelial cells in the colon," *Gastroenterology*, vol. 126, pp. 111–121, Jan. 2004.

[176] Y. Yamada, H. Mashima, T. Sakai, T. Matsuhashi, M. Jin, and H. Ohnishi, "Functional roles of TGF-$\beta$1 in intestinal epithelial cells through Smad-dependent and non-Smad pathways.," *Digestive diseases and sciences*, vol. 58, pp. 1207–17, May 2013.

[177] J. C. Hardwick, L. L. Kodach, G. J. Offerhaus, and G. R. van den Brink, "Bone morphogenetic protein signalling in colorectal cancer.," *Nature reviews. Cancer*, vol. 8, pp. 806–12, Oct. 2008.

[178] P. de Santa Barbara, G. R. van den Brink, and D. J. Roberts, "Development and differentiation of the intestinal epithelium.," *Cellular and molecular life sciences*, vol. 60, pp. 1322–32, July 2003.

[179] E. Batlle, J. Henderson, H. Beghtel, M. van den Born, E. Sancho, G. Huls, J. Meeldijk, J. Robertson, M. van de Wetering, T. Pawson, and Others, "beta-catenin and TCF mediate cell positioning in the intestinal epithelium by controlling the expression of EphB/ephrinB," *Cell*, vol. 111, no. 2, pp. 251–263, 2002.

[180] C. Alexandre, A. Baena-Lopez, and J.-P. Vincent, "Patterning and growth control by membrane-tethered Wingless.," *Nature*, vol. 505, pp. 180–5, Jan. 2014.

[181] O. Serralbo and C. Marcelle, "Migrating cells mediate long-range WNT signaling.," *Development (Cambridge, England)*, vol. 141, pp. 2057–63, May 2014.

[182] S. Seshagiri, E. W. Stawiski, S. Durinck, Z. Modrusan, E. E. Storm, C. B. Conboy, S. Chaudhuri, Y. Guan, V. Janakiraman, B. S. Jaiswal, J. Guillory, C. Ha, G. J. P. Dijkgraaf, J. Stinson, F. Gnad, M. A. Huntley, J. D. Degenhardt, P. M. Haverty, R. Bourgon, W. Wang, H. Koeppen, R. Gentleman, T. K. Starr, Z. Zhang, D. A. Largaespada, T. D. Wu, and F. J. de Sauvage, "Recurrent R-spondin fusions in colon cancer," *Nature*, Aug. 2012.

[183] T.-C. He, A. B. Sparks, C. Rago, H. Hermeking, L. Zawal, L. T. da Costa, P. J. Morin, B. Vogelstein, and K. W. Kinzler, "Identification of c-MYC as a Target of the APC Pathway," *Science*, vol. 281, pp. 1509–1512, Sept. 1998.

[184] A. W. Burgess, M. C. Faux, M. J. Layton, and R. G. Ramsay, "Wnt signaling and colon tumorigenesis–a view from the periphery.," *Experimental cell research*, vol. 317, pp. 2748–58, Nov. 2011.

[185] O. Sansom, K. Reed, and A. Hayes, "Loss of Apc in vivo immediately perturbs Wnt signaling, differentiation, and migration," *Genes & . . .*, pp. 1385–1390, 2004.

[186] O. J. Sansom, V. S. Meniel, V. Muncan, T. J. Phesse, J. a. Wilkins, K. R. Reed, J. K. Vass, D. Athineos, H. Clevers, and A. R. Clarke, "Myc deletion rescues Apc deficiency in the small intestine.," *Nature*, vol. 446, pp. 676–9, Apr. 2007.

[187] L. Levy and C. S. Hill, "Smad4 dependency defines two classes of transforming growth factor {beta} (TGF-{beta}) target genes and distinguishes TGF-{beta}-induced epithelial-mesenchymal transition from its antiproliferative and migratory responses.," *Molecular and cellular biology*, vol. 25, pp. 8108–25, Sept. 2005.

[188] A.-P. G. Haramis, H. Begthel, M. van den Born, J. van Es, S. Jonkheer, G. Offerhaus, and H. Clevers, "De novo crypt formation and juvenile polyposis on BMP inhibition in mouse intestine.," *Science*, vol. 303, pp. 1684–6, Mar. 2004.

