

DREISS: dynamics of gene expression driven by external and internal regulatory modules based on state space model

Daifeng Wang^{1,2}, Fei He⁴, Sergei Maslov⁴, Mark Gerstein^{1,2,3*}

¹Program in Computational Biology and Bioinformatics; ²Department of Molecular Biophysics and Biochemistry;

³Department of Computer Science, Yale University, New Haven, CT, USA. ⁴Brookhaven National Laboratory, Upton, NY,

USA. *Correspondence to: pi@gersteinlab.org

ABSTRACT

Motivation: Gene expression is controlled by combinatorial effects of gene regulatory factors from different biological subsystems such as general transcription factors, cellular growth factors and microRNAs. A subsystem's gene expression may be controlled by its internal regulatory factors, exclusively, or by other external subsystems, or by both. It is thus useful to distinguish the degree to which a subsystem is regulated internally or externally; e.g., to understand how external regulatory factors affect the expression of conserved genes during evolution.

Results: We developed a computational method, DREISS for dynamics of gene expression driven by external and internal regulatory modules based on state space model to help dissect the effects of different regulatory subsystems on gene expression. Given a subsystem, the "state" and "control" in the model refer to its own (internal) and another subsystem's (external) gene expression levels. The state at a time is determined by the state and control at previous time. DREISS integrates the dimensionality reduction for combating the limited time samples, and identifies the canonical temporal expression trajectories (e.g., degradation, growth, damped oscillation) representing the regulatory effects from various subsystems.

We applied DREISS to the time-series gene expression datasets of worm (*C. elegans*) and fly (*D. melanogaster*) during their embryonic development, to demonstrate its capabilities for studying the regulatory effects of evolutionary conserved vs. divergent transcription factors across distant species. We analyzed the expression dynamics of the conserved, orthologous genes (orthologs), seeing the degree to which these can be accounted for by orthologous (internal) versus species-specific (external) transcription factors (TFs). We found that between the two species, the canonical trajectories of orthologs expression driven by orthologous TFs are more similar to each other than those driven by species-specific ones. This is particularly true for genes with evolutionarily ancient functions (e.g. the ribosome), in contrast to those with more recently evolved functions (e.g., cell-cell communication). This implies

*To whom correspondence should be addressed.

ASAW
TLWS

- Daifeng Wang 5/4/15 11:55 AM
Deleted: A
- Daifeng Wang 5/4/15 11:55 AM
Deleted: network can be modularized into multiple interconnected
- Daifeng Wang 5/4/15 11:55 AM
Deleted: evolutionarily conserved
- Daifeng Wang 5/4/15 11:55 AM
Deleted: species-specific ones
- Daifeng Wang 5/4/15 11:55 AM
Deleted: the species-specific (
- Daifeng Wang 5/4/15 11:55 AM
Deleted:)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: conserved gene
- Daifeng Wang 5/4/15 11:55 AM
Deleted: the gene expression levels of it
- Daifeng Wang 5/4/15 11:55 AM
Deleted: subsystem
- Daifeng Wang 5/4/15 11:55 AM
Deleted: .
- Daifeng Wang 5/4/15 11:55 AM
Deleted: To illustrate DREISS, we
- Daifeng Wang 5/4/15 11:55 AM
Deleted: it
- Daifeng Wang 5/4/15 11:55 AM
Deleted: from the modENCODE project
- Daifeng Wang 5/4/15 11:55 AM
Deleted: different
- Daifeng Wang 5/4/15 11:55 AM
Deleted: subsystems
- Daifeng Wang 5/4/15 11:55 AM
Deleted: transcription factors (TFs)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: ones

that despite striking morphological differences, some basic embryonic-developmental processes are still tightly under the control of ancient regulation for the similar canonical trajectories of worm-fly orthologs driven by the orthologous TFs.

Availability and implementation: We make DREISS available as general-purpose tool (github.com/gersteinlab/dreiss) to compare the effects from different types of biological regulatory sub-systems in future.

DREISS
GL.

1 INTRODUCTION

Gene regulatory networks (GRNs) systematically control gene expression dynamics. Those networks are highly modular, and consist of various sub-networks. Each sub-network includes a number of regulatory factors representing a subsystem to drive particular gene regulatory functions (Kim and Tidor, 2003; Vilar, 2006). The subsystems interact with one another, and work together to carry out the entire gene regulatory function. For example, the gene expression in embryogenesis is controlled by the combinatorial effects of various regulatory subsystems comprising complex evolutionary GRNs (Peter and Davidson, 2011). These regulatory subsystems drive very diverse developmental functions, from the highly conserved (e.g. DNA replication) to the species-specific (e.g. body segmentation). For example, the orthologous genes can be regulated by both orthologous and species-specific transcription factors (TFs). The orthologous TFs constitute an “internal” regulatory network, while the species-specific TFs constitute an “external” regulatory network. Unfortunately, existing experimental gene expression data cannot decouple the expression components that are driven by different subsystems. Thus, we need computational methods to assess the contribution from each factor or subsystem from the gene expression data. In this study, we propose a novel computational method, DREISS - Decomposition of gene Regulatory network into External and Internal components based on State Space models. We identify temporal gene expression dynamic patterns for evolutionarily conserved genes during embryonic development, as driven by conserved and species-specific regulatory subsystems. This advances our current understanding of GRNs in evolution as well as differentiation during development.

Developmental GRNs control gene expression to determine developmental processes. These GRNs have been evolved, making it difficult to understand their regulatory mechanisms at the system level. Hence, one typically compares developmental gene expression across species to infer activities of developmental GRNs. For example, embryogenesis provides a platform to study the evolution of gene expression

(X)

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [1]

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [2]

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [3]

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [4]

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [5]

Daifeng Wang 5/4/15 11:55 AM
Deleted: Developmental gene regulatory networks (GRNs)

PROGRAMS

DEF

between different species. Recent work has showed that significant biological insight can be gained by cross-species comparisons of the expression profiles during embryogenesis for worms (Levin, et al., 2012), flies (Kalinka, et al., 2010), frogs (Yanai, et al., 2011) and several other vertebrates (Irie and Kuratani, 2011). It was found that the orthologous genes have minimal temporal expression divergence during the phylotypic stage, a middle phase at embryonic development across species within the same phylum. These patterns are often characterized as hourglass patterns (Casici, 2011). In addition, the conserved hourglass patterns were observed within a single species while comparing the developmental gene expression data across distant species, such as worm and fly (Gerstein, et al., 2014); i.e., the expression divergence among evolutionarily conserved genes become minimal during the phylotypic stage in both worm and fly. However, how the orthologous genes in each species eventually contribute to their species-specific phenotypes is less studied due to the lack of appropriate computational approaches. Our method is able to identify temporal gene expression dynamic patterns for evolutionarily conserved genes during embryonic development, as driven by conserved and species-specific regulatory subsystems. This advances our current understanding of GRNs in evolution and differentiation during development.

THIS IS WHAT WE WANT.

