

Analysis of single-cell gene pair coexpression landscapes by stochastic kinetic modeling reveals gene-pair interactions in development

Cameron P. Gallivan¹, Honglei Ren^{2,3} and Elizabeth L. Read^{1,2,*}

¹*Dept. of Chemical & Biomolecular Engineering, University of California, Irvine, CA, USA*

²*NSF-Simons Center for Multiscale Cell Fate, University of California, Irvine, CA, USA*

³*Mathematical and Computational Systems Biology Graduate Program, University of California, Irvine, CA, USA*

Correspondence*:
Elizabeth L. Read
elread@uci.edu

2 ABSTRACT

3 Single-cell transcriptomics is advancing discovery of the molecular determinants of cell identity,
4 while spurring development of novel data analysis methods. Stochastic mathematical models
5 of gene regulatory networks help unravel the dynamic, molecular mechanisms underlying cell-
6 to-cell heterogeneity, and can thus aid interpretation of heterogeneous cell-states revealed by
7 single-cell measurements. However, integrating stochastic gene network models with single cell
8 data is challenging. Here, we present a method for analyzing single-cell gene-pair coexpression
9 patterns, based on biophysical models of stochastic gene expression and interaction dynamics.
10 We first developed a high-computational-throughput approach to stochastic modeling of gene-pair
11 coexpression landscapes, based on numerical solution of gene network Master Equations. We
12 then comprehensively catalogued coexpression patterns arising from tens of thousands of gene-
13 gene interaction models with different biochemical kinetic parameters and regulatory interactions.
14 From the computed landscapes, we obtain a low-dimensional “shape-space” describing distinct
15 types of coexpression patterns. We applied the theoretical results to analysis of published single
16 cell RNA sequencing data and uncovered complex dynamics of coexpression among gene pairs
17 during embryonic development. Our approach provides a generalizable framework for inferring
18 evolution of gene-gene interactions during critical cell-state transitions.

19 **Keywords:** stochastic modeling, gene expression noise, gene regulatory networks, single-cell data, scRNA-seq

1 INTRODUCTION

20 In recent years, single-cell-resolution measurements have revealed unprecedented levels of cell-to-cell
21 heterogeneity within tissues. The discovery of this ever-present heterogeneity is driving a more nuanced
22 view of cell phenotype, wherein cells exist along a continuum of cell-states, rather than conforming to
23 discrete classifications. The comprehensive view of diverse cell states revealed by single cell measurements
24 is also affording new opportunities to discover molecular regulators of cell phenotype and dynamics of
25 lineage commitment (Trapnell et al. (2014); Olsson et al. (2016); Briggs et al. (2018)). For example, single
26 cell transcriptomics have revealed the widespread nature of *multilineage priming* (MLP), a phenomenon
27 wherein individual, multipotent cells exhibit “promiscuous” coexpression of genes associated with distinct

28 lineages prior to commitment (Nimmo et al. (2015)). In principle, mathematical modeling of gene regulatory
29 network dynamics can provide a theoretical foundation for understanding cell heterogeneity and gene
30 expression dynamics, by quantitatively linking molecular-level regulatory mechanisms with observed cell
31 states. However, due to the molecular complexity of gene regulatory mechanisms, it remains challenging to
32 integrate such models with single-cell data.

33 Mathematical models of gene regulatory network dynamics can account for (and at least partially
34 reproduce) observed cellular heterogeneity in two primary ways. First, gene network models are multi-
35 stable dynamical systems, meaning a given network has the potential to reach multiple stable states of
36 gene expression. These states arise from the dynamic interplay of activation, inhibition, feedback, and
37 nonlinearity (Kauffman (1969); MacArthur et al. (2009); Huang (2012)). Second, some mathematical
38 models inherently treat cellular noise. This noise, or stochasticity, is modeled in various ways depending
39 on assumptions about the source (Peccoud and Ycart (1995); Arkin et al. (1998); Kepler and Elston
40 (2001); Swain et al. (2002)). Discrete, stochastic models of gene regulation, which track discrete molecular
41 entities, regulatory-protein binding kinetics, and binding states of promoters controlling gene activity, have
42 formed the basis of biophysical theories of gene expression noise due to so-called *intrinsic* molecular noise
43 (Peccoud and Ycart (1995); Thattai and van Oudenaarden (2001); Kepler and Elston (2001); Pedraza and
44 Paulsson (2008)). Such stochastic gene-regulation mechanisms have also been incorporated into larger
45 regulatory network models using the formalism of stochastic biochemical reaction networks, and have
46 been utilized to explore how molecular fluctuations can cause heterogeneity within phenotype-states and
47 promote stochastic transitions between phenotypes (Feng and Wang (2012); Sasai et al. (2013); Zhang and
48 Wolynes (2014); Tse et al. (2015)).

49 The quantitative *landscape* of cellular states is another concept that is increasingly utilized to describe
50 cellular heterogeneity. Broadly, the cellular potential landscape (first conceptualized by Waddington
51 (Waddington (2014); Wang et al. (2011); Huang (2012)) is a function in high-dimensional space (over
52 many molecular observables, typically expression levels of different genes), that quantifies the stability
53 of a given cell-state. In analogy to potential energy (gravitational, chemical, electric, etc.), cell states of
54 higher potential are less stable than those of lower potential. The landscape concept inherently accounts for
55 cellular heterogeneity, since it holds that a continuum of states is theoretically accessible to the cell, with
56 low-potential states (in “valleys”) more likely to be observed than high-potential states. The landscape is a
57 rigorously defined function derived from the dynamics of the underlying gene network model, according to
58 some choice of mathematical formalism (Wang et al. (2011); Bhattacharya et al. (2011); Huang (2012);
59 Zhou et al. (2016)). For stochastic gene network models that inherently treat noise, the landscape is directly
60 obtained from the computed probability distribution over cell-states (Cao and Liang (2008); Micheelsen
61 et al. (2010); Feng and Wang (2012); Tse et al. (2015)).

62 Stochastic modeling of gene network dynamics has been employed in various forms for analysis of
63 single cell measurements. For example, application of noisy dynamical systems theory has shed light on
64 cell-state transitions (Mojtahedi et al. (2016); Jin et al. (2018)). Stochastic simulations of gene network
65 dynamics have been used to benchmark tools for tasks such as network reconstruction (Schaffter et al.
66 (2011); Dibaeinia and Sinha (2019)). However, we are not aware of any existing analysis methods that
67 utilize discrete-molecule, stochastic models, which fully account for intrinsic gene expression noise and its
68 impact on cell-state, to aid interpretation of noisy distributions recovered from single cell measurements.
69 There exists an opportunity to link such biophysical, stochastic models, which reproduce intrinsic noise
70 and cell heterogeneity *in silico*, to single cell datasets that characterize cell heterogeneity *in vivo*. In

71 particular, the landscape of heterogeneous cell-states computed from discrete stochastic models can be
72 directly compared to single-cell measurements.

73 In this work, we present a method for analyzing single-cell gene pair coexpression patterns that is
74 founded on biophysical theory of stochastic gene networks. In our approach, the key object linking the
75 models to the data is the gene-pair coexpression landscape, which is derived directly from the bivariate
76 distribution of expression states, and which is computed from a stochastic model or extracted from
77 single cell measurements. The rationale underlying the method is two-fold: (1) information on gene-gene
78 interactions can be inferred from the distinctive characteristics of noise in single-cell data (i.e., from
79 the “shape” of the landscape); (2) existing analysis techniques are relatively insensitive to landscape
80 shape. We first comprehensively compute and classify the landscapes produced by a family of $\sim 40,000$
81 stochastic two-gene regulatory network models. We then use the model-derived classification to analyze
82 published data from vertebrate development. In so doing, we uncover both expected and novel patterns
83 of coexpression in development. While our analysis here is proof-of-principle, and limited to two-gene
84 interactions, the conceptual framework could be expanded to include multi-body gene interactions in the
85 future.

2 METHODS

86 2.1 Discrete, Stochastic Models of Two-Gene Regulatory Networks

87 We first developed a family of stochastic models of gene-gene interactions (see Fig. 1 for model
88 schematic), which is based on previously published models (Feng and Wang (2012); Zhang and Wolynes
89 (2014)). We label two genes X and Y . Each gene encodes a protein, which acts as a transcription factor
90 (TF) that potentially regulates its own expression as well as that of the other gene. Each gene has a promoter
91 (or more generally, regulatory regions of DNA) that can be bound by any combination of its own expressed
92 protein and/or the other gene’s expressed protein. The promoter states are thus labeled as: X_{00} (neither
93 transcription factor is bound to X ’s promoter), X_{0x} (X ’s own protein is bound, resulting in auto-regulation
94 of gene expression), X_{y0} (Y ’s protein is bound to X ’s promoter, resulting in cross-regulation), X_{yx} (both
95 proteins are bound to X ’s promoter, resulting in combinatorial regulation). (The promoter states for gene
96 Y are defined in a symmetric manner.) The regulatory effect of each promoter state (i.e., the effect of
97 having none, one, or both proteins bound on the gene’s expression) is accounted by the transcription rate
98 g_{ij} corresponding to each possible promoter state: e.g., when gene X ’s promoter is unbound, it transcribes
99 at rate g_{00}^X . Binding of Y ’s protein changes the transcription rate to g_{y0}^X , which may be lower, higher, or the
100 same, depending on whether the effect of Y on X is assumed to be repressing, activating, or not impacting.
101 (All other transcription rates for each promoter state and for gene Y are defined similarly.) The model
102 involves three classes of reactions: mRNA synthesis, mRNA degradation, and promoter-state-change
103 reactions. mRNA synthesis reactions are given by:



104 where x and y denote mRNA transcripts which will be translated into the transcription factors encoded by
105 genes X and Y , respectively. mRNA degradation reactions are given by:



106 Promoter-state-change reactions are given by, e.g.:



107 which represents the change of promoter-state (and corresponding regulatory impact) on gene X when Y 's
108 transcription factor binds (forward reaction) or unbinds (reverse reaction). All other promoter-state-change
109 reactions for X and Y are defined similarly. The changes of promoter state occur with forward rates $hy^2/2$
110 or $hx^2/2$ (when the change of state occurs due to binding of transcription factor from gene Y or X , respec-
111 tively) and f (when the change of state occurs due to an unbinding event). The model tracks copy numbers
112 of individual mRNA molecules in the cell, to enable direct comparison with single cell transcriptomic
113 data, but translation of mRNA into protein is not explicitly accounted for. Instead, transcription factor
114 (protein) levels are assumed to be linearly proportional to mRNA, and this proportionality constant is
115 subsumed into the binding rate h . The quadratic dependence of the forward binding rates on x or y arises
116 from the assumption that homodimeric transcription factors regulate gene expression, which is a general
117 and convenient way to include cooperativity in the model.