[189] T. J. Freeman, J. J. Smith, X. Chen, M. K. Washington, J. T. Roland, A. L. Means, S. a. Eschrich, T. J. Yeatman, N. G. Deane, and R. D. Beauchamp, "Smad4-mediated signaling inhibits intestinal neoplasia by inhibiting expression of $\beta$-catenin.," *Gastroenterology*, vol. 142, pp. 562–571.e2, Mar. 2012.

[190] D. R. Warner, M. Pisano, E. A. Roberts, and R. M. Greene, "Identification of three novel Smad binding proteins involved in cell polarity," *FEBS Letters*, vol. 539, pp. 167–173, Mar. 2003.

[191] D. R. Warner, R. M. Greene, and M. M. Pisano, "Interaction between Smad 3 and Dishevelled in murine embryonic craniofacial mesenchymal cells.," *Orthodontics & craniofacial research*, vol. 8, pp. 123–30, May 2005.

[192] Z. Liu, Y. Tang, T. Qiu, X. Cao, and T. L. Clemens, "A dishevelled-1/Smad1 interaction couples WNT and bone morphogenetic protein

signaling pathways in uncommitted bone marrow stromal cells.," *The Journal of biological chemistry*, vol. 281, pp. 17156–63, June 2006.

[193] A. Mamidi, M. Inui, A. Manfrin, S. Soligo, E. Enzo, M. Aragona, M. Cordenonsi, O. Wessely, S. Dupont, and S. Piccolo, "Signaling crosstalk between TGFβ and Dishevelled/Par1b," *Cell death and differentiation*, vol. 19, pp. 1689–97, Oct. 2012.

[194] M. Furuhashi and K. Yagi, "Axin facilitates Smad3 activation in the transforming growth factor signaling pathway," *Molecular and Cellular Biology*, vol. 21, no. 15, pp. 5132–5141, 2001.

[195] W. Liu, H. Rui, J. Wang, S. Lin, Y. He, M. Chen, Q. Li, Z. Ye, S. Zhang, S. C. Chan, Y.-G. Chen, J. Han, and S.-C. Lin, "Axin is a scaffold protein in TGF-beta signaling that promotes degradation of Smad7 by Arkadia.," *The EMBO journal*, vol. 25, pp. 1646–58, Apr. 2006.

[196] Y. A. Zeng, M. Rahnama, S. Wang, W. Lee, and E. M. Verheyen, "Inhibition of Drosophila Wg signaling involves competition between Mad and Armadillo/beta-catenin for dTcf binding.," *PloS one*, vol. 3, p. e3893, Jan. 2008.

[197] H. Miyoshi, R. Ajima, C. T. Luo, T. P. Yamaguchi, and T. S. Stappenbeck, "Wnt5a potentiates TGF-β signaling to promote colonic crypt regeneration after tissue injury.," *Science (New York, N.Y.)*, vol. 338, pp. 108–13, Oct. 2012.

[198] S. Edlund, S. Lee, and S. Grimsby, "Interaction between Smad7 and B-catenin: importance for transforming growth factor B-induced apoptosis," *Molecular and Cellular Biology*, vol. 25, no. 4, pp. 1475–1488, 2005.

[199] D. R. Warner, R. M. Greene, and M. M. Pisano, "Cross-talk between the TGFbeta and Wnt signaling pathways in murine embryonic maxillary mesenchymal cells.," *FEBS letters*, vol. 579, pp. 3539–46, July 2005.

[200] S. Lei, A. Dubeykovskiy, A. Chakladar, L. Wojtukiewicz, and T. C. Wang, "The murine gastrin promoter is synergistically activated by transforming growth factor-beta/Smad and Wnt signaling pathways.," *The Journal of biological chemistry*, vol. 279, pp. 42492–502, Oct. 2004.

[201] D. R. Warner, P. Mukhopadhyay, G. N. Brock, V. Pihur, M. M. Pisano, and R. M. Greene, "TGFβ-1 and Wnt-3a interact to induce unique gene expression profiles in murine embryonic palate mesenchymal cells.," *Reproductive toxicology (Elmsford, N.Y.)*, vol. 31, pp. 128–33, Feb. 2011.

