

# Algorithm for cellular reprogramming

Scott Ronquist<sup>a</sup>, Geoff Patterson<sup>b</sup>, Lindsey A. Muir<sup>c</sup>, Stephen Lindsly<sup>a</sup>, Haiming Chen<sup>a</sup>, Markus Brown<sup>d</sup>, Max S. Wicha<sup>e</sup>, Anthony Bloch<sup>f</sup>, Roger Brockett<sup>g</sup>, and Indika Rajapakse<sup>a,f,1</sup>

<sup>a</sup>Department of Computational Medicine and Bioinformatics, Medical School, University of Michigan, Ann Arbor, MI 48109; <sup>b</sup>Department of Curriculum Design, IXL Learning, Raleigh, NC 27560; <sup>c</sup>Department of Pediatrics and Communicable Diseases, University of Michigan, Ann Arbor, MI 48109; <sup>d</sup>Department of Biological Sciences, University of Maryland, College Park, MD 20742; <sup>e</sup>Department of Hematology/Oncology, University of Michigan, Ann Arbor, MI 48109; <sup>f</sup>Department of Mathematics, University of Michigan, Ann Arbor, MI 48109; and <sup>g</sup>John A. Paulson School of Engineering and Applied Science, Harvard University, Cambridge, MA 02138

Edited by Steven Henikoff, Fred Hutchinson Cancer Research Center, Seattle, WA, and approved September 26, 2017 (received for review July 14, 2017)

**The day we understand the time evolution of subcellular events at a level of detail comparable to physical systems governed by Newton's laws of motion seems far away. Even so, quantitative approaches to cellular dynamics add to our understanding of cell biology. With data-guided frameworks we can develop better predictions about, and methods for, control over specific biological processes and system-wide cell behavior. Here we describe an approach for optimizing the use of transcription factors (TFs) in cellular reprogramming, based on a device commonly used in optimal control. We construct an approximate model for the natural evolution of a cell-cycle-synchronized population of human fibroblasts, based on data obtained by sampling the expression of 22,083 genes at several time points during the cell cycle. To arrive at a model of moderate complexity, we cluster gene expression based on division of the genome into topologically associating domains (TADs) and then model the dynamics of TAD expression levels. Based on this dynamical model and additional data, such as known TF binding sites and activity, we develop a methodology for identifying the top TF candidates for a specific cellular reprogramming task. Our data-guided methodology identifies a number of TFs previously validated for reprogramming and/or natural differentiation and predicts some potentially useful combinations of TFs. Our findings highlight the immense potential of dynamical models, mathematics, and data-guided methodologies for improving strategies for control over biological processes.**

cellular reprogramming | control theory | time series data | genome architecture | networks

In 1989, pioneering work by Weintraub et al. (1) successfully reprogrammed human fibroblasts into muscle cells via overexpression of transcription factor (TF) MYOD1, becoming the first study to demonstrate that the natural course of cell development could be altered. In 2007, Yamanaka and coworkers (2) changed the paradigm further by successfully reprogramming human fibroblasts into an embryonic stem-cell-like state [induced pluripotent stem cells (iPSCs)], using four TFs: POU5F1, SOX2, KLF4, and MYC. This work showed that a differentiated cell state could be reverted to a more pluripotent state. These discoveries have changed the trajectory of regenerative medicine, opening the possibility of generating needed cell types on demand for repairing damaged or diseased tissues. Ultimately, patient-derived fibroblasts could be used in autologous transplantations to minimize immune incompatibility.

These remarkable findings also demonstrate that the genome is a system capable of being controlled via an external input of TFs. In this context, determining how to push the cell from one state to another is, at least conceptually, a classical problem of control theory (3). The difficulty arises in the fact that the dynamics—and even proper representations of the cell state and inputs—are not well defined in the context of cellular reprogramming. Nevertheless, it seems natural to treat reprogramming as a problem in control theory, with the final state being the desired reprogrammed cell. In this paper, we provide such a framework based on empirical data and demonstrate the poten-

tial of this framework to provide insights into cellular reprogramming (4).

Our goal is to mathematically identify TFs that can directly reprogram human fibroblasts into a desired target cell type. As part of our methodology, we create a model for the natural dynamics of proliferating human fibroblasts, using time series data collected throughout the cell cycle. We couple data from bioinformatics with methods of mathematical control theory—a framework that we dub data-guided control (DGC). We use this model to determine a principled way to identify the best TFs for efficient reprogramming.