NOV BIO EXP

The state-space model has been widely used in engineering (Brogan, 1991) and analyzing gene expression dynamics (Bansal, et al., 2006; Huang and Ingber, 2006; Rangel, et al., 2004). It models the dynamical system output as a function of both current internal system state and external input signal. Similarly, for the orthologous genes, their expression at the next developmental stage can be predicted from their expression (internal) and species-specific regulatory factors (external) at the current stage. Unlike previous studies that calculates the expression correlation between individual genes, the state-space model predicts the temporal causal relationships at the system level; i.e., the state at a time is determined by the state and external input at previous time. The previous work applied the state-space model to study the gene expression dynamics focusing on small-scale systems, and did not explore the analytic dynamic characteristics of the inferred state-space models. The complex and large-scale biological datasets, especially temporal gene expression data, are very noisy, and high dimensional (i.e., the number of genes is much greater than the number of time samples), thereby preventing an accurate estimation of the state-space model's parameters. The dimensionality reduction techniques have thus been used to project high-dimensional genes to low-dimensional meta-genes (i.e., the selected features representing de-noised and systematic expression patterns (Chu, et al., 1998; Kim and Tidor, 2003; Saeys, et al., 2007)) as well as

- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Levin, et al., 2012)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Kalinka, et al., 2010)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Yanai, et al., 2011)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Irie and Kuratani, 2011)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: . It was found that the hourglass patterns of orthologous genes have minimal expression divergence across species within the same phylum during the phylotypic stage at embryonic development.
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Gerstein, et al., 2014)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: divergence
- Daifeng Wang 5/4/15 11:55 AM
Deleted: ... (1)
- Daifeng Wang 5/4/15 11:55 AM
Moved up [1]: The subsystems interact with one another, and work together to carry out the entire gene regulatory function. For example, the gene expression in embryogenesis is controlled by the combinatorial effects of various regulatory subsystems comprising complex evolutionary GRNs
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Peter and Davidson, 2011). These regulatory subsystems have driven
- Daifeng Wang 5/4/15 11:55 AM
Moved up [2]: very diverse developmental functions, from the highly conserved (e.g. DNA replication) to the species-specific (e.g. body segmentation).
- Daifeng Wang 5/4/15 11:55 AM
Moved up [4]: Unfortunately, existir ... (2)
- Daifeng Wang 5/4/15 11:55 AM
Moved up [3]: For example, the orth ... (3)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: Thus, we need computation: ... (4)
- Daifeng Wang 5/4/15 11:55 AM
Moved up [5]: the contribution from ... (5)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Brogan, 1991)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Bansal, et al., 2006; Huang ... (6)
- Daifeng Wang 5/4/15 11:55 AM
Deleted: work
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (Chu, et al., 1998; Kim and ... (7)

the principal dynamic patterns for those meta-genes (Wang, et al., 2012; Wang, et al., 2012). DREISS applies the dimensionality reduction to the gene expression data, develops an effective state-space model for their meta-genes, and identifies a group of canonical temporal expression trajectories representing the dynamic patterns driven by effective conserved and species-specific meta-gene regulatory networks according to the model's analytic characteristics. These dynamic patterns reveal temporal gene expression components that are controlled by conserved or species-specific GRNs.

DREISS IS GEN. IT COULD BE APP TO HERE AS AN ILL.

We applied DREISS to the gene expression data during embryonic development for two model organisms, worm (*Caenorhabditis elegans*) and fly (*Drosophila melanogaster*). In both species, we identified the expression patterns of worm-fly orthologs driven by the conserved regulatory network consisting of the worm-fly TFs (i.e., the conserved regulatory subsystems between two species), as well as the worm/fly-specific regulatory network consisting of non-orthologous TFs (i.e., the species-specific regulatory subsystem). Our results reveal that, in addition to executing conserved developmental functions between worm and fly, their orthologous genes are also regulated by species-specific TFs to involve in species-specific developmental processes. In summary, DREISS provides a framework to analyze distantly and closely related species allowing for better understanding the gene regulatory mechanisms during development.

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Wang, et al., 2012; Wang, et al., 2012)

2 MATERIALS AND METHODS

DRNL?

DREISS consists of five major steps (Figure 1):

Step 1: DREISS models temporal gene expression dynamics using state-space models in control theory.

The "state" refers to the expressions for a large group of genes of interest, such as the worm-fly orthologous genes investigated here. The "control" refers to any other group of genes that contribute to gene expressions of the "state", such as the species-specific TFs contributed to control orthologous gene expression.

Step 2: Due to the limited number of temporal samples in gene expression experiments, we do not have enough data to estimate the parameters of the state-space models that capture interactions among hundreds of genes. Therefore, DREISS projects high-dimensional gene expression space to lower-dimensional meta-gene expression spaces using dimensionality reduction techniques.

Step 3: DREISS then derives the effective state-space models for meta-genes so that model parameters can be estimated.

Daifeng Wang 5/4/15 11:55 AM
Deleted: orthologous transcription factors (TFs)

Daifeng Wang 5/4/15 11:55 AM
Deleted: been

Daifeng Wang 5/4/15 11:55 AM
Moved down [6]: Moreover, by applying DREISS to human genes, such analyses may further our understanding of human development. For example, we may gain valuable insights into the development of the human brain and pathologies such as Alzheimer's disease.

Daifeng Wang 5/4/15 11:55 AM
Deleted: TF studied here.

mm = BBSR

Step 4: DREISS then identifies the meta-gene expression dynamic patterns; i.e., canonical temporal expression trajectories driven by “state” (internal) and by “control” (external) based on the analytic solutions to estimated models.

Step 5: We calculate the coefficients of genes for the dynamic patterns of linear transformations between genes and meta-genes. DREISS also allows us to compare the dynamic expression patterns of multiple datasets with samples taken at different times. We describe DREISS in detail in each step as follows.

2.1 State-space models for temporal gene expression dynamics

A gene regulatory network is made up of various subsystems (Kim and Tidor, 2003; Vilar, 2006). These subsystems work together to execute the regulatory functions. Given a group of N_1 genes in a subsystem, their gene expression levels (X) are not only controlled by internal interactions among X , but also affected by the regulatory factors from other subsystems outside X (external regulations, the U group in this paper). For example, we can consider the worm-fly orthologous genes as the X group. The worm-fly orthologous TFs from the X group are the internal regulatory factors, and non-orthologous TFs such as worm- or fly- specific TFs are the external regulatory factors to the X group, namely Group U . Another example is that if the X group consists of protein-coding genes, then the transcription factors and microRNAs can be internal and external regulatory factors, respectively. Both internal and external regulatory factors control gene expressions in dynamic ways (i.e., their regulatory signals at the current time will affect gene expressions at future times). Thus, the regulatory mechanisms for the gene expressions form a control system. In this study, we used a state-space model (linear first-order difference equations, Figure 2A), which has been commonly used in control engineering, to formulate temporal gene expression dynamics for the gene group X (comprising N_1 genes) with external regulations from the gene group U (comprising N_2 genes) at time points $1, 2, \dots, T$ as follows:

$$X_{t+1} = AX_t + BU_t \quad (1)$$

, where the vector $X_t \in \mathbb{R}^{N_1 \times 1}$, the “state”, includes N_1 gene expression levels at time t in group X , and the vector $U_t \in \mathbb{R}^{N_2 \times 1}$, the “input or control”, includes N_2 gene expression levels at time t in group U .

The system matrix $A \in \mathbb{R}^{N_1 \times N_1}$ captures internal causal interactions among genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of A , A_{ij} describes the contribution from the j^{th} gene expression at time t to the i^{th} gene expression at the next time $t+1$). The control matrix $B \in \mathbb{R}^{N_1 \times N_2}$ captures external causal regulations from the genes in U to genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of B , B_{ij} describes the contribution from the j^{th} gene expression

EARLIER

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Kim and Tidor, 2003; Vilar, 2006)

Daifeng Wang 5/4/15 11:55 AM
Deleted: .

Daifeng Wang 5/4/15 11:55 AM
Deleted: in our study,

Daifeng Wang 5/4/15 11:55 AM
Deleted: .

Daifeng Wang 5/4/15 11:55 AM
Deleted: 2

GEN TO AND?

IT INSTANTIATES A NETWORK 5

Vertical scribble on the right margin.

Vertical scribble on the right margin.

in U at time t to the i^{th} gene expression in X at the next time $t+1$). According to the state space model (1), the gene expression dynamics in X is determined by the system matrix A and the control matrix B .