118 We assign rate constants to intracellular processes that are in line with experimental estimates from
119 vertebrates, where possible (see Table 1). (For full details of model reactions and parameter derivations, see
120 Supplement). Rates of mRNA synthesis and degradation are relatively well characterized, though they vary
121 considerably for different transcripts (Schwanhäusser et al. (2011)). Rates of promoter-state-change are
122 less well-defined, since promoter-state-changes that ultimately impact gene expression may be attributed
123 to a variety of molecular processes, including: (a) relatively fast processes of TF binding or unbinding
124 from DNA (b) relatively slow chromatin remodeling processes that may be initiated or facilitated by TF
125 binding, require multiple steps and cooperative interactions, and are generally poorly understood. In our
126 models, to account for this range of possible mechanisms, we consider a wide range of parameter values
127 h, f for promoter-state-changes. (The significance of these fast and slow regimes, termed the *adiabatic* and
128 *nonadiabatic* regimes, respectively, to cell-state stability has been studied previously by stochastic modeling
129 (Sasai and Wolynes (2003); Feng and Wang (2012); Sasai et al. (2013); Zhang and Wolynes (2014))). We
130 here define the “fast” regime as determined by measured parameter values of protein binding/unbinding
131 DNA (e.g., from Geertz et al. (2012)), occurring with timescales of minutes, seconds, or faster. We define
132 the “slow” regime more broadly as any epigenetic/chromatin changes occurring on timescales of hours,
133 days, or longer. For example, in mammalian cells, changes of chromatin state during cell-fate specification
134 were estimated to be on the order of several days (Hathaway et al. (2012); Mariani et al. (2010)), while
135 theoretical studies predicted timescales on the order of the cell cycle time (i.e., hours to days, Sasai et al.
136 (2013)).

137 We define two types of model systems. The **Mutual Inhibition/Self-Activation (MISA)** model encodes
138 a common network motif that is understood to control a variety of cell fate decisions (Graf and Enver
139 (2009); Huang (2013)) and has been extensively studied by mathematical modeling (Huang et al. (2007);

Rate Constant	Symbol	Units	Value	Comments/Source
mRNA synthesis (not repressed)	g_{hi}	mRNA/hr	0.8 - 1.4*	Schwanhäusser et al. (2011)
mRNA synthesis (repressed)	g_{lo}	mRNA/hr	0.001	see text
mRNA degradation	k	/hr	0.2^{\ddagger}	Schwanhäusser et al. (2011)
Promoter state change (unbinding)	f^{\dagger}	/hr	(fast) $10 - 10^5$ (slow) $10^{-6} - 10$	Geertz et al. (2012) see text
Promoter state change (binding)	h^{\dagger}	$\text{hr}^{-1} \text{mRNA}^{-2}$	(fast) $10 - 500$ (slow) $10^{-6} - 10$	Geertz et al. (2012) see text

Table 1. Rate Parameters used in gene regulatory network models. Parameter values are derived from experimental measurements in vertebrates, where possible. See Methods text for details. *Measured rates of mRNA synthesis varied, with a median of 2/hr Schwanhäusser et al. (2011)). We use lower values (within experimental range) to roughly match observed counts in scRNA-seq data, which may be lower than expected because of dropouts or other technical issues. \ddagger Corresponds to mRNA half-life of 3.5 hours, which is well within experimentally measured values but shorter than the median value of 9 hours, assuming that transcriptional regulators have shorter-than-average half-lives in the cell. \dagger Promoter state change rates f and k are reported in fast and slow regimes. Fast promoter state changes are assumed to occur due to TF-DNA unbinding or binding events, with rate parameters chosen based on values reported in Geertz et al. (2012) (see Supplement for details on parameter derivation and unit conversion). Slow promoter state changes are thought to involve collective changes in epigenetic marks and rearrangement of chromatin.

140 Feng and Wang (2012); Chu et al. (2017)). In contrast, the **Two-Gene Flex** model flexibly encodes a
141 variety of regulatory interactions, as described below.

142 2.1.1 Mutual Inhibition/Self-Activation Model

143 In all models, promoter activity is assumed to be either high (transcription rate g_{hi}) or low (g_{lo}) (giving a
144 relatively fast or slow rate of mRNA synthesis, respectively). To encode MISA regulatory logic, mRNA
145 synthesis rates for each promoter state are $\{g_{00}^X, g_{0x}^X, g_{y0}^X, g_{yx}^X\} = \{g_{lo}, g_{hi}, g_{lo}, g_{lo}\}$. Transcription rates for
146 gene Y are defined symmetrically, $\{g_{00}^Y, g_{0y}^Y, g_{x0}^Y, g_{yx}^Y\} = \{g_{lo}, g_{hi}, g_{lo}, g_{lo}\}$. The high rate corresponds to
147 maximal activity, whereas the low rate is effectively off (but is non-zero to allow for some leakiness in the
148 promoter). Thus, binding of the self-TF turns the gene on, but subsequent binding of the other TF turns
149 the gene off. The relative strengths and kinetics of the activating (self-regulatory) and repressing (cross-
150 regulatory) interactions are encoded in the rates of binding/unbinding of regulators. Autoregulatory binding
151 and unbinding rates (symmetric on both genes) are denoted by h_a and f_a , respectively. Cross-regulatory
152 rates are denoted by h_r and f_r . The model is thus fully specified by 7 parameters: $\{g_{lo}, g_{hi}, k, h_a, f_a, h_r, f_r\}$.
153 We computed landscapes for $\sim 22,000$ unique parameter combinations for the MISA regulatory logic (see
154 Table 1 for parameter value ranges). We studied only symmetric network motifs, but asymmetry between
155 the genes is accounted for by allowing the “on” transcription rate g_{hi} to be asymmetric between the two
156 genes (in case of asymmetry in g_{hi} , the model is specified by eight parameters).

157 2.1.2 Two-Gene Flex Model

158 The Two-Gene Flex model is identical to MISA in all ways except the regulatory logic. Instead of
159 the transcription rates being $\{g_{lo}, g_{hi}, g_{lo}, g_{lo}\}$, all 16 logical combinations of four promoter states and
160 two activity-levels are included. Within these combinations, various behavior is encoded including self-
161 activation, self-repression, mutual activation, mutual repression, no interaction (self- or cross-), and
162 dual-effects (where a TF has a distinct effect whether bound alone or in combination with the other). Note
163 that the MISA logic is contained within these 16 combinations. Note also that the promoter states for X
164 and Y are always defined symmetrically, i.e., only symmetric motifs are included. We computed landscapes

165 for ~34,000 unique parameter combinations for the Two-Gene Flex Model (including all network motif variants).

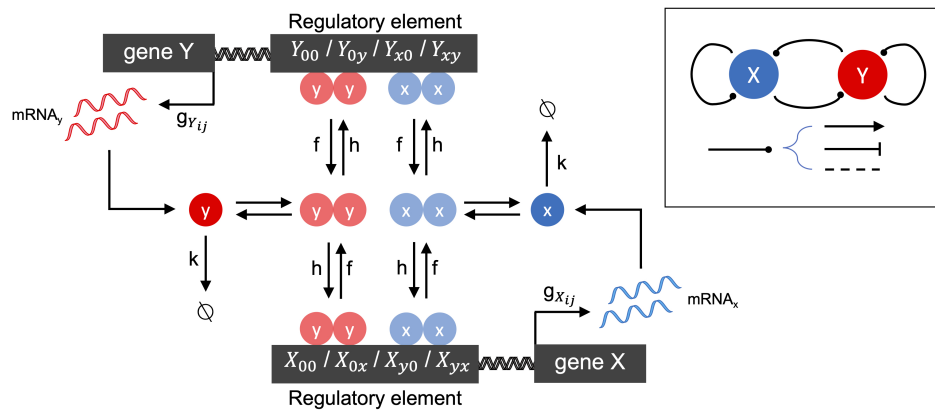


Figure 1. Schematic of the two-gene regulatory network model. The overall network motif is variable (see Inset), encoding a symmetric combination of repression (flat arrow-head), activation (pointed arrow-head) or no-impact (dashed line), mutually between the two genes labeled X and Y , and by each gene on itself (see Methods for details). The stochastic reaction kinetic model includes rate constants for mRNA synthesis (g_{ij}), mRNA degradation (k), and regulatory element state-changes due to transcription factor binding (h) and unbinding (f). Cooperative effects are included by the assumption that transcription factors bind as homodimers.

166

167 2.2 Mathematical Framework: Chemical Master Equation

168 2.2.1 Chemical Master Equation

169 Stochastic dynamics for the above-described network motifs are modeled by a Chemical Master Equation
 170 (CME) (alternatively known as a discrete space, continuous time Markov Chain). The instantaneous state
 171 of the system is given by the vector \mathbf{n} , which enumerates the mRNA copy numbers and promoter-states
 172 of both genes, i.e., $\mathbf{n} = [n_x, n_y, X_{ij}, Y_{ij}]$, where n_x is the mRNA copy number for gene X , X_{ij} is the
 173 promoter state for gene X , and so on. The CME gives the probability for the system to exist in a given
 174 state at a given time, $\mathbf{p}(\mathbf{n}, t)$. The CME can be written in vector-matrix form as a linear system

$$\frac{d\mathbf{p}(\mathbf{n}, t)}{dt} = \mathbf{K}\mathbf{p}(\mathbf{n}, t) \quad (4)$$

175 where \mathbf{K} is the reaction rate-matrix. Each off-diagonal element K_{lm} gives the rate of transitioning from
 176 state m to l (non-zero values correspond to allowed state transitions with rates according to reactions 1-3
 177 above), while the diagonal elements are the summed rates for exiting each state, $K_{ll} = -\sum_{m \neq l} K_{ml}$.
 178 Transition rates are computed according to standard stochastic chemical kinetic rate laws (Gillespie (1977)).
 179 If both types of mRNA are assumed to exist in the cell in copy numbers that never exceed $M - 1$, then the
 180 total size of the enumerated space including all possible states is $N = M \times M \times 4 \times 4$ (note that the total
 181 number of mRNA copy number states includes the state of 0 copies, thus $n_x, n_y \in \{0, 1, \dots, M - 1\}$).

182 2.2.2 Computing Gene Pair Coexpression Landscapes

183 The complete steady state probability to find a cell in state \mathbf{n} is given by the vector $\boldsymbol{\pi}(\mathbf{n}) = \mathbf{p}(\mathbf{n}, t \rightarrow \infty)$,
 184 which is obtained from Eq. 4 using eigenvalue routines in numpy and scipy (van der Walt et al. (2011))

185 (McKinney (2010)). Each individual model requires solution of an N -state system, where N is $\mathcal{O}(10^4)$
186 (e.g., assuming the probability to have mRNA exceed 25 is negligible, then $N = 10,816$). Efficient
187 computation of the landscapes over tens of thousands of model variants/parameter combinations was
188 achieved using routines compiled with the numba library (Lam et al. (2015)) and parallelization using
189 Python's multiprocessing library to distribute the workload across the available cores.

190 To mimic experimental scRNA-seq data, the probability is projected onto the mRNA subspace by
191 summation over all promoter state combinations. We hereon define the gene pair coexpression landscape
192 as the steady-state probability to find a cell with mRNA count numbers (n_x, n_y) . More precisely, the
193 *probability landscape* is the vector π with each element π_i giving the steady-state probability for the cell to
194 be found in state i with the combination of mRNA counts (n_x, n_y) from genes X and Y , and $i \in 1, \dots, M^2$.
195 Alternatively, the *quasipotential landscape* is log-transformed, given by the vector ϕ where $\phi_i = -\ln(\pi_i)$.