Previously, selection of TFs for reprogramming has been based largely on trial and error, typically relying on TF differential expression between cell types for initial predictions. Recent work has sought to predict TFs for reprogramming the cell state (5–8). Rackham et al. (7) devised a predictive method based on differential expression, as well as gene and protein network data. Our approach is fundamentally different in that we take a dynamical systems point of view, opening avenues for investigating efficiency (probability of conversion), timing (when to introduce TFs), and optimality (minimizing the number of TFs and amount of input).

Our method identifies TFs previously found to reprogram human fibroblasts into embryonic stem-cell-like cells, muscle cells, and many additional target cell types. Furthermore, our analysis predicts the points in the cell cycle at which the introduction of TFs might most efficiently affect a desired change of cell

## Significance

**Reprogramming the human genome toward any desirable state is within reach; application of select transcription factors drives cell types toward different lineages in many settings. We introduce the concept of data-guided control in building a universal algorithm for directly reprogramming any human cell type into any other type. Our algorithm is based on time series genome transcription and architecture data and known regulatory activities of transcription factors, with natural dimension reduction using genome architectural features. Our algorithm predicts known reprogramming factors, top candidates for new settings, and ideal timing for application of transcription factors. This framework can be used to develop strategies for tissue regeneration, cancer cell reprogramming, and control of dynamical systems beyond cell biology.**

Author contributions: R.B. and I.R. designed research; S.R., G.P., A.B., R.B., and I.R. performed research; S.R., G.P., S.L., H.C., M.B., M.S.W., A.B., R.B., and I.R. analyzed data; and S.R., G.P., L.A.M., S.L., A.B., R.B., and I.R. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

This is an open access article distributed under the [PNAS license](#).

<sup>1</sup>To whom correspondence should be addressed. Email: [indika@umich.edu](mailto:indika@umich.edu).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1712350114/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1712350114/-DCSupplemental).



identity plus a rank one matrix chosen to match the data for each time step  $k$ ; we impose the condition that without inputs we have  $x_{k+1} - \bar{x} = A_k(x_k - \bar{x})$ .

Define a time-dependent state transition matrix  $A_k$  as

$$A_k := I_{\tilde{N}} + \frac{(x_{k+1} - x_k)x_k^T}{x_k^T x_k} \in \mathbb{R}^{\tilde{N} \times \tilde{N}}, \quad k = 1, 2, 3, 4, 5, \quad [3]$$

where  $I_{\tilde{N}}$  is the  $\tilde{N} \times \tilde{N}$  identity matrix. Let the measured values of the state of the unforced evolution be  $x_1, x_2, \dots, x_5$ ; let the controls be labeled  $u_1, u_2, \dots, u_5$ ; let the values of the state with the controls acting be  $z_2, z_3, \dots, z_6$ . Letting  $z$  denote the deviation from the cell-cycle average, we have

$$z_{k+1} = \left( I + \frac{1}{x_k^T x_k} (x_{k+1} - x_k)x_k^T \right) z_k + Bu_k,$$

where  $A_k$  is as above. Solving this difference equation, we have

$$z_k = \prod_{i=1}^{k-1} A_i x_1 + \sum_{i=1}^{k-1} \prod_{j=i}^{k-1} A_j B u_i$$

with the understanding that  $A_0 = I$ . This explicit expression shows that the effect of the  $u(i)$  cannot be inferred from the sum of the  $u_i$  because different  $u_i$  are weighted in different ways, dependent on the stage in the cell cycle at which it is applied. This is a significant point relating to the model and plays a significant role in determining the optimal times for inserting TFs.

**Input Matrix and Input Signal:  $B, u_k$ .** With the natural TAD-level dynamics established in the context of our control Eq. 1, we turn our attention to quantifying methods for control.

A TF can regulate a gene positively or negatively by binding to a specific DNA sequence near a gene and encouraging or discouraging transcription. The degree to which a TF activates or represses gene expression depends on the specific TF–gene interaction, which is influenced by a variety of factors that are difficult to quantify. Let  $w_{i,m}$  be the theoretical regulation weight of TF  $m$  on gene  $i$ , where  $w_{i,m} > 0$  ( $w_{i,m} < 0$ ) if TF  $m$  activates (represses) gene  $i$ , and  $m = 1, \dots, M$ , where  $M$  is the total number of well-characterized TFs. Weights that are bigger in absolute value,  $|w_{i,m}| \gg 0$ , indicate stronger transcriptional influence, and weights equal to zero,  $w_{i,m} = 0$ , indicate no influence.