2.2 Dimensionality reduction from genes to meta-genes

The temporal gene expression experiments normally have limited time samples (for example, there may only be a dozen time points), which are far less than the [time](#) samples needed to estimate the large matrices A and B when X and U have hundreds or thousands of genes. Thus, we project high dimensional temporal gene expressions to much lower dimensional meta-gene expression levels using dimensionality reduction (Figure [2B](#)). Those meta-gene expression levels should capture original gene expression patterns, such as the ones having the greatest degree of co-variation. We calculate the meta-gene expression levels as follows:

$$\tilde{X}_t = W_X^* X_t; \tilde{U}_t = W_U^* U_t \quad (2)$$

, where $\tilde{X}_t \in \mathfrak{R}^{M_1 \times 1}$, the “meta-gene state”, includes M_1 ($\ll N_1$ and $< T$) meta-gene expression levels; i.e., the values of first M_1 singular vectors from singular value decomposition (SVD) of matrix $[X_1 X_2 \dots X_T]$ at time t in group X ; the vector $\tilde{U}_t \in \mathfrak{R}^{M_2 \times 1}$, the “meta-gene input or control”, includes M_2 ($\ll N_2$ and $< T$) meta-gene expression levels (i.e., the values of the first M_2 singular vectors from SVD of matrix $[U_1 U_2 \dots U_T]$ at time t in group U ; $W_X \in \mathfrak{R}^{N_1 \times M_1}$ is the linear projection matrix of SVD from M_1 meta-gene expression space to N_1 gene expression space in X , $W_U \in \mathfrak{R}^{N_2 \times M_2}$ is the linear projection matrix of SVD from M_2 meta-gene expression space to N_2 gene expression space in U), and $(.)^*$ is a pseudo-inverse operation; i.e., $W^* W = I$, where I is the identity matrix.

2.3 Estimation of effective state-space model for meta-gene expression dynamics

Next, we can obtain the effective state-space model for meta-genes using linear projections W_X and W_U between genes and meta-genes as follows (Figure [2C](#)). By replacing (1) using (2), we obtain that

$$W_X \tilde{X}_{t+1} = A W_X \tilde{X}_t + B W_U \tilde{U}_t \quad (3)$$

, and by multiplying the pseudo-inverse of W_X , $W_X^* \in \mathfrak{R}^{M_1 \times N_1}$ s.t. $W_X^* W_X = I$ where I is an identity matrix, at both sides of (3),

$$\tilde{X}_{t+1} = \underbrace{W_X^* A W_X}_{\tilde{A}} \tilde{X}_t + \underbrace{W_X^* B W_U}_{\tilde{B}} \tilde{U}_t \Rightarrow \tilde{X}_{t+1} = \tilde{A} \tilde{X}_t + \tilde{B} \tilde{U}_t \quad (4)$$

Daifeng Wang 5/4/15 11:55 AM

Deleted: are

Daifeng Wang 5/4/15 11:55 AM

Deleted: 3

Daifeng Wang 5/4/15 11:55 AM

Deleted: 4

, where the effective meta-gene system matrix $\tilde{A} = W_X^* A W_X \in \mathfrak{R}^{M_1 \times M_1}$ captures internal causal interactions among meta-genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of \tilde{A} (\tilde{A}_{ij}) describes the contribution from the j^{th} meta-gene expression at time t to i^{th} meta-gene expression at next time $t+1$), and the effective control matrix $\tilde{B} = W_X^* B W_U \in \mathfrak{R}^{M_1 \times M_2}$ captures external causal regulations from meta-genes in U to meta-genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of \tilde{B} , \tilde{B}_{ij} describes the contribution from the j^{th} meta-gene expression in U at time t to i^{th} meta-gene expression in X at next time $t+1$). Equation (4) describes the effective state space model for the meta-genes in X , whose expression dynamics is determined by \tilde{A} and \tilde{B} . Because the meta-gene dimension, M_1 (M_2) is less than T , and much less than N_1 (N_2), we can estimate \tilde{A} and \tilde{B} as follows.

We rewrite Equation (4) as a matrix product on the right side:

$$\tilde{X}_{t+1} = \tilde{A}\tilde{X}_t + \tilde{B}\tilde{U}_t = \begin{bmatrix} \tilde{A} & \tilde{B} \end{bmatrix} \begin{bmatrix} \tilde{X}_t \\ \tilde{U}_t \end{bmatrix}. \quad (5)$$

By applying Equation (5) to time points, 2, 3, ..., T , we then obtain that

$$\underbrace{\begin{bmatrix} \tilde{X}_2 & \tilde{X}_3 & \cdots & \tilde{X}_T \end{bmatrix}}_Z = \begin{bmatrix} \tilde{A} & \tilde{B} \end{bmatrix} \underbrace{\begin{bmatrix} \tilde{X}_1 & \tilde{X}_2 & \cdots & \tilde{X}_{T-1} \\ \tilde{U}_1 & \tilde{U}_2 & \cdots & \tilde{U}_{T-1} \end{bmatrix}}_Y \quad (6)$$

, where $Z \in \mathfrak{R}^{M_1 \times (T-1)}$ and $Y \in \mathfrak{R}^{(M_1+M_2) \times (T-1)}$.

The effective internal system matrix \tilde{A} and external control matrix \tilde{B} can be estimated by

$$\begin{bmatrix} \tilde{A} & \tilde{B} \end{bmatrix} = ZY^* \quad (7)$$

, where $Y^* \in \mathfrak{R}^{(T-1) \times (M_1+M_2)}$ is the pseudo-inverse of Y ; i.e.

$YY^* = I$, with $M_1 < N_1, M_2 < N_2, M_1 + M_2 < T, t = 1, 2, \dots, T$.

2.4 Identification of internally and externally driven principal expression dynamic patterns of meta-genes (canonical temporal expression trajectories)

HOW DO YOU KNOW ENF HAS DATA

Daifeng Wang 5/4/15 11:55 AM

Deleted: are

According to the analytic solution to Equation (4), the components of meta-gene expressions in X driven by effective internal regulations ($\tilde{X}_{t+1}^I = \tilde{A}\tilde{X}_t^I$) are linear combinations of M_1 dynamic patterns determined by the eigenvalues of the effective system matrix \tilde{A} , as follows:

$$\tilde{X}_t^I = \sum_{p=1}^{M_1} \lambda_p^t \tilde{V}_p^A \quad (8)$$

, where λ_p (\tilde{V}_p^A) is the p^{th} eigenvalue (eigenvector) of \tilde{A} , which determines the p^{th} dynamic pattern driven by effective internal regulations, defined as the p^{th} internal principal dynamic pattern (iPDP) =

$[\lambda_p^1 \lambda_p^2 \dots \lambda_p^T]$. If an eigenvalue λ is complex when is \tilde{A} asymmetric, then its conjugate $\bar{\lambda}$ is also an eigenvalue, so we sum its iPDP and its conjugate eigenvalue, $\bar{\lambda}$'s iPDP as a unified iPDP with real elements equal to $[\lambda_p^1 + \bar{\lambda}_p^1 \lambda_p^2 + \bar{\lambda}_p^2 \dots \lambda_p^T + \bar{\lambda}_p^T]$. Similarly, the components of meta-gene expressions in X driven by effective external regulations from U , i.e., $\tilde{X}_{t+1}^E = \tilde{B}\tilde{X}_t^E$ are linear combinations of M_2 dynamic patterns determined by the eigenvalues of the effective system matrix \tilde{B} as follows:

$$\tilde{X}_t^E = \sum_{q=1}^{M_2} \sigma_q^t \tilde{V}_q^B \quad (9)$$

, where σ_q (\tilde{V}_q^B) is the q^{th} eigenvalue(eigenvector) of \tilde{B} , which determines q^{th} dynamic pattern driven by effective external regulations, defined as q^{th} external principal dynamic pattern (ePDP) = $[\sigma_q^1 \sigma_q^2 \dots \sigma_q^T]$.

If an eigenvalue σ is complex, then its conjugate $\bar{\sigma}$ is also an eigenvalue, so we sum its ePDP and its conjugate eigenvalue, $\bar{\sigma}$'s ePDP as a unified ePDP with real elements equal to

$[\sigma_p^1 + \bar{\sigma}_p^1 \sigma_p^2 + \bar{\sigma}_p^2 \dots \sigma_p^T + \bar{\sigma}_p^T]$.