196 2.3 scRNA-seq Data Acquisition, and Landscape Estimation

197 Experimental data is obtained from the published single cell RNA sequencing (scRNA-seq) measurements
198 of Briggs et al. (2018). The dataset "Corrected_combined.annotated_counts.tsv" was used which provides
199 the normalized transcriptome profiles for *Xenopus tropicalis* at single cell resolution for ten different stages
200 of embryonic development, with labelled cell types and parent cell types. We analyzed 1380 gene pairs,
201 which were identified as putative MLP pairs in Briggs et al. (2018), based on their estimated changes in
202 coexpression over the course of development. Gene pairs were identified by their developmental stage
203 and lineage branch point in which coexpression was maximal. Cell types from other stages were then
204 included in the lineage if they were a parent (preceding in development) cell type or daughter (descendant
205 later in development) cell type. After selecting the desired gene pair and cell/tissue/cluster type of interest,
206 gene pair counts were combined and summed resulting in ten gene pair landscapes, one for each stage of
207 development, in cells of the relevant lineage.

208 To directly compare computed coexpression landscapes with experimental data, we extracted cell count
209 matrices for each gene pair, and where necessary, truncated to mRNA count numbers $\leq M - 1$ (truncation
210 eliminated less than 0.5% of cells in the data, across all gene pairs and cell stages). This produces an
211 $M \times M$ (including zeros) count matrix that serves as a sampled estimator of the steady-state distribution,
212 $\tilde{\pi}(\mathbf{n})$, of the same size as computed landscapes. In order to compute the sampled quasipotential landscape,
213 we use $\tilde{\phi}(\mathbf{n}) = -\ln\tilde{\pi}(\mathbf{n})$, after replacing the not-observed count-combinations with a low but non-zero
214 estimate of these probabilities (since log of zero is undefined). We use a general estimate of 1E-6 for
215 non-observed counts, both because it is in line with the predictions of the theoretical models for the low
216 probability edges of the distributions, and because it is less than the lowest estimable probability (i.e.,
217 observation of one cell in a given matrix position, given total cell counts on the order of 10^5 , would
218 correspond to an estimated probability of 1E-5).

219 2.4 Dimensionality Reduction for Landscape Shape-space

220 We apply Principal Component Analysis (PCA) to the theoretically computed landscapes over the model
221 sets to achieve a reduced-dimension description of landscape shape. All PCA training and dimensionality
222 reduction was performed using the decomposition module of the python package scikit learn. Each unique
223 model is treated as a replicate and the steady-state probability π_i (or alternatively, quasipotential ϕ_i) of
224 each of the $M \times M$ possible mRNA copy-number states (n_x, n_y) is treated as a feature.

225 The principal components obtained from the model set were then used to fit the experimental data, where
226 each landscape from each gene-pair/stage is a replicate.

227 2.5 Clustering of Developmental Landscape-Shape Trajectories

228 By viewing the time-ordered coexpression landscapes of a given gene pair in PCA space, termed
229 “landscape-shape trajectories”, one can gain insight into the genes’ roles in development. The trajectories
230 were hierarchically clustered based on their geometric distance in PCA space. More specifically, the
231 *fcluster* method in scikit-learn package was used in hierarchical clustering (McKinney (2010)), and the
232 geometric distance between trajectories A and B were defined as the sum of the pair-wised Euclidean
233 distance between two corresponding stages, i.e.

$$\|A - B\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n (A_{i,j} - B_{i,j})^2} \quad (5)$$

234 where $\|\cdot\|_F$ is the Frobenius norm, A and B are two trajectories represented by m by n matrices, m is the
235 number of developmental stages in single cell data, n is the number of PCA components used in clustering.

3 RESULTS

236 3.1 Stochastic two-gene network models show a variety of coexpression landscape 237 shapes, distinguishable by Principal Component Analysis

238 Our modeling framework enabled efficient computation of coexpression landscapes resulting from
239 discrete, stochastic gene network models. This in turn enabled us to compute landscapes for tens of
240 thousands of parameter sets, encompassing both various relative strengths and kinetics of regulatory
241 interactions, as well as different schemes of regulatory logic among the two genes (see Methods). This
242 approach afforded a comprehensive view of theoretically predicted landscape shapes resulting from
243 gene-gene interactions (within the assumptions of the current model system).

244 We applied Principal Component Analysis to the computed probability landscapes for Two-Gene Flex,
245 in order to find a low-dimensional description of their shapes (Fig2). The first two PCA components
246 encompass 98% percent of total covariance, and all models fall within a triangular region of this 2D
247 subspace. The vertices of the triangle correspond generally to landscapes with: (1) very low expression of
248 both genes (i.e., transcript levels of X/Y are *lo/lo*, Fig2E), (2) high simultaneous expression of both genes
249 (*hi/hi*, Fig2C), and (3) expression of only one gene at a time (*hi/lo* and *lo/hi*, Fig2A). Landscapes located
250 away from the vertices are thus well-described by some linear combination of these three shapes, consistent
251 with PCA, and supported by visual inspection. In all, the results reveal that two-gene interaction motifs can
252 encode a wide variety of patterns of coexpression, including mixtures of all combinations of *lo/lo*, *hi/hi*
253 and *lo/hi*, *hi/lo* phenotypes (e.g., Fig2B). At the same time, this variety of shapes is well-described by a
254 small number of principal components (which form a basis for what we term the “shape-space”), and we
255 hereon use the magnitudes along these components as measures of landscape shape.

256 3.2 Shape measures of coexpression landscapes distinguish different types of mutual 257 gene-gene interactions

258 We sought to understand how different regulatory motifs contributed to landscape shape. Projecting
259 the landscapes arising from each network motif separately revealed distinctive patterns (i.e., occupying
260 distinct, but overlapping, regions of the PCA triangle) (approximately 2,000 landscapes were computed
261 for each network motif, i.e., $\sim 2,000$ models that share the regulatory logic but have different kinetic
262 parameters). We grouped all motifs according to their region of occupancy within the PCA triangle, and

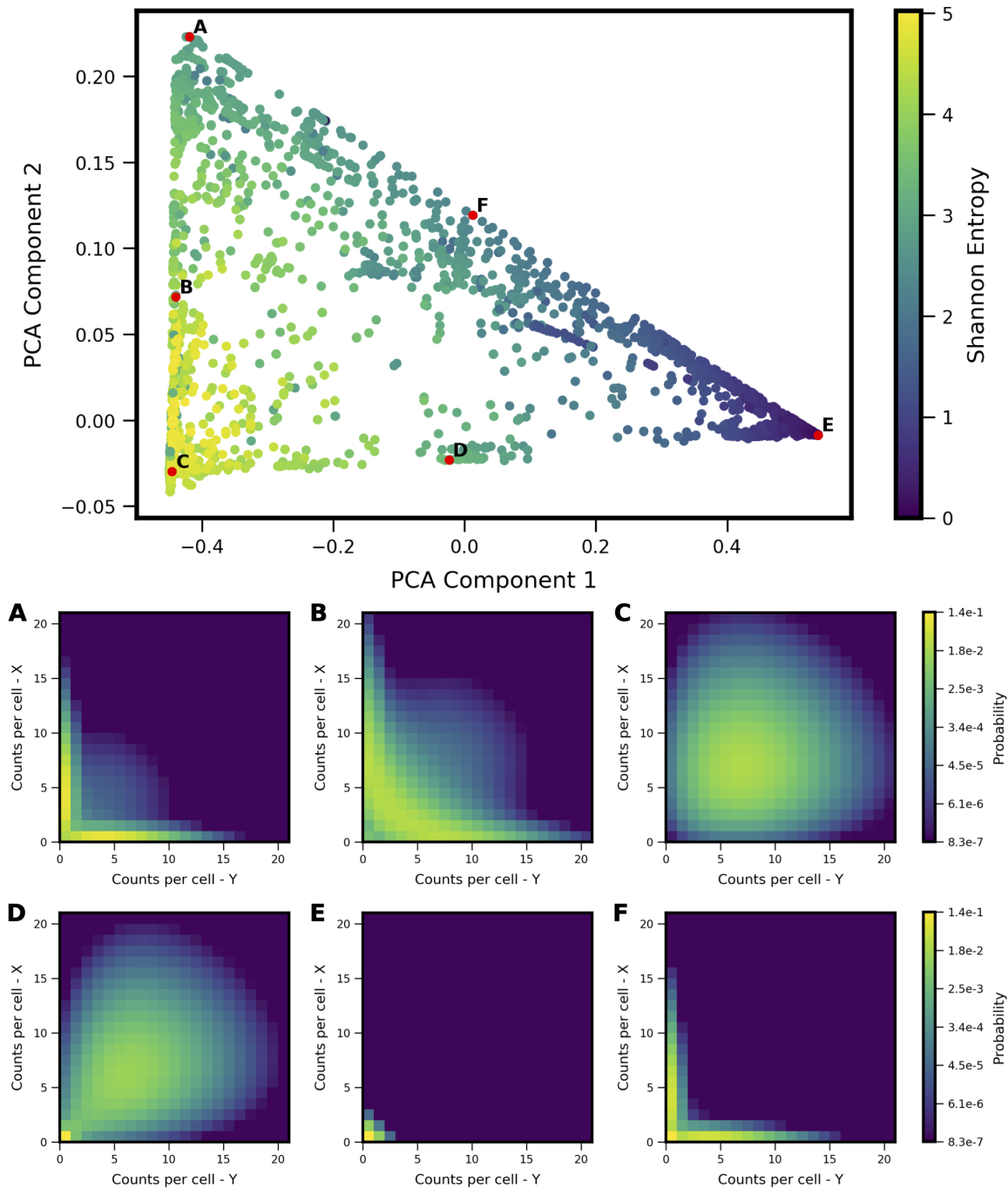


Figure 2. Shape-space of simulated Two-Gene Flex coexpression landscapes analyzed by PCA. Co-expression landscapes were computed for 34,097 unique two-gene stochastic network models with varying regulatory interactions and kinetic rate parameters (see Model schematic in Fig. 1). (Top) All model landscapes projected onto the first two principal components. Each dot corresponds to one model, colored by the model's Shannon Entropy. (Bottom) Representative quasipotential landscapes $\phi(n)$ (see Text) of individual models from different regions of PCA component-space. Color of each discrete grid space in $\{x, y\}$ corresponds to computed probability (in log-scale) to find a single cell with the corresponding numbers of $\{x, y\}$ transcripts.