Extensive TF perturbation experiments would be needed to determine  $w_{i,m}$  for each TF  $m$  on each gene  $i$ . Instead, we propose a simplified method to approximate  $w_{i,m}$  from existing, publicly available data for TF binding sites (TFBSs), gene accessibility, and average activator/repressor activity. To determine the number of possible binding sites a TF  $m$  recognizes near gene  $i$ , the reference genome was scanned for the locations of potential TFBSs following methods outlined by Neph et al. (14) (SI Appendix). Position frequency matrices (PFMs), which give information on TF–DNA binding probability, were

obtained for 547 TFs from a number of publicly available sources ( $M = 547$ ). Let  $c_{i,m}$  be the number of TF  $m$  TFBSs found within  $\pm 5$  kb of the transcriptional start site (TSS) of gene  $i$  (SI Appendix, Fig. S2).

Although many TFs can do both in the right circumstances, most TFs have a tendency toward either activator or repressor activity (15). That is, if TF  $m$  is known to activate (repress) most genes, we can say with some confidence that TF  $m$  is an activator (repressor), so  $w_{i,m} \geq 0$  ( $w_{i,m} \leq 0$ ) for all  $i$ . To determine a TF's function, we performed a literature search for all 547 TFs and labeled 299 as activators and 124 as repressors (SI Appendix). The remaining TFs were labeled unknown for lack of conclusive evidence and were evaluated as both an activator and a repressor in separate calculations. Here, we define  $a_m$  as the activity of TF  $m$ , with 1 and  $-1$  denoting activator and repressor, respectively.

TFBSs are cell-type invariant since they are based strictly on the linear genome. However, it is known that for a given cell type, certain areas of the genome may be opened or closed, depending on epigenetic aspects. To capture cell-type-specific regulatory information, we obtained publicly available gene accessibility data (DNase-seq) on human fibroblasts (GSM1014531). DNase-seq extracts cell-type-specific chromatin accessibility information genome-wide by testing the genome's sensitivity to the endonuclease DNase I and sequencing the nondigested genome fragments. These data are used for our initial cell type to determine which genes are available to be controlled by TFs (16). Here, we define  $s_i$  to be the DNase I sensitivity information (accessibility; open/close) of gene  $i$  in the initial state, with 1 and 0 denoting accessible and inaccessible, respectively (SI Appendix).

We approximate  $w_{i,m}$  as

$$w_{i,m} := a_m s_i c_{i,m}, \quad [4]$$

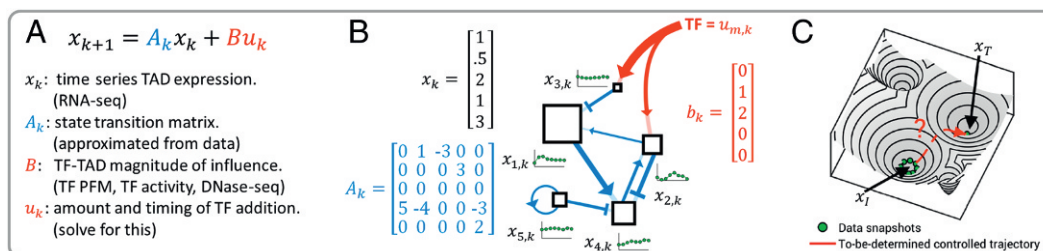
so that the magnitude of influence is equal to the number of observed consensus motifs  $c_{i,m}$ , except when the gene is inaccessible ( $s_i = 0$ ) in which case  $w_{i,m} = 0$ .

Since we are working off a TAD-dimensional model, our input matrix  $B$  must match this dimension. Let  $b_m$  be a 2,245-dimensional vector, where the  $j$ th component is

$$b_{j,m} := \sum_{i \text{ s.t. } \text{tad}(i)=j} w_{i,m}, \quad [5]$$

and define a matrix  $B = [b_1 \ b_2 \ \dots \ b_M]$ .

The amount of control input is captured in  $u_k$ , which is an  $\mathbb{R}^{M \times 1}$  vector representing the quantity of the external TFs we are inputting to the system (cell) at time  $k$ . This can be controlled by the researcher experimentally through manipulation of the TF concentration (17). In this light, we restrict our analysis to  $u_k \geq 0$  for all  $k$ , as TFs cannot be subtracted from the cell.  $u_{m,k}$  is defined as the amount of TF  $m$  to be added at time point  $k$ .