[202] S. M. Hussein, E. K. Duff, and C. Sirard, "Smad4 and beta-catenin co-activators functionally interact with lymphoid-enhancing factor to regulate graded expression of Msx2.," *The Journal of biological chemistry*, vol. 278, pp. 48805–14, Dec. 2003.

[203] B. Zhou, Y. Liu, M. Kahn, D. K. Ann, A. Han, H. Wang, C. Nguyen, P. Flodby, Q. Zhong, M. S. Krishnaveni, J. M. Liebler, P. Minoo, E. D. Crandall, and Z. Borok, "Interactions between β-catenin and transforming growth factor-β signaling pathways mediate epithelial-mesenchymal transition and are dependent on the transcriptional co-activator cAMP-response element-binding protein (CREB)-binding protein (CBP).," *The Journal of biological chemistry*, vol. 287, pp. 7026–38, Mar. 2012.

[204] E. Labbé, A. Letamendia, and L. Attisano, "Association of Smads with lymphoid enhancer binding factor 1/T cell-specific factor mediates cooperative signaling by the transforming growth factor-beta and wnt pathways.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 97, pp. 8358–63, July 2000.

[205] M. C. Hu and N. D. Rosenblum, "Smad1, beta-catenin and Tcf4 associate in a molecular complex with the Myc promoter in dysplastic renal tissue and cooperate to control Myc transcription.," *Development (Cambridge, England)*, vol. 132, pp. 215–25, Jan. 2005.

[206] S. Daopin, K. A. Piez, Y. Ogawa, and D. R. Davies, "Crystal structure of transforming growth factor-beta 2: an unusual fold for the superfamily.," *Science (New York, N.Y.)*, vol. 257, pp. 369–73, July 1992.

[207] S. Keller, J. Nickel, J.-L. Zhang, W. Sebald, and T. D. Mueller, "Molecular recognition of BMP-2 and BMP receptor IA.," *Nature structural & molecular biology*, vol. 11, pp. 481–8, May 2004.

[208] B. M. Chacko, B. Y. Qin, A. Tiwari, G. Shi, S. Lam, L. J. Hayward, M. De Caestecker,

and K. Lin, "Structural basis of heteromeric smad protein assembly in TGF-beta signaling.," *Molecular cell*, vol. 15, pp. 813–23, Sept. 2004.

[209] N. BabuRajendran, P. Palasingam, K. Narasimhan, W. Sun, S. Prabhakar, R. Jauch, and P. R. Kolatkar, "Structure of Smad1 MH1/DNA complex reveals distinctive rearrangements of BMP and TGF-beta effectors.," *Nucleic acids research*, vol. 38, pp. 3477–88, June 2010.

[210] J. Pei, M. Tang, and N. V. Grishin, "PROMALS3D web server for accurate multiple protein sequence and structure alignments.," *Nucleic acids research*, vol. 36, pp. W30–4, July 2008.

[211] E. Paradis, J. Claude, and K. Strimmer, "A{PE}: analyses of phylogenetics and evolution in {R} language," *Bioinformatics*, vol. 20, pp. 289–290, 2004.

[212] S. F. Altschul, T. L. Madden, A. A. Schäffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman, "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.," *Nucleic acids research*, vol. 25, pp. 3389–402, Sept. 1997.

[213] A. I. Roig, U. Eskiocak, S. K. Hight, S. B. Kim, O. Delgado, R. F. Souza, S. J. Spechler, W. E. Wright, and J. W. Shay, "Immortalized epithelial cells derived from human colon biopsies express stem cell markers and differentiate in vitro.," *Gastroenterology*, vol. 138, pp. 1012–21.e1–5, Mar. 2010.

[214] A. Quaroni, K. J. Isselbacher, and E. Ruoslahti, "Fibronectin synthesis by epithelial crypt cells of rat small intestine.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 75, pp. 5548–52, Nov. 1978.

[215] J. Fogh, W. C. Wright, and J. D. Loveless, "Absence of HeLa cell contamination in 169 cell lines derived from human tumors.," *Journal of the National Cancer Institute*, vol. 58, pp. 209–14, Feb. 1977.