263 discovered logical consistency among the groups (see Fig. 3). For example, all motifs with some type of
264 mutual activation were found to co-occupy a region of PCA shape space in the lower part of the triangle
265 (3A). This result is consistent with the intuition that motifs with mutual activation cannot produce the

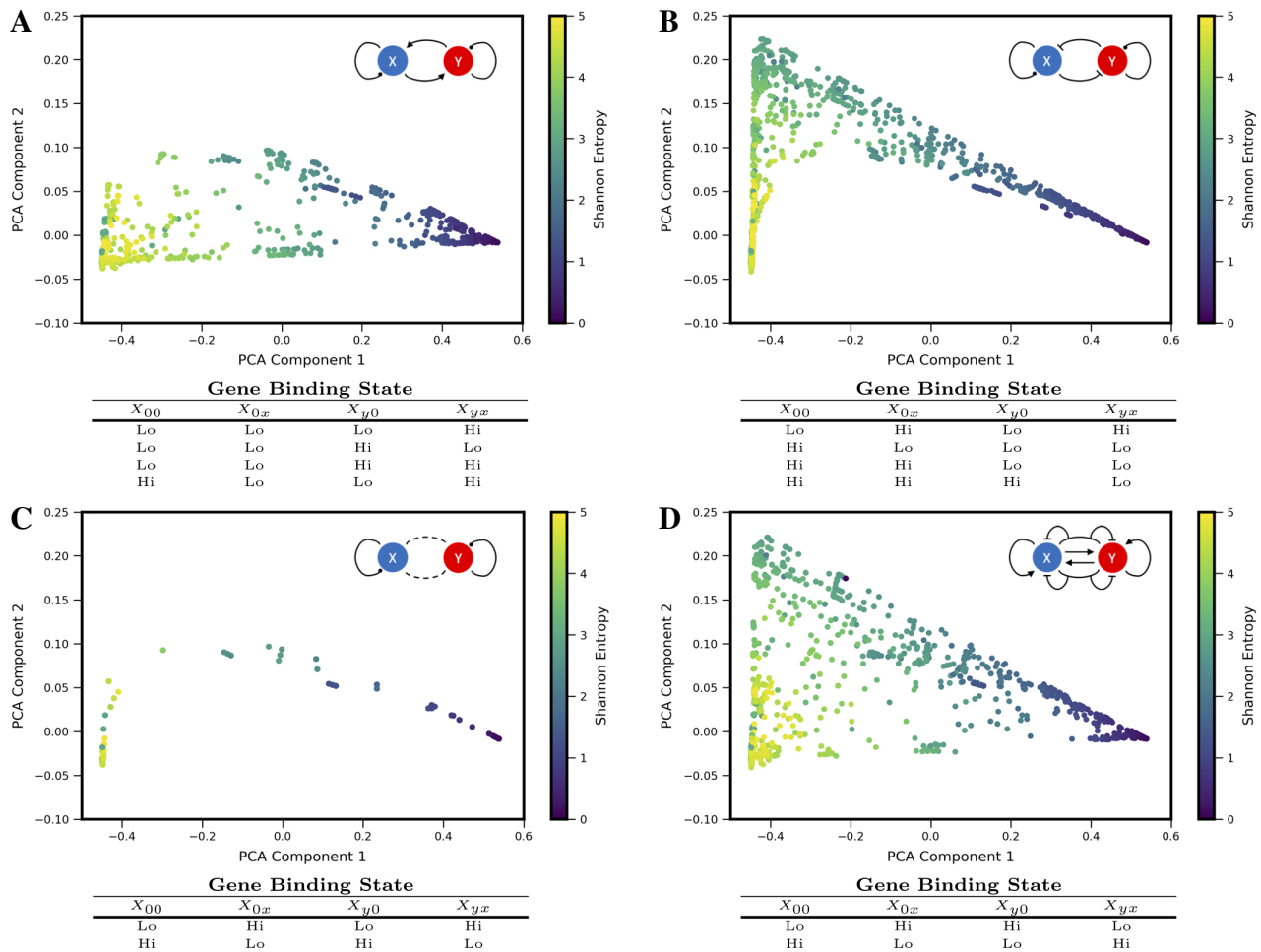


Figure 3. Coexpression landscapes computed from the Two-Gene Flex models show distinctive shapes that depend on the regulatory logic of gene-gene interactions. The Two-Gene Flex model encodes 16 logical combinations (2^4) of gene-gene interactions, corresponding to four possible promoter-binding states and two possible levels of transcription activity (low and high). These 16 model variants can be grouped into motif classes: (A) Models with mutual activation. (B) Models with mutual repression. (C) No mutual gene-gene interactions. (D) “Incoherent” models, where the combinatorial-binding state has the opposite behavior of both of the singly-bound states (see text). Within each motif class, different kinetic parameters serve to modify the relative strength of interactions (i.e., different weights on the edges). Each motif class occupies a distinct, but overlapping, region of the shape space (with the exception of the Incoherent motif, which can reach all areas of the shape space).

266 apparent bistability seen in landscapes at the hi/lo-lo/hi vertex of the triangle. The other three motif
 267 groupings include motifs with some type of mutual repression, motifs with no inter-gene interactions, and
 268 incoherent motifs with dual-interactions (when the regulator bound by itself has the opposite effect of the
 269 regulator bound in combination with the other TF). Note that two of the sixteen logical combinations of
 270 promoter binding-states in the Two-Gene Flex models are not included here, since they effectively encode
 271 no gene-gene interactions (the “always on” or “always off” logic, $\{g_{hi}, g_{hi}, g_{hi}, g_{hi}\}$ or $\{g_{lo}, g_{lo}, g_{lo}, g_{lo}\}$).
 272 Note that here we assess all kinetic parameter combinations associated to one regulatory motif; these
 273 parameters tune the strength of different interactions. As such, the analysis of 3 assumes fixed network
 274 topologies but variable weights on network edges, accounting for the overlap between different motifs.
 275 These results indicate that landscape shape can to some extent be used to distinguish regulatory interactions
 276 between pairs of genes, despite variable and/or unknown kinetics governing the interactions.

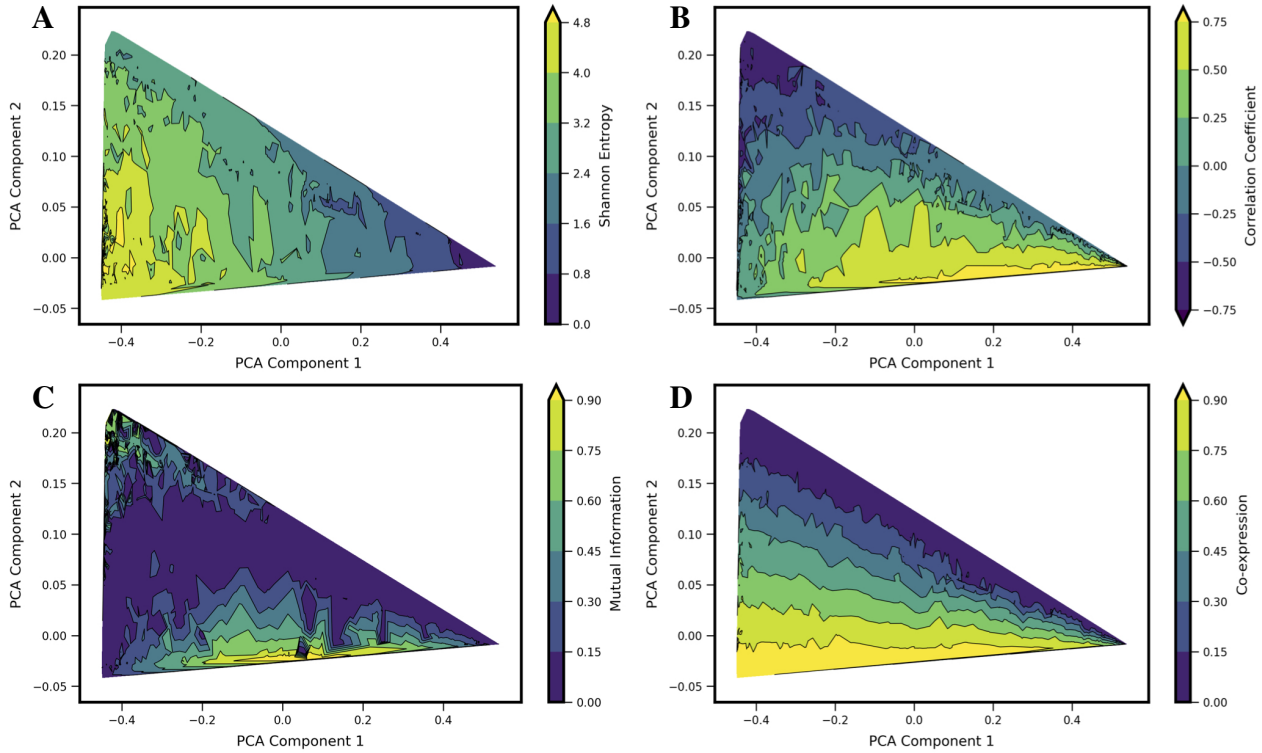


Figure 4. Comparison of four standard metrics of gene-gene coexpression with landscape shape. Metrics include: (A) Shannon Entropy. (B) Correlation Coefficient. (C) Mutual Information. (D) Co-expression Index (see text for details). Each metric was computed for each computed model landscape, using the same set of 34,097 Two-Gene Flex models as in Figs. 2 and 3. Contour plots show each metric as a function of principal components 1 and 2, obtained by local averaging and interpolation over the results from individual model landscapes. Taken together with Fig. 2, the results show how these metrics correspond with landscape shape.

277 3.3 Commonly used pairwise metrics are relatively insensitive to coexpression 278 landscape shape

279 In order to analyze how previously-applied measures of gene-gene interactions align with landscape
280 shape, we computed a set of metrics for each model landscape and visualized the resultant values projected
281 onto the PCA subspace. We chose four metrics: Shannon Entropy, Pearson Correlation Coefficient, Mutual
282 Information, and a Coexpression Index (see Fig. 4, note Shannon Entropy is visualized also in Figs. 2
283 and 3). The first three of these are obtained directly from the computed bivariate probability distributions
284 according to standard definitions; the Coexpression Index has been used previously (Briggs et al. (2018))
285 and is given by the conditional probability to find cells with non-zero counts of both mRNA x and y
286 (conditioned on the cells having non-zero counts of at least one of genes X or Y). Here, for a given model
287 j , we derive this metric from the probability landscape π over count-states i by:

$$m_{j,\text{Coex.Index}} = \frac{\sum_{i \in n_x > 0 \cap n_y > 0} \pi_i}{\sum_{i \in n_x > 0 \cup n_y > 0} \pi_i}. \quad (6)$$

288 We estimate the value of each metric as a function of landscape shape (that is, we estimate the function
289 $m(c_1, c_2)$, where m is a given metric and (c_1, c_2) are the coordinate values in PCA components 1 and 2). For
290 each of the four metrics, we estimate and visualize this function by local averaging and interpolation over
291 the computed results for each individual model landscape. We found that each metric aligns in distinctive,

292 and generally intuitive, ways with the PCA landscape shape space. High or low values of each metric were
293 to some extent localized to particular sub-regions of the triangle, and thus could be understood to be arising
294 from landscapes of similar shape. However, numerous examples can also be found of models colocated (or
295 nearly colocated) in the triangle but having different values of a given metric, so the functional dependence
296 $m(c_1, c_2)$ is noisy.