**Fig. 2.** DGC overview. (A) Summary of control equation variables. (B) Each TAD is a node in a dynamic network. The blue connections represent the edges of the network and are determined from time series fibroblast RNA-seq data. The green plots represent the expression of each TAD changing over time. The red arrows indicate additional regulation imposed by exogenous TFs. (C) A conceptual illustration of the problem: Can we determine TFs to push the cell state from one basin to another?



With all variables of our control Eq. 1 defined, we can now attempt to predict which TFs will most efficiently achieve cellular reprogramming from some  $x_I$  (initial state; fibroblast in our setting) to  $x_T$  (target state; any human cell type for which compatible RNA-seq data are available) through manipulation of  $u_k$ . An overview of our DGC framework is given in Fig. 2.

**Selection of TFs.** Our general procedure for scoring TFs is explained as follows. Eq. 1 has an explicit solution that is given below. The first few terms are

$$\begin{aligned} z_2 &= A_1 x_1 + B u_1 \\ z_3 &= A_2 A_1 x_1 + A_2 B u_1 + B u_2 \\ z_4 &= A_3 A_2 A_1 x_1 + A_3 A_2 B u_1 + A_3 B u_2 + B u_3 \\ &\vdots \end{aligned}$$

This shows how  $z_4$  depends on  $u_1$ ,  $u_2$ , and  $u_3$ .

If  $x_T$  is a target condition, then the Euclidean distance  $\|\cdot\|$  can be used to measure how close a state is to the target state. We define

$$d = \|x_T - z_6(u)\|, \quad [6]$$

where the notation  $z_6(u)$  is used to emphasize the dependence of  $z_6$  on  $u$ . Considering all possible input signals, one can compute the optimal control that finds the minimum distance for a given initial and target cell type. Let  $u_*$  denote the optimal  $u$  used to minimize  $d$  and  $d_*$  denote this minimum distance value.

When appropriate, we write  $z_6(u)$  to emphasize the fact that the final state depends on the input. The Euclidean distance  $\|\cdot\|$  can be used to measure how close a given state is to the target. If there were no restrictions on the  $u$  terms, the control that minimizes the distance between  $z_6$  and the target could be computed without difficulty. However, there are reasons for restricting the number of different TFs used in any one trial. Transfection of cells with too many TFs can lower the efficiency of transfection and even lead to cell death. Moreover, many confirmed direct reprogramming experiments use  $\leq 4$  TFs to achieve reprogramming. For these reasons, we modify the optimization problem by adding the constraint that there are no more than a fixed number of TFs (components of  $u$ ) used in a given trial.

Let  $\hat{p}$  be a set of integers that identifies the subset of the components of  $u$  (read: TFs) that are allowed to be nonzero. For example,  $\hat{p} = \{1, 4, 7\}$  refers to TFs 1, 4, and 7. Let  $p$  be the number of elements in  $\hat{p}$ . Given a set of TFs,  $\hat{p}$ , we determine the quantity and timing of TF input,  $u_{m,k}$ , that minimizes the difference between  $z_6$  and the target cell state,  $x_T$ . Mathematically, this can be written as

$$\begin{aligned} &\underset{u}{\text{minimize}} \quad \|x_T - z_6(u)\| \\ &\text{subject to} \quad \begin{cases} u_{m,k} \geq 0, & k = 1, \dots, 5 \\ u_{m,k} = 0, & \text{if } m \notin \hat{p} \\ u_{m,k+1} \geq u_{m,k} \end{cases} \end{aligned} \quad [7]$$

We use MATLAB's *fmincon* function to solve Eq. 7, which gives  $u_{m,k}$  and  $d_*$ .

Let  $d_0 := \|x_T - x_0\|$  be the distance between the final state and target state with no control input. Define a score  $\mu := d_0 - d_*$ , which can be interpreted as the improvement provided by a particular choice of  $u$ . This can be calculated for each  $\hat{p}$  and sorted (high to low) to determine which TF or TF combination is the best candidate for direct reprogramming between  $x_0$  and  $x_T$ .

We consider different scenarios for the type of input regime in the results. The first one assumes the input signal is constant  $u_1 = u_k = \bar{u}$ , intended to mimic empirical regimes where TFs are given at a single time point. Later, we also consider inputting TFs at different times  $\hat{k}$ , which can be viewed mathematically

as requiring  $u_{m,k} = 0$  for all  $k < \hat{k}$ , and  $u_{m,k}$  is a constant value for all  $k \geq \hat{k}$ . This is intended to mimic inputting a TF at time  $\hat{k}$ , which will continue to express at a constant level until time point  $k = 6$ .