Both internal and external principal dynamic patterns (PDPs) represent the canonical temporal expression trajectories, which can be increasing, damped oscillation and so on depending on PDP's eigenvalues (Table 1).

Table 1. Classification of canonical temporal expression trajectories for PDP eigenvalue types

PDP eigenvalue	Real					Complex (radius)			
	>1	$=1$	$<1 \ \& \ >0$	$<0 \ \& \ >-1$	$=-1$	<-1	>1	$=1$	<1
Canonical temporal expression trajectory	increasing	flat	decreasing	vibrating early	vibrating late	vibrating	undamped oscillation	damped oscillation	oscillation

Daifeng Wang 5/4/15 11:55 AM
Deleted: degradation (real eigenvalues <1), growth (real eigenvalues >1), flat (real eigenvalues $=1$), damped oscillation (radius of complex eigenvalues <1), undamped oscillation (radius of complex eigenvalues >1), and oscillation (radius of complex eigenvalues $=1$).

GRAPHIC

2.5 Identification of gene coefficients of principal expression dynamic patterns

Because genes and meta-genes have linear relationships in terms of their expression levels as Equation (2), the components of gene expression levels in X driven by internal regulations, X_t^I can be also expressed as linear combinations of M_1 iPDPs:

$$X_t^I = W_X \tilde{X}_t^I = \sum_{p=1}^{M_1} \lambda_p^I \underbrace{W_X \tilde{V}_p^A}_{C_p^A} = \sum_{p=1}^{M_1} \lambda_p^I C_p^A \quad (10)$$

, where $C_p^A = W_X \tilde{V}_p^A \in \mathbb{R}^{M_1 \times 1}$ includes the gene coefficients for p^{th} iPDP. The gene expression components driven by external regulations from U can be also expressed as linear combinations of M_2 ePDPs:

$$X_t^E = W_X \tilde{X}_t^E = \sum_{q=1}^{M_2} \sigma_q^E \underbrace{W_X \tilde{V}_q^B}_{C_q^B} = \sum_{q=1}^{M_2} \sigma_q^E C_q^B \quad (11)$$

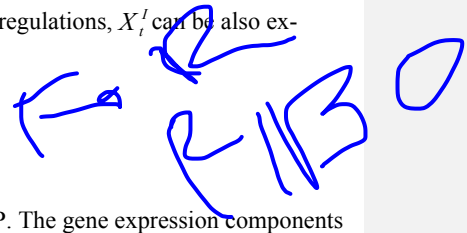
, where $C_q^B = W_X \tilde{V}_q^B \in \mathbb{R}^{M_2 \times 1}$ includes the gene coefficients for q^{th} ePDP.

3 RESULTS

Gene expression data during embryogenesis provide information about dynamics of genomic functions throughout developmental processes, from the conserved functions such as DNA replication to the species-specific functions such as body segmentations, but hardly reveal evolutionary gene regulatory subsystems that drive those developmental functions (Peter and Davidson, 2011). Thus, in order to understand the relationships between those subsystems and their driving genomic functions, we apply DREISS to worm and fly gene expression datasets during embryogenesis in modENCODE and identify various developmental genomic functions of worm-fly orthologous gene pairs driven by two different evolutionary regulatory subsystems, conserved (worm-fly TFs) and non-conserved (worm/fly TFs). As model organisms for developmental biology, both worm and fly have been used to study embryogenesis for decades.

3.1 Applications to worm and fly embryonic developmental data in modENCODE: orthologous genes, transcription factors and gene expression datasets

DREISS enables us to compare expression dynamic patterns between two or more temporal gene expression datasets even though they have different numbers of samples, as well as differences in the times at which those samples were collected. For example, we can apply DREISS to two different datasets of



Daifeng Wang 5/4/15 11:55 AM
Deleted: GRNs

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Peter and Davidson, 2011)

Daifeng Wang 5/4/15 11:55 AM
Deleted: evolutionary GRNs

Daifeng Wang 5/4/15 11:55 AM
Deleted: , we need advanced computational approaches to identify the effects of different evolutionary GRNs from gene expression data. In this study

Daifeng Wang 5/4/15 11:55 AM
Deleted: GRNs

Daifeng Wang 5/4/15 11:55 AM
Deleted: We found that the conserved GRNs drive similar genomic functions, but non-conserved GRNs drive species-specific functions of orthologous genes between worm and fly, implying that, in addition to having ancient conserved functions, orthologous genes have been regulated by evolutionarily younger GRNs to execute species-specific functions in evolution.

the same group of genes, and identify similar/different dynamic patterns driven by internal regulations captured by the eigenvalues of the effective system matrices between two datasets.

In this paper, we apply DREISS to 3,153 one-to-one orthologous genes between worm (*Caenorhabditis elegans*) and fly (*Drosophila melanogaster*) as Group X for their expression dynamics during embryonic development (Gerstein, et al., 2014). We refer to species-specific TFs as external regulations; i.e., Group U . We found that worm-fly orthologs have similar internal dynamic patterns, which may be mainly driven by conserved TFs, but have very different external dynamic patterns driven by species-specific TFs between worm and fly embryonic developmental stages. We focus on comparing internal dynamic patterns along with orthologous gene coefficients between worm and fly. The datasets are summarized as follows.

We define Group X as 3,153 one-to-one orthologous genes between worm and fly during embryonic development, and Group U as all species-specific TFs (509 worm-specific TFs, 442 fly-specific TFs) (Reece-Hoyes, et al., 2005; Shazman, et al., 2014). We used their temporal gene expression levels (as measured by the RPKM values in RNA-seq) during embryonic development from the modENCODE project (Gerstein, et al., 2014). The worm embryonic development dataset includes $T=25$ time stages at 0, 0.5, 1, 1.5, ..., 12 hours, and the fly dataset includes $T=12$ time stages at 0, 2, 4, ..., 22 hours, but $t=1,2,\dots,25$ for worm and $t=1,2,\dots,12$ for fly are used in this paper, representing the relative time points for the entire embryonic development processes. Because $M_1 + M_2 < T$ in Equation (7), we choose $M_1 = M_2 = 5$ meta-genes for fly ($T=12$), and find that five meta-genes of Group X and five meta-genes of Group U capture ~98% of the co-variation of orthologous gene expressions and fly-specific TF gene expressions, respectively. In order to compare worm and fly, we also choose $M_1 = M_2 = 5$ meta-genes for worm, which capture ~98% of the co-variation of orthologous gene expressions and worm-specific TF gene expressions.

3.2 Metagenes of orthologous genes between worm and fly have similar internal but different external principal dynamic patterns during embryonic development

We find that the meta-gene canonical temporal expression trajectories driven by conserved regulatory networks (i.e., internal principal dynamic patterns, iPDPs) include four major patterns in both worm and

Daifeng Wang 5/4/15 11:55 AM

Deleted: (Gerstein, et al., 2014)

Daifeng Wang 5/4/15 11:55 AM

Deleted: (Table 1).

Daifeng Wang 5/4/15 11:55 AM

Deleted: (Reece-Hoyes, et al., 2005; Shazman, et al., 2014)

Daifeng Wang 5/4/15 11:55 AM

Deleted: (Gerstein, et al., 2014)

fly embryonic development: 1) a highly varied pattern late, 2) a fast decaying pattern early, 3) a slowly increasing pattern, and 4) an oscillating pattern (Figure 3A). In contrast to the iPDP similarities, we find that worm and fly have very different external principal dynamic patterns (ePDPs) (Figure 3B); i.e., the canonical temporal expression trajectories driven by species-specific TFs. The meta-gene dynamic patterns driven by the worm-specific regulatory network; i.e., worm ePDPs consist of a varied pattern at late embryonic development, a varied pattern at early embryonic development, a fast increasing and then unvarying pattern, a decaying pattern, and an increasing pattern at late embryonic development. The fly ePDPs, however, have two fast decaying patterns at early embryonic development, a fast increasing pattern at late embryonic development, and a highly increasing oscillation pattern. Moreover, to see the eigenvalue variations across orthologous genes, we left one gene out, and then calculated eigenvalues, which gave the eigenvalue variations shown as error bars in Figure 3. The iPDP eigenvalues vary less than ePDP in both worm and fly.