297 For Shannon entropy, the highest values are generally seen near the hi/hi vertex of the triangle, while
298 the lowest values are seen near the lo/lo vertex. This reflects the amount of disorder in the hi/hi state
299 of expression, in which a broad range of count-values are possible for each gene, whereas in the in the
300 lo/lo vertex, count values are always zero or near-zero. The noise in expression levels can be quantified
301 more precisely for the subset of models in the “slow-binding” regime ($h, f \ll g, k$). In this parameter
302 regime, cells show distinctive high (“hi”) and/or low (“lo”) expression states with mean counts g_{hi}/k and
303 g_{lo}/k , respectively, and the disorder in each expression state can be quantified as Poisson birth/death noise
304 (Al-Radhawi et al. (2019)), such that variance scales linearly with the expression rate g . Sources of disorder
305 contributing to higher values of Shannon entropy include both noisy expression within a given phenotype
306 state and the ability for cells to exist in multiple different phenotype states (i.e., the breadth of a valley in
307 the potential landscape, and the number of different valleys). Notably, in the parameter regimes studied
308 here, the highest Shannon Entropy models are single-phenotype (hi/hi), indicating that the noise in this one
309 state contributes more disorder than does noise from multiple phenotype-states. As such, models with two
310 or more accessible states have intermediate values of Shannon entropy.

311 A strongly negative correlation coefficient between the two genes is found near the lo/hi-hi/lo vertex
312 of the triangle, which is occupied by models showing bistability (cells can express one gene or the other,
313 but not both simultaneously) resulting from mutual repression in the network motif. Landscapes with
314 high positive correlation tend to be those that combine expression in the hi/hi and lo/lo quadrants of the
315 two dimensional subspace (see, e.g. 4B and 2D), resulting from mutual activation in the network motif.
316 Mutual Information aligns somewhat with large absolute values of correlation coefficients, but cannot
317 distinguish high positive from high negative correlation. Mutual Information values near zero co-localize
318 with Correlation coefficients near zero. This arc-shaped region bisecting the triangle also overlaps with the
319 models lacking interactions between the two genes (see Fig. 3C).

320 The Coexpression Index shows the smoothest functional dependence on PCA components (c_1, c_2). Of
321 note, the model-subspace of high coexpression is not fully overlapping with the subspace of high correlation
322 coefficients. This reflects the fact that high simultaneous expression occurs in both genes in an uncorrelated
323 manner, since the noise arises from aforementioned birth-death noise of mRNA transcription/degradation.

324 None of the four metrics are by themselves able to fully differentiate between landscape shapes. For
325 example, model landscapes with similarly high values of Mutual Information include both hi/lo-lo/hi
326 landscapes from mutual repression motifs and hi/hi-lo/lo landscapes from mutual activation motifs. (see,
327 e.g., Fig. 4A and B). Model landscapes with similar intermediate values of Coexpression Index also
328 encompass a variety of landscape shapes, including some that arise from different network motifs (see,
329 e.g., Fig. 4C and D). Taken together, these results show that these four single metrics are not reliable
330 determinants of landscape shape. They furthermore show that a given value for commonly used measures,
331 as obtained from experimental data, can potentially arise from a variety of regulatory scenarios.

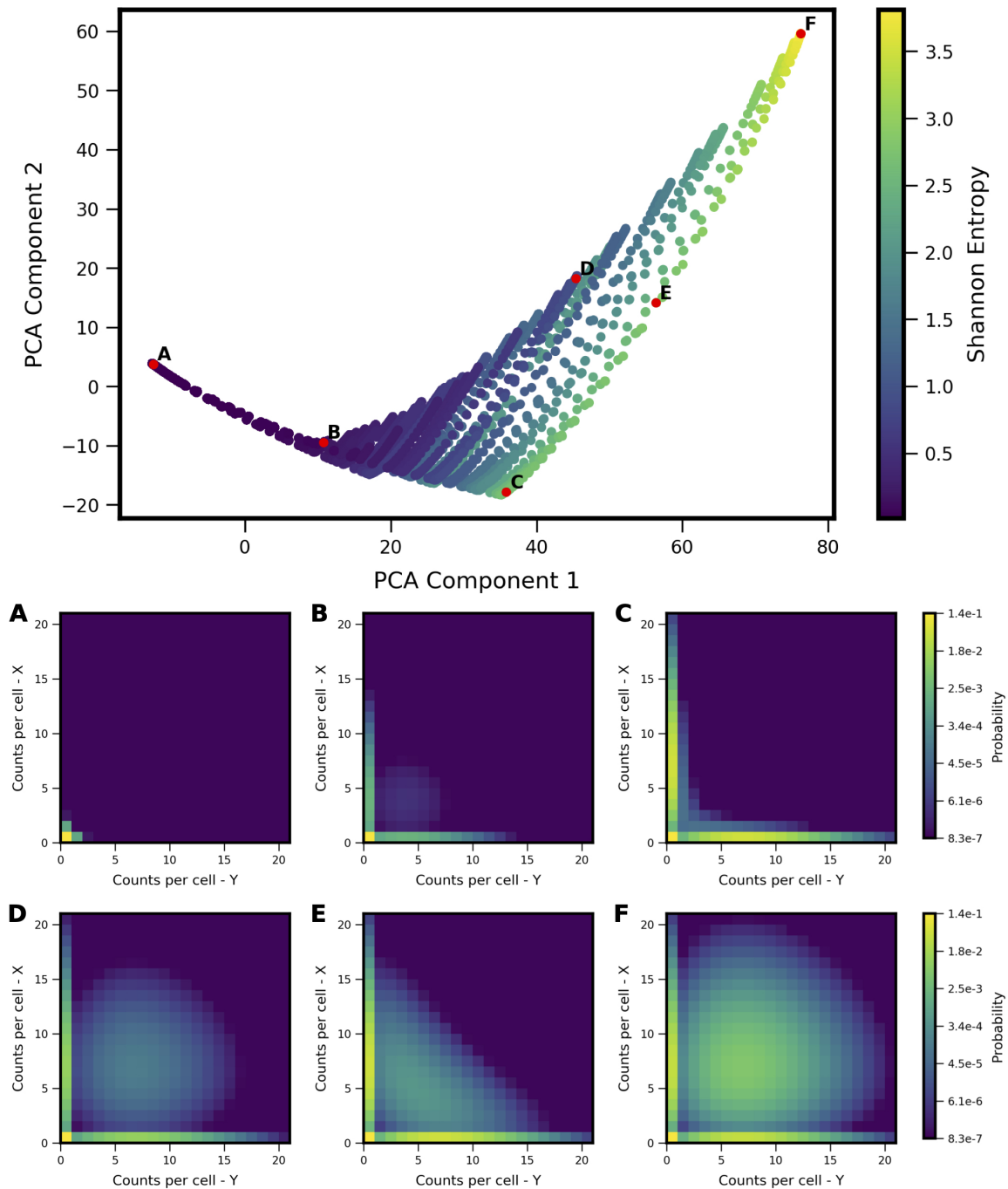


Figure 5. Shape-space of simulated MISA coexpression landscapes analyzed by PCA. Coexpression landscapes were computed for 22,718 unique two-gene stochastic network models with MISA logic and varying kinetic rate parameters. Promoter-state change rates were restricted to the fast regime (see Table 1). (Top) All model landscapes projected onto the first two principal components. Each dot corresponds to one model, colored by the model's Shannon Entropy. (Bottom) Representative quasipotential landscapes $\phi(n)$ (see Text) of individual models from different regions of PCA component-space. Color of each discrete grid space in $\{x, y\}$ corresponds to computed probability (in log-scale) to find a single cell with the corresponding numbers of $\{x, y\}$ transcripts. (Analogous to Figure 2).

3.4 Stochastic theory-based analysis of coexpression landscapes from single-cell experiments reveals distinct developmental “landscape shape” trajectories

We applied the landscape shape analysis framework, developed above on the basis of theoretical models, to publicly available single cell RNA sequencing data in vertebrate development. We applied the analysis to putative MLP gene pairs in *Xenopus tropicalis* development (Briggs et al. (2018)). To carry out the analysis, we first analyzed the landscape shape-space for a restricted set of theoretical models, which encode only the MISA interaction motif. The MISA motif has been previously discovered to operate at critical cell-fate branch points (Graf and Enver (2009)) and has potential to enable both antagonistic expression and coexpression of genes in individual cells (depending on kinetic parameters), as is characteristic of MLP gene-pairs. We first generated a MISA-specific set of models for training the PCA shape analysis. In addition to restriction of the network motif, there were two other differences between the MISA-model training set (Fig.5) and the Two-gene Flex-model training set (Fig.2). For MISA, we utilized quasipotential landscapes, rather than probability landscapes, in order to increase sensitivity to rarer cell-states (i.e., weaker landscape features). We furthermore restricted the kinetic parameters h, f to the fast (adiabatic) regime (see Table 1), in order to use the models to analyze time-resolved data. That is, the experiments measure embryos at different developmental stages, which are roughly 1-3 hours apart in time. We compare the steady-state landscapes from stochastic models to the experiment-derived landscapes at different timepoints by applying a quasi-steady-state assumption: we assume that the promoter-binding states (which govern gene activity) reach equilibrium faster than the progression of developmental stage, which is valid only in the adiabatic regime. Despite these modifications to the model training set, the projection of models onto the PCA subspace for MISA (Fig. 5) shows qualitative similarity to that of Two-gene Flex ((Fig. 2), including delineation of a subregion of a triangle (note that the triangle is inverted between the two figures, which is an arbitrary result of eigenvector sign invariance). However, antagonistic expression of the two genes is a stronger feature across models in the MISA training set, such that the hi/hi vertex of the triangle for MISA still shows considerable probability for cells to antagonistically express one gene or the other (Fig. 5F).