**Remark.** Subsets of TFs were chosen for each calculation based on the following criteria:  $\geq 10$ -fold expression increase in target state compared with initial state and  $\geq 10$  RPKM in target state. These criteria are used to select differentially expressed TFs and TFs that are sufficiently active in the target state.

## Results

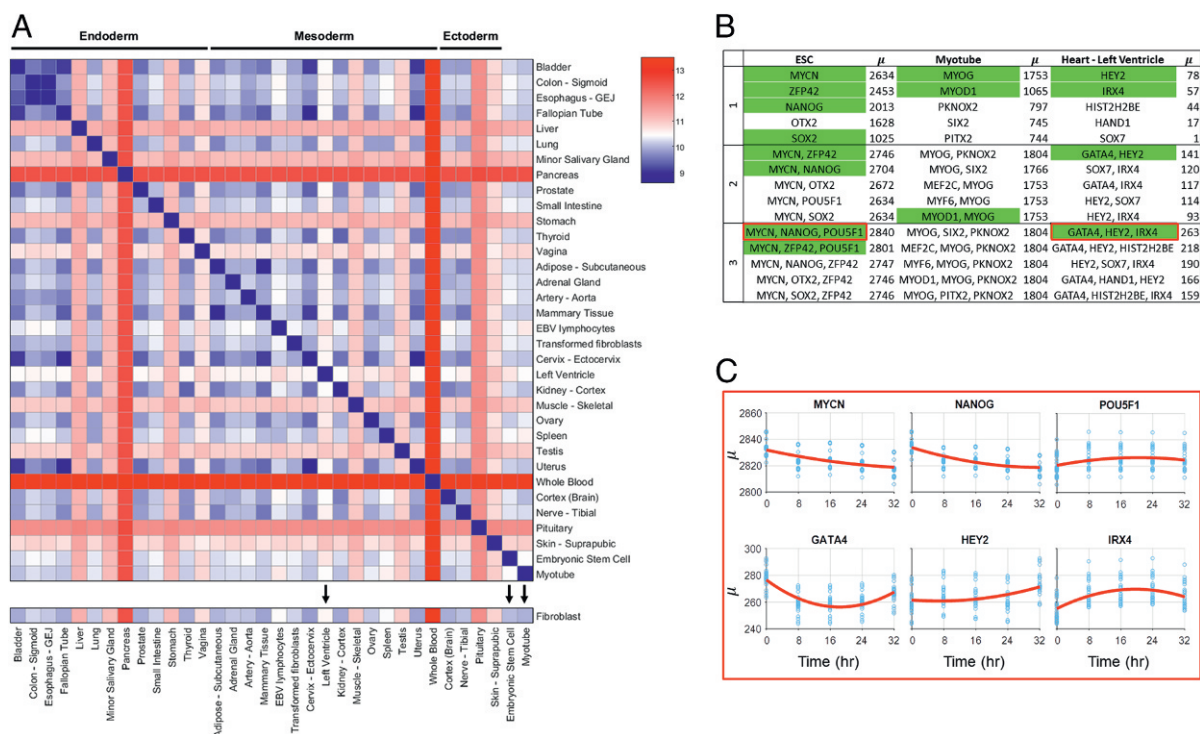
**Quantitative Measure Between Cell Types.** To best use our algorithm to predict TFs for reprogramming, compatible data on target cell types must be collected. For this, we explore a number of publicly available databases where RNA-seq has been collected, along with RNA-seq data collected in our laboratory. The ENCODE Consortium has provided data on myotubes and embryonic stem cells (ESCs) (*SI Appendix*) (18). The GTEx portal provides RNA-seq data on a large variety of different human tissue types (19). Although each GTEx experiment is performed on tissue samples, thus containing multiple different cell types, we use these data as more general cell-state targets.

To give a numerical structure to cell-type differences, conceptually similar to Waddington's epigenetic landscape, we calculate  $d_0$  between all cell types collected. Fig. 3A shows  $d_0$  values for 32 tissue samples collected from the GTEx portal, along with ESC, myotube, and our fibroblast data (additional cell-type  $d_0$  values shown in *SI Appendix*). GTEx RNA-seq data are scaled to keep total RPKM difference between time series fibroblast and GTEx fibroblast RNA-seq minimal (*SI Appendix*).