The above results suggest that the conserved regulatory networks from orthologous meta-genes between worm and fly have similar effects to orthologous meta-genes, given their similar iPDPs (i.e., both have four patterns, as described above). The species-specific regulatory networks from species-specific meta-genes (i.e., worm-specific or fly specific TFs) have effects that differ from orthologous meta-genes for their different ePDPs.

3.3 Orthologous genes have correlated coefficients between worm and fly for their matched internal principal dynamic patterns

In both worm and fly, we obtain the similar four types of internally driven canonical temporal expression trajectories; i.e., internal principal dynamic patterns (iPDPs), so we are interested in seeing how individual orthologous genes relate to those dynamic patterns. We find that the worm-fly orthologous genes have correlated coefficients over each of four iPDPs. Based on Equation (10), we can obtain the coefficients of orthologous genes for each iPDP. We find that their coefficients are significantly correlated between worm and fly iPDPs with a similar pattern (Figure 4): $r=0.33$ ($p<2.2e-16$) for the highly varied pattern at late embryonic development, $r=0.66$ ($p<2.2e-16$) for the fast decaying pattern at early embryonic development, $r=0.67$ ($p<2.2e-16$) for the slowly increasing pattern during embryonic development, and $r=0.73$ ($p<2.2e-16$) for the oscillation pattern during embryonic development. This implies that, not only do the orthologous meta-genes have similar internal (conserved) regulatory effects (i.e.,

DON'T WE HAVE TO MATCH UP iPDP

- Daifeng Wang 5/4/15 11:55 AM
Deleted: (iPDP with the real eigenvalue No. 1);
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (iPDP with the real eigenvalue No. 2);
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (iPDP with the real eigenvalue No. 3);
- Daifeng Wang 5/4/15 11:55 AM
Deleted: iPDP with the complex eigenvalue (
- Daifeng Wang 5/4/15 11:55 AM
Deleted: 5
- Daifeng Wang 5/4/15 11:55 AM
Deleted: S1
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 1),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 2),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 3),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 4),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 5).
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDPs with real eigenvalue No. 1 and 2),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with real eigenvalue No. 3),
- Daifeng Wang 5/4/15 11:55 AM
Deleted: (ePDP with complex eigenvalue).
- Daifeng Wang 5/4/15 11:55 AM
Deleted: Figures 5 and S1.

Daifeng Wang 5/4/15 11:55 AM
Deleted: 6

similar iPDPs), but the orthologous genes also have similar internally-driven expression dynamics between worm and fly because they have significantly correlated coefficients for iPDPs. The ePDPs between worm and fly generally do not have similar matches, but if we flip worm ePDP No. 3, and compare with fly ePDPs No. 4 and No. 5, they are roughly representing the fast decaying patterns. We found that the orthologous gene coefficient correlations between those ePDP patterns are much lower ($r=0.12$ for worm ePDP No.3 vs. fly ePDP No. 4, and $r=0.18$ for worm ePDP No. 3 vs. fly ePDP No. 5).

DESC
CALC

3.4 Ribosomal genes have significantly larger coefficients for internal than external principal dynamic patterns, but signaling genes exhibit the opposite trend

The ribosome produces proteins, which is an ancient process and conserved across the worm and fly, which diverged roughly a billion years ago. The ribosomal genes are highly expressed during embryogenesis, since intensive cell division and migration require a large amount of proteins to be synthesized.

We collected 195 ribosome-related genes based on the GO annotations. We compared the iPDP and ePDP coefficients of ribosomal genes, and found that the iPDP coefficients are significantly larger than ePDP ones in both worm (KS-test $p<0.001$) and fly (KS-test $p<2.2e-16$) as shown in Figure 5A. This means that the ribosomal gene expression is significantly more driven by the conserved regulatory network than by the species-specific regulatory network, which is consistent with ribosomal genes having conserved functions during embryonic development.

Daifeng Wang 5/4/15 11:55 AM

Deleted: ~200 ribosomal

Daifeng Wang 5/4/15 11:55 AM

Deleted: 7A

The orthologous genes related to signal transduction for cell-cell communication (a significantly more recent evolutionary adaptation relative to ribosomes) exhibit the opposite trend. We found that 320 signaling genes from GO annotations have significantly larger ePDP coefficients than iPDP ones in both worm (KS-test $p<7e-4$) and fly (KS-test $p<6e-4$), as shown in Figure 5B. This result implies that the signaling gene expression is significantly more driven by the species-specific regulatory network than by the conserved regulatory network, which is consistent with the signaling genes typically being associated with species-specific functions, such as body plan establishment and cell differentiation.

Daifeng Wang 5/4/15 11:55 AM

Deleted: 7B

3.5 DNA replication and Proteasome machinery are enriched in orthologous genes with high coefficients for the dynamic patterns with fast growing canonical trajectories

We next turn to the biological meaning of individual canonical temporal expression trajectory for iPDPs and ePDPs. For the fast-decaying pattern (2nd iPDP), we found that the DNA replication is significantly

enriched in Top 300 (~10%) orthologous genes that have the most negative coefficients for this pattern, in both worm ($p < 1.6e-8$) and fly ($p < 4.5e-6$). The very negative coefficients for the fast decaying pattern mean high positive coefficients for a fast-growing pattern, (flipped 2nd iPDP), showing a drastic increase at the beginning of embryogenesis, then remain flat during the late embryogenesis (red curves in Figure 6). Most of the cell division of embryogenesis in both worm and fly happens approximately within the first 300 minutes. Then, the cell elongation and migration start to dominate the development (Bate and Martinez Arias, 1993; Baugh, et al., 2003). The mRNA abundance of the genes involved in DNA replication may change accordingly. This is well reflected by the second iPDP. Interestingly, the original expression patterns of those top orthologous genes actually do not have fast-growing patterns (black curves in Figure 6), probably because of the combined effects of both conserved and species-specific GRN. Maternal mRNAs, which are pre-loaded before fertilization, may also mask the fast growing pattern of DNA replication genes. This pattern could only be observed after we separated the effect of two types of TFs using DREISS. In addition, we did not find any enrichment of DNA replication in top genes of other iPDPs and ePDPs. Therefore, the iPDP patterns identified by our method reveal basic cellular process of both species (i.e. DNA replication), which should mainly be controlled by the conserved regulatory network.

Besides a fast growing pattern driven by conserved TFs, we also identified a fast growing pattern driven by non-conserved TFs for those two species. The Top 300 orthologous genes (~10%) with fast-growing worm ePDP and fly ePDP (i.e., driven by species-specific regulatory networks) shared 36 orthologous genes. 10 of them encode genes in the proteasome complex ($p\text{-value} < 1.2e-9$). Protein degradation is not only a key process in apoptosis, but also throughout the whole process of development (DeRenzo and Seydoux, 2004). For example, eliminating proteins that are no longer needed is a vital process during embryo development; e.g., the maternal proteins need to be cleaned as the embryogenesis proceeds. Previous reports also showed that different species usually have different maternal mRNA in the oocyte, which indicates that species-specific strategies might be utilized to regulate the protein degradation process (Shen-Orr, et al., 2010). In our study, after separating the effect of conserved and non-conserved regulatory networks, the protein degradation is significantly enriched in the genes majorly driven by species-specific TFs.