We extracted two-gene coexpression quasipotential landscapes corresponding to distinct developmental stages from the dataset of Briggs, et al. We then projected the landscapes onto the PCA subspace, and thereby derived developmental trajectories through landscape shape-space. By way of illustration, we first present developmental trajectories for three representative gene pairs (Fig. 6). *Gata5* and *pax8* were identified (in Briggs, et al.) as being antagonistically expressed within the intermediate mesoderm lineage, in cardiac mesoderm and pronephric mesenchyme cell subtypes, respectively. In contrast, *lhx1* and *pax8* were shown to co-express in cells of the pronephric mesenchyme. Finally, the gene pair *sox2* and *brachyury* (t) has been identified as influencing the cell fate decision between the neural plate and the dorsal marginal zone (Wardle and Smith (2004)), and was identified as presenting MLP behavior, characterized by high coexpression at some stage of development, followed by antagonistic expression at a later stage (Briggs et al. (2018)). We found that these three gene pairs showed distinctive trajectories through PCA subspace. All of the genes showed low expression early in development (stage 8) and their landscapes were colocated near the lo/lo vertex in the model subspace. Their trajectories then diverged: *gata5-pax8* travels along the bistable edge of the triangle, increasing expression of both genes over the course of development, but in largely non-overlapping subpopulations of cells. In contrast, *lhx1-pax8* shows strong coexpression starting at stage 14, and continues thereafter to move toward increasing values of PCA component 2, which coincides with increasing coexpression. (*lhx1-pax8* landscapes for some of the measured developmental stages fall slightly outside the area reached by MISA models in the training set, suggesting that the

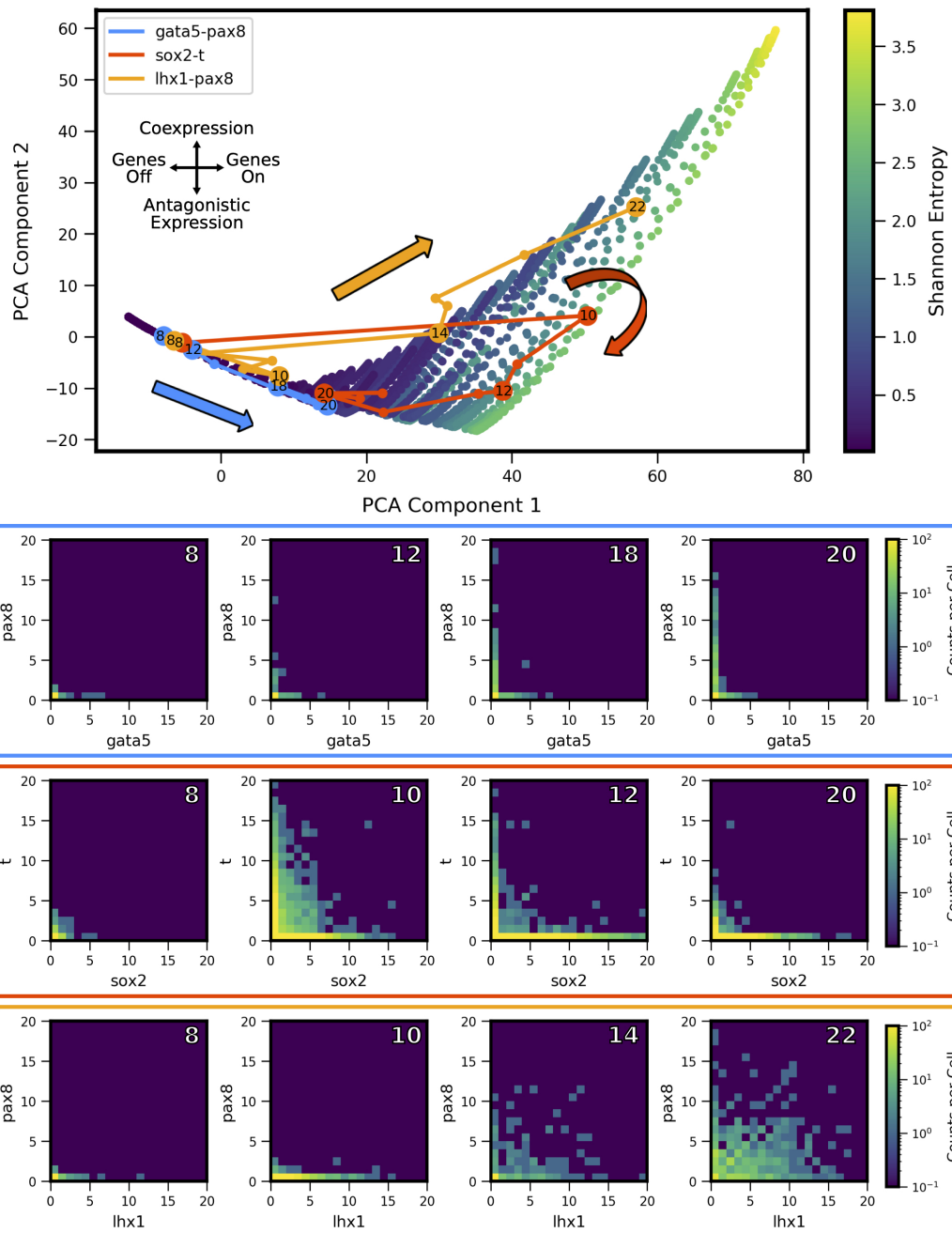


Figure 6. Landscape-shape trajectories of three representative gene pairs from scRNA-seq measurements in *Xenopus tropicalis* embryonic development. (Top) Developmental trajectories of three different gene pairs, plotted in principal component-space. (Bottom) Coexpression quasipotential landscapes extracted from experimental measurements for the three gene pairs at different labeled stages of embryonic development (white numbers indicate developmental stage). The experiment-derived landscapes were trained on the principal components generated from the simulated MISA dataset of Fig. 5. Principal component 1 corresponds to overall level of expression, while component 2 separates antagonistic vs coexpression (see Fig.7). The landscape of *gata5-pax8* (blue) shows increasing antagonistic expression, consistent with movement along the lower left edge of the triangle in PCA shape-space. *Sox2-t* (red) shows high coexpression at stage 10, followed by later antagonistic expression, corresponding to a partial loop through PCA space, consistent with Multilineage Priming behavior. *Lhx1-pax8* (orange) shows consistently increasing coexpression, corresponding to a mostly steady increase in principal components 1 and 2. (Data from Briggs et al. (2018)).

376 interaction is likely not well described by a MISA motif). Finally, *sox2-t* shows a cyclic pattern in the
377 shape subspace, where landscapes move towards hi/hi, and then back towards the antagonistic lo/hi-hi/lo
378 region, landing in a similar area to *gata5-pax8*. Relating these landscape-shape dynamics to the stochastic
379 MISA model parameters suggests that the gene-pairs undergo changes in the relative balance of mutual
380 inhibition versus self-activation as development progresses (see Fig. S1).

381 The experiment-derived developmental trajectories can be further understood by considering the features
382 extracted by individual (by definition orthogonal) PCA components. Visualization of the first three PCA
383 eigenvectors (Fig.7) reveals that the first component (69.3% of covariance across the training set) can be
384 summarized as separating landscapes with more or less expression overall, regardless of whether expression
385 occurs in individual genes or both simultaneously. By contrast, the second component (15.6% of covariance)
386 separates landscapes with coexpression versus antagonistic expression. The third component (6.8% of
387 covariance) distinguishes landscapes with asymmetry between the two genes (subsequent components that
388 describe less of the covariance displayed more complex shapes, and are not shown here). Comparison
389 of the PCA scores versus developmental stage (Fig.7, right) to the experiment-derived landscapes of
390 Fig.6 confirms visually that the PCA components extract the above-described features. For example, all
391 three gene pairs show varying degrees of asymmetry (imbalance in expression levels of the two genes).
392 *Gata5-pax8* shows generally increasing positive amplitude of asymmetry, corresponding to stronger *pax8*
393 expression. At later stages, the other two gene-pairs show asymmetry in the other direction, corresponding
394 to negative amplitude in component 3. *Sox2-t* exhibits a switch in asymmetry between stage 10 ($t > \text{sox2}$)
395 and later stages ($\text{sox2} > t$).

396 Developmental trajectories through the coexpression shape-space were compiled for 1380 gene pairs
397 (putative MLP pairs in *Xenopus tropicalis* identified by Briggs et al. (2018)). By applying the developmental
398 trajectory clustering procedure described in Methods, we found that the trajectories of multiple gene pairs
399 across different lineages display conserved patterns of coexpression dynamics. Twenty-four clusters were
400 identified (see Supplemental Figs. S2 and S3), four of which are shown in Fig. 8; these clusters are chosen
401 as representative of the different types of dynamic patterns obtained. The clusters display a variety of
402 behaviors. For example, the cluster of Fig. 8B shows behavior that is consistent with MLP, i.e., genes are
403 first increasingly coexpressed in single cells, followed by a switch towards antagonistic expression, similar
404 to the cycle in PCA space delineated by *sox2-t* in Fig.6. Surprisingly, we also observed clusters that show
405 “inverted MLP” behavior (Fig.8A) where the genes initially turn on in non-overlapping subsets of cells (i.e.,
406 increasing antagonism), but later show increasing coexpression in single cells. A number of the analyzed
407 gene pairs showed generally antagonistic expression (Fig.8C), reminiscent of *gata5-pax8*. Others showed
408 behavior consistent with the dynamics of MLP (i.e., first coexpression, later antagonistic expression), but
409 with coexpression being only weakly detectable (Fig.8D). The gene pairs represented in these clusters
410 include (but are not limited to) regulators of embryonic development including *zic3*, *hoxc10*, and *neurog1*.
411 The full list of clusters and their associated gene pairs are listed in the Supplementary File 1.

4 DISCUSSION

412 In this work, we comprehensively studied theoretically predicted single-cell gene-gene coexpression
413 landscapes based on a class of stochastic gene regulation models, and applied the theory to analyze
414 two-gene coexpression landscapes from single cell measurements. From a training set of tens of thousands
415 of computed, theoretical landscapes, we identify Principal Components of landscape covariance that serve
416 as simple “fingerprints” of landscape shape and reflect underlying gene-gene interaction dynamics. We
417 then apply the theoretically-derived framework to scRNA-seq data from vertebrate development. In so