[216] B. Neumann, M. Held, U. Liebel, and H. Erfle, "High-throughput RNAi screening by time-lapse imaging of live human cells," *Nature . . .*, vol. 3, no. 5, pp. 385–390, 2006.

[217] D. W. Young, A. Bender, J. Hoyt, E. McWhinnie, G.-W. Chirn, C. Y. Tao, J. a. Tallarico, M. Labow, J. L. Jenkins, T. J. Mitchison, and Y. Feng, "Integrating high-content screening and ligand-target prediction to identify mechanism of action.," *Nature chemical biology*, vol. 4, pp. 59–68, Jan. 2008.

[218] Y. Feng, T. J. Mitchison, A. Bender, D. W. Young, and J. a. Tallarico, "Multi-parameter phenotypic profiling: using cellular effects to characterize small-molecule compounds.," *Nature reviews. Drug discovery*, vol. 8, pp. 567–78, July 2009.

[219] D. Houle, D. R. Govindaraju, and S. Omholt, "Phenomics: the next challenge.," *Nature reviews. Genetics*, vol. 11, pp. 855–66, Dec. 2010.

[220] M. Held, M. H. a. Schmitz, B. Fischer, T. Walter, B. Neumann, M. H. Olma, M. Peter, J. Ellenberg, and D. W. Gerlich, "CellCognition: time-resolved phenotype annotation in high-throughput live cell imaging.," *Nature methods*, vol. 7, pp. 747–54, Sept. 2010.

[221] Y. Futamura, M. Kawatani, S. Kazami, K. Tanaka, M. Muroi, T. Shimizu, K. Tomita, N. Watanabe, and H. Osada, "Morphobase, an encyclopedic cell morphology database, and its use for drug target identification.," *Chemistry & biology*, vol. 19, pp. 1620–30, Dec. 2012.

[222] G. Danuser, "Computer Vision in Cell Biology," *Cell*, vol. 147, pp. 973–978, Nov. 2011.

[223] M. L. Schultz, L. E. Lipkin, M. J. Wade, P. F. Lemkin, and G. M. Carman, "High resolution shading correction.," *The journal of histochemistry and cytochemistry : official journal of the Histochemistry Society*, vol. 22, pp. 751–4, July 1974.

[224] V. Madisetti and D. B. Williams, eds., *The Digital Signal Processing: Handbook*. C R C Press LLC, 1998.

[225] B. Likar, J. B. A. Maintz, M. A. Viergever, and F. Perus, "Retrospective shading correction based on entropy minimization," *Journal of Microscopy*, vol. 197, pp. 285–295, Mar. 2000.

[226] M. A. Model and J. K. Burkhardt, "A standard for calibration and shading correction of a fluorescence microscope.," *Cytometry*, vol. 44, pp. 309–16, Aug. 2001.

[227] Y. Hiraoka, J. W. Sedat, and D. A. Agard, "The use of a charge-coupled device for quantitative optical microscopy of biological structures.," *Science (New York, N.Y.)*, vol. 238, pp. 36–41, Oct. 1987.

[228] R. D. Goldman and D. L. Spector, eds., *Live Cell Imaging: a laboratory manual.* Cold Spring Harbor, NY: Cold Spring Harbor Laboratory Press, 1 ed., 2005.

[229] S. Inoué and K. R. Spring, *Video Microscopy: The Fundamentals.* Springer, 1997.

[230] J. M. Murray, P. L. Appleton, J. R. Swedlow, and J. C. Waters, "Evaluating performance in three-dimensional fluorescence microscopy.," *Journal of microscopy*, vol. 228, pp. 390–405, Dec. 2007.

[231] R. Fiolka, Y. Belyaev, H. Ewers, and A. Stemmer, "Even illumination in total internal reflection fluorescence microscopy using laser light.," *Microscopy research and technique*, vol. 71, pp. 45–50, Jan. 2008.

[232] S. Herbert and R. Henriques, "Enhanced epifluorescence microscopy by uniform and intensity optimized illumination.," *Cytometry. Part A : the journal of the International Society for Analytical Cytology*, vol. 81, pp. 278–80, Apr. 2012.