HW
FLIP

Daifeng Wang 5/4/15 11:55 AM
Deleted: means

Daifeng Wang 5/4/15 11:55 AM
Deleted: ,

Daifeng Wang 5/4/15 11:55 AM
Deleted: 8

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Bate and Martinez Arias, 1993; Baugh, et al., 2003)

Daifeng Wang 5/4/15 11:55 AM
Deleted: 8

Daifeng Wang 5/4/15 11:55 AM
Deleted: (DeRenzo and Seydoux, 2004)

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Shen-Orr, et al., 2010)

Besides the 36 shared genes in the fast-growing pattern driven by species-specific TFs, there are additional observations that we find interesting. Among the Top 300 worm orthologous genes with fast-growing ePDPs, genes involved in calcium ion binding (p-value<2e-6), GTP binding (p-value<7e-3) and neuron differentiation (p-value<0.05) are over-represented, which implies that they are activated in the early stage of embryogenesis by worm-specific TFs. This observation indicates the GRN of these genes have evolved after the speciation. Proteins involved in calcium ion binding or GTP binding usually play a role in cell signal transduction (Aspenstrom, 2004). In fact, the genes involved in Wnt signaling and MAPK signaling both exhibit a two-fold change.

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Aspenstrom, 2004)

In contrast, the Top 300 fly genes with a fast-growing ePDP show no enrichment in signaling transduction or cell differentiation. Instead, functions associated with respiration, such as oxidative phosphorylation, are enriched (p-value<5e-10). The enrichment of energy generation in the Top 300 fly genes with a fast-growing ePDP is probably indicative of the large energy requirement during fly embryogenesis (Tennessen, et al., 2014), which did not provide the evolutionary conservation of this energy-related gene regulation. Our result reveals that the fly genes associated with respiration are more up-regulated by fly-specific TFs relative to conserved TFs, and that this up-regulation evolved after the separation of worm and fly. In addition, the lack of signaling enrichment might be due to different sampling time points. It is well-known that the Wnt signaling in worms starts as early as at the 4-cell stage, when one cell receives the signal and starts differentiation (Sawa and Korswagen, 2013). The time-series worm transcriptome data used in our study may have the resolution to detect those processes. However, since each of the first 10 cell cycles takes less than 10 minutes in the fly embryo (Gilbert, 2000), the 2 hour time interval in fly data may not have the resolution to capture the early regulatory events, such as Wnt signaling.

Daifeng Wang 5/4/15 11:55 AM
Deleted: It is well-known the Wnt signaling in worms starts as early as at the 4-cell stage, when one cell receives the signal and starts differentiation (Sawa and Korswagen, 2013). The separation of regulatory effects showed that the expression of genes involved in signaling is more controlled by the species-specific TFs.

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Tennessen, et al., 2014)

Daifeng Wang 5/4/15 11:55 AM
Deleted: Since

Daifeng Wang 5/4/15 11:55 AM
Deleted: (Gilbert, 2000)

Daifeng Wang 5/4/15 11:55 AM
Deleted: .

4 DISCUSSION

In this paper, we developed a novel computational method, DREISS, which decomposes time-series expression data of a group of genes into the components driven by the regulatory network inside the group (internal regulatory subsystem), and the components driven by the external regulatory network consisting of regulators outside the group (external regulatory subsystem). We applied DREISS to the time-series gene expression datasets for worm and fly embryonic developments from the modENCODE project (Gerstein, et al., 2014), and compared the worm-fly orthologous gene expression dynamic patterns


Daifeng Wang 5/4/15 11:55 AM
Deleted: (Gerstein, et al., 2014)

MORE GENERAL

driven by the conserved regulatory network (i.e., regulation effects from orthologous TFs), with the patterns driven by the species-specific regulatory networks (i.e., regulation effects from worm or fly specific TFs). We found that the conserved TFs drive similar genomic functions, but non-conserved TFs drive species-specific functions of orthologous genes between worm and fly, implying that, in addition to having ancient conserved functions, orthologous genes have been regulated by evolutionarily younger GRNs to execute species-specific functions in evolution. This work can be extended to study the regulatory effects from orthologous TFs and species-specific TFs to species-specific genes. For example, one can find the expression dynamic patterns of worm/fly specific genes driven by specific TFs, and identify the genes with strong patterns associated with worm/fly specific functions, such as body formations. To the best of our knowledge, DREISS is the first method to reveal how the evolution of GRNs affects gene expression during embryogenesis.

We emphasize that DREISS is a general-purpose method (a free downloadable tool at github.com/gersteinlab/dreiss). Users can define the internal group (X) and external group (U) according to their interests. For example, if users want to identify the protein-coding expression patterns driven by miRNAs, they can define miRNAs as an external group and protein-coding genes as an internal group. Additionally, DREISS can be applied to more than two datasets, such as comparing worm, fly and human embryonic stem cell developmental data, and finding their conserved and specific expression patterns in development. The expression patterns driven by human-specific regulatory factors potentially help us understand human-specific developmental processes along with associated human genes. Moreover, by applying DREISS to human genes, such analyses may further our understanding of human development. For example, we may gain valuable insights into the development of the human brain and pathologies such as Alzheimer's disease.

Due to the limited time samples in gene expression datasets, DREISS uses the simple linear state space model (i.e. the first order linear invariant difference equation) to model the temporal gene expression dynamics, and identify principal temporal dynamic patterns. This model assumes that the gene regulatory networks controlling temporal gene expression dynamics don't change across the entire biological process such as (A, B) in Equation (1). Thus, based on the analytic analysis, the principal dynamic patterns (PDPs) must follow a small set of canonical temporal trajectories (Table 1). With dramatically increasing gene expression data, however, we can extend DREISS to more advanced models such as



Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [6]

Daifeng Wang 5/4/15 11:55 AM
Deleted: -

[switched and hybrid system models, non-linear models \(Schaft and Schumacher, 2000\)](#), which allows that the gene regulatory networks are time varying, and try to find the more temporal gene expression patterns capturing the more complex gene regulatory activities.

FIGURE CAPTIONS

Figure 1 DREISS workflow. 1: DREISS models temporal gene expression dynamics using state-space models in control theory. The “state” refers to the expressions for a large group of genes of interest, such as the worm-fly orthologous genes investigated here. The “control” refers to any other group of genes that contribute to gene expressions of the “state”, such as the species-specific TF studied here. **2:** it then projects high-dimensional gene expression space to lower-dimensional meta-gene expression spaces using dimensionality reduction techniques. **3:** it derives the effective state-space models for meta-genes so that model parameters can be estimated. **4:** it then identifies the meta-gene expression dynamic patterns; i.e., canonical temporal expression trajectories driven by “state” (internal) and by “control” (external) based on the analytic solutions to estimated models. **5:** it finally calculates the coefficients of genes for the dynamic patterns of linear transformations between genes and meta-genes.

Figure 2 State space model for genes and the effective model for meta-genes. A) linear state space model for a given subsystem’s gene expression; i.e., linear first-order difference equations in Equation (2), is used to formulate temporal gene expression dynamics for a given subsystem, the gene group X (comprising N_1 genes) with external regulations from the gene group U (comprising N_2 genes) at time points $1, 2, \dots, T$. The vector $X_t \in \mathfrak{R}^{N_1 \times 1}$, the “state”, includes N_1 gene expression levels at time t in group X , and the vector $U_t \in \mathfrak{R}^{N_2 \times 1}$, the “input or control”, includes N_2 gene expression levels at time t in group U . The system matrix $A \in \mathfrak{R}^{N_1 \times N_1}$ captures internal causal interactions among genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of A, A_{ij} describes the contribution from the j^{th} gene expression at time t to the i^{th} gene expression at the next time $t+1$). The control matrix $B \in \mathfrak{R}^{N_1 \times N_2}$ captures external causal regulations from the genes in U to genes in X (i.e., the $i^{\text{th}}, j^{\text{th}}$ element of B, B_{ij} describes the contribution from the j^{th} gene expression in U at time t to the i^{th} gene expression in X at the next time $t+1$). **B)** Meta-gene expression levels. The meta-gene expression levels are obtained by $\tilde{X}_t = W_X^* X_t; \tilde{U}_t = W_U^* U_t$, where $\tilde{X}_t \in \mathfrak{R}^{M_1 \times 1}$, the “meta-gene state”, includes M_1 ($\ll N_1$ and $\ll T$) meta-gene expression levels; i.e., the values of first M_1 singular vectors from singular value decomposition (SVD) of matrix $[X_1 X_2 \dots X_T]$ at time t in group X ; the vector $\tilde{U}_t \in \mathfrak{R}^{M_2 \times 1}$, the “meta-gene input or control”, includes M_2 ($\ll N_2$ and $\ll T$) meta-gene expression levels (i.e., the values of the first M_2 singular vectors from SVD of matrix $[U_1 U_2 \dots U_T]$ at time t in group U ; $W_X \in \mathfrak{R}^{N_1 \times M_1}$ is the linear projection matrix of SVD from M_1 meta-gene expression space to N_1 gene expression space in X ,