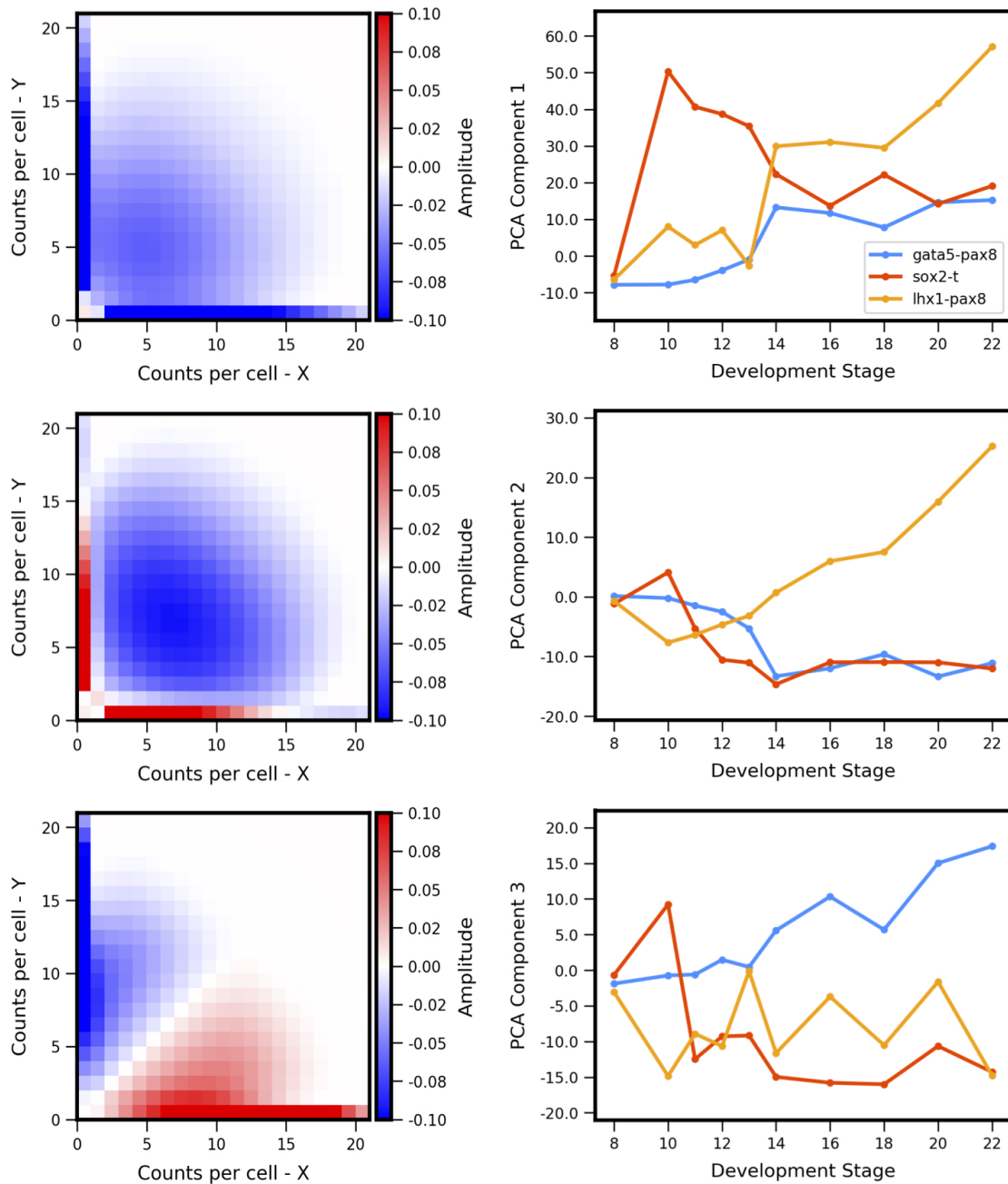


Figure 7. Principal components of landscape shape features. (Left Column) The reshaped PCA principal axes in feature space which represent the maximum variance in the data, specifically which features of the coexpression landscape that each component is accounting for. **(Right Column)** Magnitude or positive/negative value shift in observed variance for the respective component for each gene pair, versus developmental stage. Each component summarizes a landscape shape features: **(Top Row)** The overall amount of gene expression, **(Middle Row)** Antagonistic Expression vs Coexpression of the two genes, and **(Bottom Row)** degree of asymmetric expression between the two genes.

418 doing, we uncover distinctive and novel developmental trajectories of gene-gene coexpression. Specifically,
 419 our framework reveals a nuanced picture of multilineage priming, where the relative balance between
 420 expression of gene pairs simultaneously (in the same cells) versus antagonistically (in different cells)
 421 within a lineage shows complex dynamics during development, for example, revealing that simultaneous
 422 coexpression occurs either earlier or later than antagonism. Based on the results, we propose that the

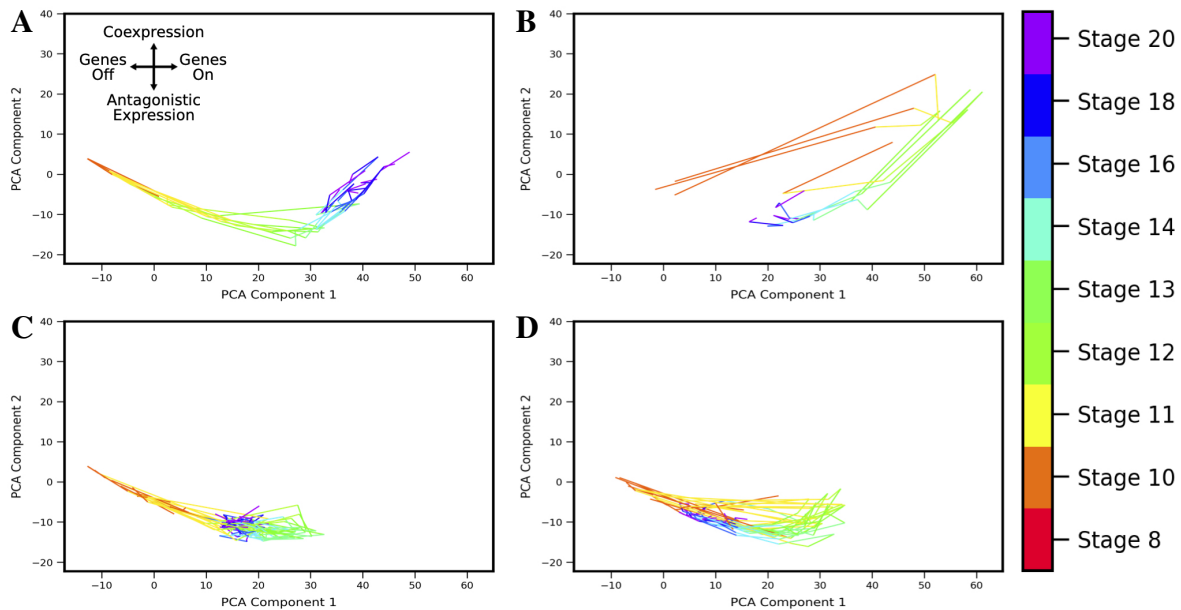


Figure 8. Landscape shape trajectory clustering reveals conserved patterns of gene-pair coexpression dynamics during development. Four representative trajectory clusters showing distinct dynamics are presented (full list of 24 clusters and associated gene pairs in Supplement). Gene pairs in cluster **A** display behavior of an “inverted MLP”: first undergoing increasing antagonistic expression which then switches to increasing coexpression around stage 13. Gene pairs in cluster **B** follow the typical MLP behavior, with highest coexpression taking place around stage 10 followed by antagonistic expression at later stages. Cluster **C** shows consistent antagonistic expression (negative component 2), with nonmonotonic overall expression (a switch-back in component 1 around stage 12). **D** shows cyclic behavior similar to **B**, with highest coexpression at stage 12, but overall expression and relative amount of coexpression is lower.

423 framework developed here can be generalized to other single cell datasets and stochastic network models
424 to analyze the evolution of gene-gene regulatory interactions over the course of development.

425 The theoretical framework applied here—discrete, stochastic reaction kinetic modeling—is well-suited to aid
426 interpretation of single cell measurements: first, because it inherently captures cell population heterogeneity
427 and second, because of the direct correspondence between the computed quantities (e.g., probability to find
428 a given number of mRNAs in a cell) and experimentally-measured transcript counts in scRNA-seq. The
429 theoretical models can partially reproduce true cell population heterogeneity, but also neglect many sources
430 of noise, both biological and technical. We employ models that treat intrinsic noise but neglect sources of
431 persistent cell-to-cell variability (i.e., extrinsic noise) (Swain et al. (2002)), which is known to contribute to
432 noise in gene expression. For example, one source of extrinsic noise would be asynchronicity between
433 cells, where individual cells might be at different stages of progression in development. Here, we opted to
434 use a relatively simplistic model framework (i.e., no additional noise assumptions beyond intrinsic noise of
435 biomolecular interactions, relatively few reactions describing molecular mechanisms of gene regulation,
436 etc.) to minimize the number of model parameters while still enabling study of a variety of “rules” for
437 gene regulatory logic. The framework presented here could be expanded in the future by integration of
438 additional types of mechanistic assumptions and noise sources in the stochastic models.

439 The models also neglect technical noise/measurement errors arising from experiments (Grün et al. (2014)).
440 For example, scRNA-seq measurements face a well-known technical issue of drop-outs (Kharchenko et al.
441 (2014)), which we have not included in our modeling. Future efforts may improve the presented modeling
442 framework by inclusion of these additional sources of noise, or by additional data-processing steps for

443 imputation of missing datapoints (Gong et al. (2018)). However, such an approach would also present
444 challenges by necessarily introducing additional assumptions about cell population heterogeneity, which is
445 still not fully understood. Given the danger of false signals (Andrews and Hemberg (2019)), we opted here
446 to utilize minimal data processing in comparing our theoretical results to a public dataset. We also note that
447 the discrete stochastic modeling framework advanced in this work has potential to shed new light on the
448 drop-outs issue: a relatively large proportion of “zeros” arises naturally from discrete stochastic models,
449 depending on the regulatory interactions among genes, suggesting that perhaps biological variability plays a
450 larger role in producing dropouts than has previously been supposed. Overall, despite the lack of additional
451 biological/technical noise sources in our models, we note that our computed landscapes qualitatively
452 reproduce the noise characteristics of the scRNA-seq measurements, in that they showed similarly broad
453 distributions of coexpression. Thus we conclude that the simplistic models employed here are sufficient
454 for the current application, which focused on characterization of coexpression landscape shape and its
455 evolution in development, but we also foresee that incorporation of additional noise sources in the model
456 might improve the practical utility of our proposed coexpression-shape-based analysis.

457 We focused here on two-gene models and pairwise interactions, because (1) certain gene-pairs are known
458 to play a critical role in development (Graf and Enver (2009)) (2) the edges (pairwise interactions) are
459 the elemental units or building blocks of larger regulatory networks. However, the focus on pairwise
460 interactions has potential drawbacks: it does not elucidate how gene-pair interactions are modified when
461 embedded in a larger network. In the same vein, it does not differentiate between direct or indirect
462 interactions between genes (e.g., by direct transcriptional regulation versus molecular intermediaries). In
463 principle, the framework presented here could be expanded to treat “3-body” (or higher order) interactions
464 among genes, though this presents several computational challenges. For example, solution of the CME
465 becomes intractable already for 3-gene networks, such that advanced approximation methods (Zhang and
466 Wolynes (2014)) or more costly simulations (Tse et al. (2018)) become necessary. Nevertheless, expansion
467 of the approach to higher-order interactions is feasible, and recent work has revealed how such an approach
468 might proceed, for example, by incorporating developments in multivariate information measures (Chan
469 et al. (2017)).

CONFLICT OF INTEREST STATEMENT

470 The authors declare that the research was conducted in the absence of any commercial or financial
471 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

472 Conceptualization and study design, CG and ER; Coding, CG; Analysis, CG, HR, and ER; Visualization,
473 CG and HR; Writing, review, and editing, CG, HR, ER. All authors read and approved the final version of
474 the manuscript.

FUNDING

475 This work was (partially) supported by a NSF grant DMS 1763272 and a grant from the Simons
476 Foundation (594598, QN). CPG was provided support by a GAANN fellowship funded by the U.S.
477 Department of Education.