[233] L. R. van den Doel, A. D. Klein, S. L. Ellenberger, H. Netten, F. R. Boddeke, L. J. van Vliet, and I. T. Young, "Quantitative evaluation of light microscopes based on image processing techniques," *Bioimaging*, vol. 6, pp. 138–149, Sept. 1998.

[234] M.-A. Bray and A. Carpenter, "Advanced Assay Development Guidelines for Image-Based High Content Screening and Analysis," in *Assay Guidance Manual* (G. Sittampalam, N. Gal-Edd, M. Arkin, D. Auld, C. Austin, B. Bejcek, M. Glicksman, J. Inglese, V. Lemmon, Z. Li, J. McGee, O. McManus, L. Minor, A. Napper, T. Riss, O. Trask, and J. Weidner, eds.), Bethesda (MD): Eli Lilly & Company and the National Center for Advancing Translational Sciences, 2013.

[235] A. D. Coster, C. Wichaidit, S. Rajaram, S. J. Altschuler, and L. F. Wu, "A simple image correction method for high-throughput microscopy," *Nature Methods*, vol. 11, pp. 602–602, May 2014.

[236] L. Hodgson, F. Shen, and K. Hahn, "Biosensors for characterizing the dynamics of rho family GTPases in living cells.," *Current protocols in cell biology / editorial board, Juan S. Bonifacino ... [et al.]*, vol. Chapter 14, pp. Unit 14.11.1–26, Mar. 2010.

[237] J. M. Zwier, G. J. Van Rooij, J. W. Hofstraat, and G. J. Brakenhoff, "Image calibration in fluorescence microscopy.," *Journal of microscopy*, vol. 216, pp. 15–24, Oct. 2004.

[238] D. E. Wolf, C. Samarasekera, and J. R. Swedlow, "Quantitative analysis of digital microscope images.," *Methods in cell biology*, vol. 81, pp. 365–96, Jan. 2007.

[239] J. C. Waters, "Accuracy and precision in quantitative fluorescence microscopy.," *The Journal of cell biology*, vol. 185, pp. 1135–48, July 2009.

[240] J. Lindblad and E. Bengtsson, "A comparison of methods for estimation of intensity nonuniformities in 2D and 3D microscope images of fluorescence stained cells," *In Proceedings of the 12th Scandinavian Conference on Image Analysis (SCIA), Bergen, Norway*, pp. 264–271, 2001.

[241] J. Lindblad, C. Wählby, E. Bengtsson, and A. Zaltsman, "Image analysis for automatic segmentation of cytoplasms and classification of Rac1 activation.," *Cytometry. Part A : the journal of the International Society for Analytical Cytology*, vol. 57, pp. 22–33, Jan. 2004.

[242] Z. Yin, A. Sadok, H. Sailem, A. McCarthy, X. Xia, F. Li, M. A. Garcia, L. Evans, A. R. Barr, N. Perrimon, C. J. Marshall, S. T. C. Wong, and C. Bakal, "A screen for morphological complexity identifies regulators of switch-like transitions between discreteÂăcell shapes," *Nature Cell Biology*, vol. 15, pp. 860–871, June 2013.

[243] C. A. Schneider, W. S. Rasband, and K. W. Eliceiri, "NIH Image to ImageJ: 25 years of image analysis," *Nature Methods*, vol. 9, pp. 671–675, June 2012.

[244] J. Schindelin, I. Arganda-Carreras, E. Frise, V. Kaynig, M. Longair, T. Pietzsch, S. Preibisch, C. Rueden, S. Saalfeld, B. Schmid, J.-Y. Tinevez, D. J. White, V. Hartenstein, K. Eliceiri, P. Tomancak, and A. Cardona, "Fiji: an open-source platform for biological-image analysis.," *Nature methods*, vol. 9, pp. 676–82, July 2012.

[245] M.-A. Bray, A. N. Fraser, T. P. Hasaka, and A. E. Carpenter, "Workflow and metrics for image quality control in large-scale high-content screens.," *Journal of biomolecular screening*, vol. 17, pp. 266–74, Feb. 2012.