Daifeng Wang 5/4/15 11:55 AM

Deleted: Linear

Daifeng Wang 5/4/15 11:55 AM

Formatted: Font:Times New Roman, Not Bold

Daifeng Wang 5/4/15 11:55 AM

Formatted: Normal

Daifeng Wang 5/4/15 11:55 AM

Deleted: . The state-space model

Daifeng Wang 5/4/15 11:55 AM

Deleted: .

Daifeng Wang 5/4/15 11:55 AM

Formatted: Font:Not Bold

Daifeng Wang 5/4/15 11:55 AM

Deleted: .

$W_U \in \mathfrak{R}^{N_2 \times M_2}$ is the linear projection matrix of SVD from M_2 meta-gene expression space to N_2 gene expression space in U , and $(\cdot)^*$ is a pseudo-inverse operation; i.e., $W^*W=I$, where I is the identity matrix. Effective state space model for meta-genes. The effective state-space model for meta-genes, Equation (4) is obtained by using linear projections W_X and W_U between genes and meta-genes from Equations (1-3). The effective meta-gene system matrix $\tilde{A} = W_X^*AW_X \in \mathfrak{R}^{M_1 \times M_1}$ captures internal causal interactions among meta-genes in X (i.e., the i^{th} , j^{th} element of \tilde{A} (\tilde{A}_{ij}) describes the contribution from the j^{th} meta-gene expression at time t to i^{th} meta-gene expression at next time $t+1$), and the effective control matrix $\tilde{B} = W_X^*BW_U \in \mathfrak{R}^{M_1 \times M_2}$ captures external causal regulations from meta-genes in U to meta-genes in X (i.e., the i^{th} , j^{th} element of \tilde{B} , \tilde{B}_{ij} describes the contribution from the j^{th} meta-gene expression in U at time t to i^{th} meta-gene expression in X at next time $t+1$). Equation (4) describes the effective state space model for the meta-genes in X , whose expression dynamics are determined by \tilde{A} and \tilde{B} . Because the meta-gene dimension, M_1 (M_2) is less than T , and much less than N_1 (N_2), we can estimate \tilde{A} and \tilde{B} as follows.

Figure 3 Principal dynamic patterns of orthologous genes between worm and fly during embryonic development.

A) Metagenes of orthologous genes have similar internal driven principal dynamic patterns. Meta-gene canonical temporal expression trajectories driven by conserved regulatory networks (i.e., internal principal dynamic patterns, iPDPs) include four major patterns in both worm and fly embryonic development: 1) a highly varied pattern late (iPDP with the real eigenvalue No. 1); 2) a fast decaying pattern early (iPDP with the real eigenvalue No. 2); 3) a slowly increasing pattern (iPDP with the real eigenvalue No. 3); and 4) an oscillating pattern (iPDP with the complex eigenvalue). **B)** Metagenes of orthologous genes have different external driven principal dynamic patterns. Worm and fly have very different external principal dynamic patterns (ePDPs); i.e., the canonical temporal expression trajectories driven by species-specific TFs. The meta-gene dynamic patterns driven by the worm-specific regulatory network; i.e., worm ePDPs consist of a varied pattern at late embryonic development (real eigenvalue No. 1), a varied pattern at early embryonic development (real eigenvalue No. 2), a fast increasing and then unvarying pattern (real eigenvalue No. 3), a decaying pattern (real eigenvalue No. 4), and an increasing pattern at late embryonic development (real eigenvalue No. 5). The fly ePDPs, however, have two fast decaying patterns at early embryonic development (real eigenvalue No. 1 and 2), a fast increasing pattern at late embryonic development (real eigenvalue No. 3), and a highly increasing oscillation pattern (complex eigenvalue).

Figure 4 Orthologous genes have correlated coefficients between worm and fly for their matched internal principal dynamic patterns. The worm-fly orthologous genes have correlated coefficients over each of four

Daifeng Wang 5/4/15 11:55 AM
Deleted: .

Daifeng Wang 5/4/15 11:55 AM
Formatted: Font:Not Bold

Daifeng Wang 5/4/15 11:55 AM
Deleted: as follows: By replacing

Daifeng Wang 5/4/15 11:55 AM
Deleted:) using (2), we obtain that
 $W_X \tilde{X}_{t+1} = AW_X \tilde{X}_t + BW_U \tilde{U}_t$ (

Daifeng Wang 5/4/15 11:55 AM
Deleted:), and by multiplying the pseudo-inverse of W_X .

$W_X^* \in \mathfrak{R}^{M_1 \times N_1}$ s.t. $W_X^*W_X = I$ where I is an identity matrix, at both sides of (3),

Daifeng Wang 5/4/15 11:55 AM
Formatted: Font:Times New Roman

Daifeng Wang 5/4/15 11:55 AM
Deleted: .

Daifeng Wang 5/4/15 11:55 AM
Deleted: 5

Daifeng Wang 5/4/15 11:55 AM
Formatted: Font:Not Bold

Daifeng Wang 5/4/15 11:55 AM
Deleted: between worm and fly

Daifeng Wang 5/4/15 11:55 AM
Formatted: Font:Not Bold

Daifeng Wang 5/4/15 11:55 AM
Deleted: during embryonic development.

Daifeng Wang 5/4/15 11:55 AM
Formatted: Font:Bold

Daifeng Wang 5/4/15 11:55 AM
Moved (insertion) [7]

Daifeng Wang 5/4/15 11:55 AM
Deleted:

Daifeng Wang 5/4/15 11:55 AM
Deleted: 6

iPDPs. Their coefficients are significantly correlated between worm and fly iPDPs with a similar pattern: $r=0.33$ ($p<2.2e-16$) for the highly varied pattern at late embryonic development, $r=0.66$ ($p<2.2e-16$) for the fast decaying pattern at early embryonic development, $r=0.67$ ($p<2.2e-16$) for the slowly increasing pattern during embryonic development, and $r=0.73$ ($p<2.2e-16$) for the oscillation pattern during embryonic development.

Figure 5 Ribosomal genes have significantly larger coefficients for internal than external principal dynamic patterns, but signaling genes exhibit the opposite trend. **A)** The iPDP and ePDP coefficients of ribosomal genes are compared: the iPDP coefficients are significantly larger than ePDP ones in both worm (KS-test $p<0.001$) and fly (KS-test $p<2.2e-16$); **B)** The iPDP and ePDP coefficients of signaling genes (cell-cell communication) are compared: they have significantly larger ePDP coefficients than iPDP ones in both worm (KS-test $p<7e-4$) and fly (KS-test $p<6e-4$).

Figure 6 DNA replication is enriched in orthologous genes with high coefficients for the dynamic patterns with fast growing canonical trajectories. For the fast-decaying pattern (2nd iPDP), we found that the DNA replication is significantly enriched in Top 300 (~10%) orthologous genes that have the most negative coefficients for this pattern, in both worm ($p<1.6e-8$) and fly ($p<4.5e-6$). The very negative coefficients for the fast decaying pattern means high positive coefficients for a fast-growing pattern, showing a drastic increase at the beginning of embryogenesis, then remain flat during the late embryogenesis (red curves). The original expression patterns of those top orthologous genes actually do not have fast-growing patterns (black curves).