REFERENCES

- 478 Al-Radhawi, M. A., Vecchio, D. D., and Sontag, E. D. (2019). Multi-modality in gene regulatory networks
479 with slow promoter kinetics. *PLOS Computational Biology* 15, e1006784. doi:10.1371/journal.pcbi.
480 1006784
- 481 Andrews, T. S. and Hemberg, M. (2019). False signals induced by single-cell imputation. *F1000Research*
482 7, 1740. doi:10.12688/f1000research.16613.2
- 483 Arkin, A., Ross, J., and McAdams, H. H. (1998). Stochastic kinetic analysis of developmental pathway
484 bifurcation in phage lambda-infected *Escherichia coli* cells. *Genetics* 149, 1633–1648
- 485 Bhattacharya, S., Zhang, Q., and Andersen, M. E. (2011). A deterministic map of Waddington’s epigenetic
486 landscape for cell fate specification. *BMC Systems Biology* 5, 85. doi:10.1186/1752-0509-5-85
- 487 Briggs, J. A., Weinreb, C., Wagner, D. E., Megason, S., Peshkin, L., Kirschner, M. W., et al. (2018).
488 The dynamics of gene expression in vertebrate embryogenesis at single-cell resolution. *Science* 360,
489 eaar5780. doi:10.1126/science.aar5780
- 490 Cao, Y. and Liang, J. (2008). Optimal enumeration of state space of finitely buffered stochastic molecular
491 networks and exact computation of steady state landscape probability. *BMC Systems Biology* 2, 30.
492 doi:10.1186/1752-0509-2-30
- 493 Chan, T. E., Stumpf, M. P., and Babbie, A. C. (2017). Gene Regulatory Network Inference from Single-Cell
494 Data Using Multivariate Information Measures. *Cell Systems* 5, 251–267.e3. doi:10.1016/j.cels.2017.08.
495 014
- 496 Chu, B. K., Tse, M. J., Sato, R. R., and Read, E. L. (2017). Markov State Models of gene regulatory
497 networks. *BMC Systems Biology* 11, 14. doi:10.1186/s12918-017-0394-4
- 498 Dibaeinia, P. and Sinha, S. (2019). *A single-cell expression simulator guided by gene regulatory networks*.
499 preprint, Bioinformatics. doi:10.1101/716811
- 500 Feng, H. and Wang, J. (2012). A new mechanism of stem cell differentiation through slow
501 binding/unbinding of regulators to genes. *Scientific Reports* 2, 550. doi:10.1038/srep00550
- 502 Geertz, M., Shore, D., and Maerkl, S. J. (2012). Massively parallel measurements of molecular interaction
503 kinetics on a microfluidic platform. *Proceedings of the National Academy of Sciences* 109, 16540–16545.
504 doi:10.1073/pnas.1206011109
- 505 Gillespie, D. T. (1977). Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical*
506 *Chemistry* 81, 2340–2361. doi:10.1021/j100540a008
- 507 Gong, W., Kwak, I.-Y., Pota, P., Koyano-Nakagawa, N., and Garry, D. J. (2018). DrImpute: imputing
508 dropout events in single cell RNA sequencing data. *BMC Bioinformatics* 19, 220. doi:10.1186/
509 s12859-018-2226-y
- 510 Graf, T. and Enver, T. (2009). Forcing cells to change lineages. *Nature* 462, 587–594. doi:10.1038/
511 nature08533
- 512 Grün, D., Kester, L., and van Oudenaarden, A. (2014). Validation of noise models for single-cell
513 transcriptomics. *Nature Methods* 11, 637–640. doi:10.1038/nmeth.2930
- 514 Hathaway, N., Bell, O., Hodges, C., Miller, E., Neel, D., and Crabtree, G. (2012). Dynamics and Memory
515 of Heterochromatin in Living Cells. *Cell* 149, 1447–1460. doi:10.1016/j.cell.2012.03.052
- 516 Huang, S. (2012). The molecular and mathematical basis of Waddington’s epigenetic landscape: A
517 framework for post-Darwinian biology? *BioEssays* 34, 149–157. doi:10.1002/bies.201100031
- 518 Huang, S. (2013). Hybrid T-Helper Cells: Stabilizing the Moderate Center in a Polarized System. *PLoS*
519 *Biology* 11, e1001632. doi:10.1371/journal.pbio.1001632
- 520 Huang, S., Guo, Y.-P., May, G., and Enver, T. (2007). Bifurcation dynamics in lineage-commitment in
521 bipotent progenitor cells. *Developmental Biology* 305, 695–713. doi:10.1016/j.ydbio.2007.02.036

- 522 Jin, S., MacLean, A. L., Peng, T., and Nie, Q. (2018). scEpath: energy landscape-based inference of
523 transition probabilities and cellular trajectories from single-cell transcriptomic data. *Bioinformatics* 34,
524 2077–2086. doi:10.1093/bioinformatics/bty058
- 525 Kauffman, S. (1969). Metabolic stability and epigenesis in randomly constructed genetic nets. *Journal of*
526 *Theoretical Biology* 22, 437–467. doi:10.1016/0022-5193(69)90015-0
- 527 Kepler, T. B. and Elston, T. C. (2001). Stochasticity in Transcriptional Regulation: Origins, Consequences,
528 and Mathematical Representations. *Biophysical Journal* 81, 3116–3136. doi:10.1016/S0006-3495(01)
529 75949-8
- 530 Kharchenko, P. V., Silberstein, L., and Scadden, D. T. (2014). Bayesian approach to single-cell differential
531 expression analysis. *Nature Methods* 11, 740–742. doi:10.1038/nmeth.2967
- 532 Lam, S. K., Pitrou, A., and Seibert, S. (2015). Numba: A llvm-based python jit compiler. In *Proceedings*
533 *of the Second Workshop on the LLVM Compiler Infrastructure in HPC* (New York, NY, USA: ACM),
534 LLVM '15, 7:1–7:6. doi:10.1145/2833157.2833162
- 535 MacArthur, B. D., Ma'ayan, A., and Lemischka, I. R. (2009). Systems biology of stem cell fate and cellular
536 reprogramming. *Nature Reviews Molecular Cell Biology* 10, 672–681. doi:10.1038/nrm2766
- 537 Mariani, L., Schulz, E. G., Lexberg, M. H., Helmstetter, C., Radbruch, A., Löhning, M., et al. (2010).
538 Short-term memory in gene induction reveals the regulatory principle behind stochastic IL-4 expression.
539 *Molecular Systems Biology* 6, 359. doi:10.1038/msb.2010.13
- 540 McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python*
541 *in Science Conference*, eds. S. van der Walt and J. Millman. 51 – 56
- 542 Micheelsen, M. A., Mitarai, N., Sneppen, K., and Dodd, I. B. (2010). Theory for the stability and regulation
543 of epigenetic landscapes. *Physical Biology* 7, 026010. doi:10.1088/1478-3975/7/2/026010
- 544 Mojtahedi, M., Skupin, A., Zhou, J., Castaño, I. G., Leong-Quong, R. Y. Y., Chang, H., et al. (2016). Cell
545 Fate Decision as High-Dimensional Critical State Transition. *PLOS Biology* 14, e2000640. doi:10.1371/
546 journal.pbio.2000640
- 547 Nimmo, R. A., May, G. E., and Enver, T. (2015). Primed and ready: understanding lineage commitment
548 through single cell analysis. *Trends in Cell Biology* 25, 459–467. doi:10.1016/j.tcb.2015.04.004
- 549 Olsson, A., Venkatasubramanian, M., Chaudhri, V. K., Aronow, B. J., Salomonis, N., Singh, H., et al.
550 (2016). Single-cell analysis of mixed-lineage states leading to a binary cell fate choice. *Nature* 537,
551 698–702. doi:10.1038/nature19348
- 552 Peccoud, J. and Ycart, B. (1995). Markovian Modeling of Gene-Product Synthesis. *Theoretical Population*
553 *Biology* 48, 222–234. doi:10.1006/tpbi.1995.1027
- 554 Pedraza, J. M. and Paulsson, J. (2008). Effects of Molecular Memory and Bursting on Fluctuations in Gene
555 Expression. *Science* 319, 339–343. doi:10.1126/science.1144331
- 556 Sasai, M., Kawabata, Y., Makishi, K., Itoh, K., and Terada, T. P. (2013). Time Scales in Epigenetic
557 Dynamics and Phenotypic Heterogeneity of Embryonic Stem Cells. *PLoS Computational Biology* 9,
558 e1003380. doi:10.1371/journal.pcbi.1003380
- 559 Sasai, M. and Wolynes, P. G. (2003). Stochastic gene expression as a many-body problem. *Proceedings of*
560 *the National Academy of Sciences* 100, 2374–2379. doi:10.1073/pnas.2627987100
- 561 Schaffter, T., Marbach, D., and Floreano, D. (2011). GeneNetWeaver: in silico benchmark generation
562 and performance profiling of network inference methods. *Bioinformatics* 27, 2263–2270. doi:10.1093/
563 bioinformatics/btr373
- 564 Schwanhäusser, B., Busse, D., Li, N., Dittmar, G., Schuchhardt, J., Wolf, J., et al. (2011). Global
565 quantification of mammalian gene expression control. *Nature* 473, 337–342. doi:10.1038/nature10098

- 566 Swain, P. S., Elowitz, M. B., and Siggia, E. D. (2002). Intrinsic and extrinsic contributions to stochasticity
567 in gene expression. *Proceedings of the National Academy of Sciences* 99, 12795–12800. doi:10.1073/
568 pnas.162041399
- 569 Thattai, M. and van Oudenaarden, A. (2001). Intrinsic noise in gene regulatory networks. *Proceedings of*
570 *the National Academy of Sciences* 98, 8614–8619. doi:10.1073/pnas.151588598
- 571 Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., et al. (2014). The dynamics
572 and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature*
573 *Biotechnology* 32, 381–386. doi:10.1038/nbt.2859
- 574 Tse, M., Chu, B., Roy, M., and Read, E. (2015). DNA-Binding Kinetics Determines the Mechanism of
575 Noise-Induced Switching in Gene Networks. *Biophysical Journal* 109, 1746–1757. doi:10.1016/j.bpj.
576 2015.08.035
- 577 Tse, M. J., Chu, B. K., Gallivan, C. P., and Read, E. L. (2018). Rare-event sampling of epigenetic
578 landscapes and phenotype transitions. *PLOS Computational Biology* 14, e1006336. doi:10.1371/journal.
579 pcbi.1006336
- 580 van der Walt, S., Colbert, S. C., and Varoquaux, G. (2011). The numpy array: A structure for efficient
581 numerical computation. *Computing in Science Engineering* 13, 22–30. doi:10.1109/MCSE.2011.37
- 582 Waddington, C. H. (2014). *The strategy of the genes*. OCLC: 861212016
- 583 Wang, J., Zhang, K., Xu, L., and Wang, E. (2011). Quantifying the Waddington landscape and biological
584 paths for development and differentiation. *Proceedings of the National Academy of Sciences* 108,
585 8257–8262. doi:10.1073/pnas.1017017108
- 586 Wardle, F. C. and Smith, J. C. (2004). Refinement of gene expression patterns in the early *Xenopus* embryo.
587 *Development* 131, 4687–4696. doi:10.1242/dev.01340
- 588 Zhang, B. and Wolynes, P. G. (2014). Stem cell differentiation as a many-body problem. *Proceedings of*
589 *the National Academy of Sciences* 111, 10185–10190. doi:10.1073/pnas.1408561111
- 590 Zhou, J. X., Samal, A., d’Hérouël, A. F., Price, N. D., and Huang, S. (2016). Relative stability of network
591 states in Boolean network models of gene regulation in development. *Biosystems* 142-143, 15–24.
592 doi:10.1016/j.biosystems.2016.03.002