[246] N. Otsu, "A threshold selection method from gray-level histograms," *Automatica*, vol. 9, no. 1, pp. 62–66, 1975.

[247] N. Malo, J. A. Hanley, S. Cerquozzi, J. Pelletier, and R. Nadon, "Statistical practice in high-throughput screening data analysis.," *Nature biotechnology*, vol. 24, pp. 167–75, Feb. 2006.

[248] P. Dragiev, R. Nadon, and V. Makarenkov, "Systematic error detection in experimental high-throughput screening.," *BMC bioinformatics*, vol. 12, p. 25, Jan. 2011.

[249] P. Dragiev, R. Nadon, and V. Makarenkov, "Two effective methods for correcting experimental high-throughput screening data.," *Bioinformatics (Oxford, England)*, vol. 28, pp. 1775–82, July 2012.

[250] J.-P. Carralot, A. Ogier, A. Boese, A. Genovesio, P. Brodin, P. Sommer, and T. Dorval, "A novel specific edge effect correction method for RNA interference screenings.," *Bioinformatics (Oxford, England)*, vol. 28, pp. 261–8, Jan. 2012.

[251] R. Zhong, M. S. Kim, M. A. White, Y. Xie, and G. Xiao, "SbacHTS: Spatial background noise correction for High-Throughput RNAi Screening.," *Bioinformatics (Oxford, England)*, vol. 29, pp. 2218–2220, July 2013.

[252] A. E. Carpenter, T. R. Jones, M. R. Lamprecht, C. Clarke, I. H. Kang, O. Friman, D. A. Guertin, J. H. Chang, R. A. Lindquist, J. Moffat, P. Golland, and D. M. Sabatini, "CellProfiler: image analysis software for identifying and quantifying cell phenotypes.," *Genome biology*, vol. 7, p. R100, Jan. 2006.

[253] L. Vincent and P. Soille, "Watersheds in digital spaces: an efficient algorithm based on immersion simulations," *IEEE transactions on pattern analysis and . . .*, 1991.

[254] N. Malpica, C. O. de Solórzano, J. J. Vaquero, a. Santos, I. Vallcorba, J. M. García-Sagredo, and F. del Pozo, "Applying watershed algorithms to the segmentation of clustered nuclei.," *Cytometry*, vol. 28, pp. 289–97, Aug. 1997.

[255] C. Ku, Y. Wang, B. Pavie, S. Altschuler, and L. Wu, "On identifying information from image-based spatial polarity phenotypes in neutrophils," in *Biomedical Imaging: From Nano to Macro, 2010 IEEE International Symposium on*, no. Figure 1, pp. 1029–1032, IEEE, 2010.

[256] J. B. Pawley, ed., *Handbook Of Biological Confocal Microscopy.* Springer, 3 ed., 2006.

[257] S. Rajaram, B. Pavie, L. F. Wu, and S. J. Altschuler, "PhenoRipper: software for rapidly profiling microscopy images.," *Nature methods*, vol. 9, pp. 635–7, July 2012.

[258] P. N. Dean and J. H. Jett, "Mathematical analysis of DNA distributions derived from flow microfluorometry.," *The Journal of cell biology*, vol. 60, pp. 523–7, Feb. 1974.

[259] M. H. Fox, "A model for the computer analysis of synchronous DNA distributions obtained by flow cytometry.," *Cytometry*, vol. 1, pp. 71–7, July 1980.

[260] G. Schatz, "The faces of Big Science," *Nature Reviews Molecular Cell Biology*, vol. 15, pp. 423–426, May 2014.

[261] B. Alberts, M. W. Kirschner, S. Tilghman, and H. Varmus, "Rescuing US biomedical research from its systemic flaws.," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 111, pp. 5773–7, Apr. 2014.

[262] J. P. A. Ioannidis, "Why most published research findings are false.," *PLoS medicine*, vol. 2, p. e124, Aug. 2005.

[263] J. Jarvik, S. Adler, C. Telmer, V. Subramaniam, and A. Lopez, "CD-tagging: a new approach to gene and protein discovery and analysis," *BioTechniques*, vol. 20, no. 5, pp. 896–904, 1996.