REFERENCES

- Aspenstrom, P. (2004) Integration of signalling pathways regulated by small GTPases and calcium, *Biochimica et biophysica acta*, **1742**, 51-58.
- Bansal, M., Della Gatta, G. and di Bernardo, D. (2006) Inference of gene regulatory networks and compound mode of action from time course gene expression profiles, *Bioinformatics*, **22**, 815-822.
- Bate, M. and Martinez Arias, A. (1993) *The Development of Drosophila melanogaster*. Cold Spring Harbor Laboratory Press, Plainview, N.Y.
- Baugh, L.R., et al. (2003) Composition and dynamics of the Caenorhabditis elegans early embryonic transcriptome, *Development*, **130**, 889-900.
- Brogan, W.L. (1991) *Modern control theory*. Prentice Hall, Englewood Cliffs, N.J.
- Casci, T. (2011) *Development: Hourglass theory gets molecular approval*, *Nature reviews. Genetics*, **12**, 76.
- Chu, S., et al. (1998) The transcriptional program of sporulation in budding yeast, *Science*, **282**, 699-705.
- DeRenzo, C. and Seydoux, G. (2004) A clean start: degradation of maternal proteins at the oocyte-to-embryo transition, *Trends in cell biology*, **14**, 420-426.
- Gerstein, M.B., et al. (2014) Comparative analysis of the transcriptome across distant species, *Nature*, **512**, 445-448.
- Gilbert, S.F. (2000) *Developmental biology*. Sinauer Associates, Sunderland, Mass.
- Huang, S. and Ingber, D.E. (2006) A non-genetic basis for cancer progression and metastasis: self-organizing attractors in cell regulatory networks, *Breast disease*, **26**, 27-54.
- Irie, N. and Kuratani, S. (2011) Comparative transcriptome analysis reveals vertebrate phylotypic period during organogenesis, *Nature communications*, **2**, 248.

Daifeng Wang 5/4/15 11:55 AM

Deleted: 7

Daifeng Wang 5/4/15 11:55 AM

Deleted: :

Daifeng Wang 5/4/15 11:55 AM

Deleted: :

Daifeng Wang 5/4/15 11:55 AM

Deleted: 8

Daifeng Wang 5/4/15 11:55 AM

Deleted: Figure S1 Metagenes of orthologous genes between worm and fly have different external driven principal dynamic patterns during embryonic development

Daifeng Wang 5/4/15 11:55 AM

Moved up [7]: . Worm and fly have very different external principal dynamic patterns (ePDPs); i.e., the canonical temporal expression trajectories driven by species-specific TFs. The meta-gene dynamic patterns driven by the worm-specific regulatory network; i.e., worm ePDPs consist of a varied pattern at late embryonic development (real eigenvalue No. 1), a varied pattern at early embryonic development (real eigenvalue No. 2), a fast increasing and then unvarying pattern (real eigenvalue No. 3), a decaying pattern (real eigenvalue No. 4), and an increasing pattern at late embryonic development (real eigenvalue No. 5). The fly ePDPs, however, have two fast decaying patterns at early embryonic development (real eigenvalue No. 1 and 2), a fast increasing pattern at late embryonic development (real eigenvalue No. 3), and a highly increasing oscillation pattern (complex eigenvalue).

Daifeng Wang 5/4/15 11:55 AM

Deleted: -

-
- Kalinka, A.T., *et al.* (2010) Gene expression divergence recapitulates the developmental hourglass model, *Nature*, **468**, 811-814.
- Kim, P.M. and Tidor, B. (2003) Subsystem identification through dimensionality reduction of large-scale gene expression data, *Genome research*, **13**, 1706-1718.
- Levin, M., *et al.* (2012) Developmental milestones punctuate gene expression in the *Caenorhabditis* embryo, *Developmental cell*, **22**, 1101-1108.
- Peter, I.S. and Davidson, E.H. (2011) Evolution of gene regulatory networks controlling body plan development, *Cell*, **144**, 970-985.
- Rangel, C., *et al.* (2004) Modeling T-cell activation using gene expression profiling and state-space models, *Bioinformatics*, **20**, 1361-1372.
- Reece-Hoyes, J.S., *et al.* (2005) A compendium of *Caenorhabditis elegans* regulatory transcription factors: a resource for mapping transcription regulatory networks, *Genome biology*, **6**, R110.
- Saeyns, Y., Inza, I. and Larranaga, P. (2007) A review of feature selection techniques in bioinformatics, *Bioinformatics*, **23**, 2507-2517.
- Sawa, H. and Korswagen, H.C. (2013) Wnt signaling in *C. elegans*, *WormBook : the online review of C. elegans biology*, 1-30.
- [Schaft, A.J.v.d. and Schumacher, J.M. \(2000\) *An introduction to hybrid dynamical systems. Lecture notes in control and information sciences. Springer, London ; New York.*](#)
- Shazman, S., *et al.* (2014) OnTheFly: a database of *Drosophila melanogaster* transcription factors and their binding sites, *Nucleic acids research*, **42**, D167-171.
- Shen-Orr, S.S., Pilpel, Y. and Hunter, C.P. (2010) Composition and regulation of maternal and zygotic transcriptomes reflects species-specific reproductive mode, *Genome biology*, **11**, R58.
- Tennessen, J.M., *et al.* (2014) Coordinated metabolic transitions during *Drosophila* embryogenesis and the onset of aerobic glycolysis, *G3*, **4**, 839-850.
- Vilar, J.M. (2006) Modularizing gene regulation, *Molecular systems biology*, **2**, 2006 0016.
- Wang, D., *et al.* (2012) Principal-oscillation-pattern analysis of gene expression, *PLoS one*, **7**, e28805.
- Wang, D., *et al.* (2012) Eigen-genomic system dynamic-pattern analysis (ESDA): modeling mRNA degradation and self-regulation, *IEEE/ACM transactions on computational biology and bioinformatics / IEEE, ACM*, **9**, 430-437.
- Yanai, I., *et al.* (2011) Mapping gene expression in two *Xenopus* species: evolutionary constraints and developmental flexibility, *Developmental cell*, **20**, 483-496.

Page 3: [1] Deleted **Daifeng Wang** **5/4/15 11:55 AM**

Gene regulatory networks are highly modular, and consist of various sub-networks. Each sub-network includes a number of regulatory factors representing a subsystem to drive particular gene regulatory functions (Kim and Tidor, 2003; Vilar, 2006).

Page 3: [2] Moved to page 2 (Move #4) Daifeng Wang **5/4/15 11:55 AM**

Unfortunately, existing experimental gene expression data cannot decouple the expression components that are driven by different subsystems.

Page 3: [3] Moved to page 2 (Move #3) Daifeng Wang **5/4/15 11:55 AM**

For example, the orthologous genes can be regulated by both orthologous and species-specific transcription factors (TFs). The orthologous TFs constitute an “internal” regulatory network, while the species-specific TFs constitute an “external” regulatory network.

Page 3: [4] Deleted **Daifeng Wang** **5/4/15 11:55 AM**

Thus, we need computational methods to derive

Page 3: [5] Moved to page 2 (Move #5) Daifeng Wang **5/4/15 11:55 AM**

the contribution from each factor or subsystem from the gene expression data. In this study, we propose a novel computational method, DREISS - Decomposition of gene Regulatory network into External and Internal components based on State Space models. We identify temporal gene expression dynamic patterns for evolutionarily conserved genes during embryonic development, as driven by conserved and species-specific regulatory subsystems. This advances our current understanding of GRNs in evolution

Page 3: [6] Deleted **Daifeng Wang** **5/4/15 11:55 AM**

(Bansal, et al., 2006; Huang and Ingber, 2006; Rangel, et al., 2004)

Page 3: [7] Deleted **Daifeng Wang** **5/4/15 11:55 AM**

(Chu, et al., 1998; Kim and Tidor, 2003; Saeys, et al., 2007)