[264] J. W. Jarvik, G. W. Fisher, C. Shi, L. Hennen, C. Hauser, S. Adler, and P. B. Berget, "In vivo functional proteomics: mammalian genome annotation using CD-tagging.," *BioTechniques*, vol. 33, pp. 852–4, 856, 858–60 passim, Oct. 2002.

[265] A. Sigal, R. Milo, A. Cohen, N. Geva-Zatorsky, Y. Klein, I. Alaluf, N. Swerdlin, N. Perzov, T. Danon, Y. Liron, T. Raveh, A. E. Carpenter, G. Lahav, and U. Alon, "Dynamic proteomics in individual human cells uncovers widespread cell-cycle dependence of nuclear proteins.," *Nature Methods*, vol. 3, no. 7, pp. 525–531, 2006.

[266] A. Sigal, T. Danon, A. Cohen, R. Milo, N. Geva-Zatorsky, G. Lustig, Y. Liron, U. Alon, and N. Perzov, "Generation of a fluorescently

labeled endogenous protein library in living human cells.," *Nature Protocols*, vol. 2, pp. 1515–27, Jan. 2007.

[267] I. Issaeva, A. a. Cohen, E. Eden, C. Cohen-Saidon, T. Danon, L. Cohen, and U. Alon, "Generation of double-labeled reporter cell lines for studying co-dynamics of endogenous proteins in individual human cells.," *PloS One*, vol. 5, p. e13524, Jan. 2010.

[268] J. Livet, T. T. A. Weissman, H. Kang, R. W. R. Draft, J. Lu, R. A. R. Bennis, J. J. R. Sanes, and J. W. J. Lichtman, "Transgenic strategies for combinatorial expression of fluorescent proteins in the nervous system.," *Nature*, vol. 450, pp. 56–62, Nov. 2007.

[269] H. J. Snippert and H. Clevers, "Tracking adult stem cells.," *EMBO reports*, vol. 12, pp. 113–22, Feb. 2011.

[270] G. Neurohr and D. W. Gerlich, "Assays for mitotic chromosome condensation in live yeast and mammalian cells.," *Chromosome Research*, vol. 17, pp. 145–54, Jan. 2009.

[271] T. Kitamura, Y. Koshino, F. Shibata, T. Oki, H. Nakajima, T. Nosaka, and H. Kumagai, "Retrovirus-mediated gene transfer and expression cloning: powerful tools in functional genomics.," *Experimental hematology*, vol. 31, pp. 1007–14, Nov. 2003.

[272] J. M. Coffin, S. H. Hughes, H. E. Varmus, and V. Vogt, *Retroviral Virions and Genomes*. Cold Spring Harbor Laboratory Press, 1997.

[273] X. Wu, Y. Li, B. Crise, and S. M. Burgess, "Transcription start regions in the human genome are favored targets for MLV integration.," *Science (New York, N.Y.)*, vol. 300, pp. 1749–51, June 2003.

[274] R. S. Mitchell, B. F. Beitzel, A. R. W. Schroder, P. Shinn, H. Chen, C. C. Berry, J. R. Ecker, and F. D. Bushman, "Retroviral DNA integration: ASLV, HIV, and MLV show distinct target site preferences.," *PLoS Biology*, vol. 2, p. E234, Aug. 2004.

[275] C. Cattoglio, D. Pellin, E. Rizzi, G. Maruggi, G. Corti, F. Miselli, D. Sartori, A. Guffanti, C. Di Serio, A. Ambrosi, G. De Bellis, and F. Mavilio, "High-definition mapping of retroviral integration sites identifies active regulatory elements in human multipotent hematopoietic progenitors.," *Blood*, vol. 116, pp. 5507–17, Dec. 2010.

[276] D. G. Gibson, H. O. Smith, C. A. Hutchison, J. C. Venter, and C. Merryman, "Chemical synthesis of the mouse mitochondrial genome.," *Nature methods*, vol. 7, pp. 901–3, Nov. 2010.

[277] T. Roe, T. Reynolds, G. Yu, and P. Brown, "Integration of murine leukemia virus DNA depends on mitosis.," *The EMBO journal*, vol. 12, no. 5, pp. 2099–2108, 1993.