**TF Scores.** To assess our method's predictive power, a subset of target cell types is presented here that has validated either TF reprogramming methods or TFs highly associated with the target cell type. Additional predicted TFs for reprogramming are included in *SI Appendix*. We note that although experimentally validated TFs provide the best current standard for comparison, we believe experimental validation with our predicted TFs may provide more efficient and comprehensive reprogramming results. For all reprogramming regimes presented in this section, fibroblast is used as the initial cell type due to the availability of synchronized time series data, and all TFs are introduced at  $k = 1$  (11).

For conversion of fibroblast to myotubes, the top predicted single-input TFs are MYOG and MYOD1, both of which are known to be crucial for myogenesis. While MYOD1 is the classic master regulator reprogramming TF for myotube conversion, activation of downstream factor MYOG is necessary for full conversion (20). For fibroblast to ESC conversion, a number of TFs known to be necessary for pluripotency are predicted, including MYCN, ZFP42, NANOG, and SOX2 (2). With the knowledge that no single TF has been shown to fully reprogram a fibroblast to an embryonic state, combinations of TFs are more informative for this analysis. The top-scoring combination of three TFs is MYCN, NANOG, and POU5F1—three well-known markers for pluripotency (2). Interestingly, POU5F1 scores poorly when input individually, but is within the top set of three TFs when used in combination with MYCN and NANOG. Left ventricle reprogramming includes TFs that are known to be necessary for natural differentiation in the top score for all one to three combinations. These include GATA4 (a known TF in fibroblast to cardiomyocyte reprogramming), HEY2, and IRX4 (21–23).

**Time-Dependent TF Addition.** Fibroblast to ESC conversion was of particular interest in our analysis as this is a well-studied



**Fig. 3.** Quantitative measure between cell types and TF scores. (A)  $d_0$  values between GTEx tissue types and ESC, myotube, and fibroblast. Tissue types and cell types with black arrows have predicted TFs for reprogramming from fibroblasts shown in B. (B) Table of predicted TFs for a subset of cell and tissue types. Top five TFs for combinations of one to three are shown. Green labeled TFs are highly associated with the differentiation process of the target cell type and/or validated for reprogramming. These TFs are discussed in the main text. (C) Time-dependent scores for selected combinations of three TFs for fibroblast to ESC and fibroblast to "heart - left ventricle." x axis refers to time of TF addition, and y axis refers to  $\mu$ .

regime with a number of validated TFs (with a variety of reported efficiencies), and this conversion is promising for its regenerative medicine application. High-scoring TFs yield many that are known markers for pluripotency, but the top combination of three, MYCN, NANOG, and POU5F1, has not been used specifically together, to our knowledge. Here, we analyzed how the TF combination would score if input at different points throughout the cell cycle.

Time-dependent analysis of the top-scoring ESC TFs reveals that scores vary widely, depending on the time of input. MYCN and NANOG show a strong preference for input at the beginning of the cell cycle, while POU5F1 shows a slight preference for input toward the end of the cell cycle, with the highest score achieved when MYCN and NANOG are input at 0 h and POU5F1 is input at 32 h. Analysis on how the time of input control affects  $\mu$  is shown in Fig. 3C. Time-dependent analysis was also conducted for the top combination of three TFs for fibroblast to left ventricle. This analysis predicted that the best reprogramming results would occur if GATA4 is given immediately (0 h), with IRX4 and HEY2 given later (24 and 32 h, respectively).

## Discussion

The results from this algorithm show promise in their prediction of known reprogramming TFs and demonstrate the importance of including time series data for gene network dynamics. Time of input control has shown to have an impact on the end cell state, in line with what has been shown in natural differentiation (24).

While we believe that this is the best model currently available for predicting TFs for reprogramming, we are aware of its limitations and assumptions. TAD-based dimension reduction is based on the observation that genes within them correlate in expression over time, although we lack definitive proof of regu-

lation by shared transcriptional machinery (11). This assumption was deemed necessary for dimension reduction in the context of deriving transition matrix  $A_k$ . With finer time steps in RNA-seq data, the assumption may not be necessary for TF prediction, at the cost of increased computation time. Additionally, a 5-kb window flanking the TSS of each gene was used to ensure that all potential regulators are found, at the cost of potential inclusion of false positive motifs.

Although this program can score TFs relative to other TFs in a given reprogramming regime, it is difficult to predict a  $\mu$  threshold that would guarantee conversion. Additionally, rigorous experimental testing will be required to validate these findings and determine how our  $u$  vector translates to TF concentration. This is a product of the large number of assumptions that we have made to develop the initial framework for a reprogramming algorithm. With finer resolution in the time series gene expression, more subtle aspects of the genomic network may be observed, allowing for better prediction.

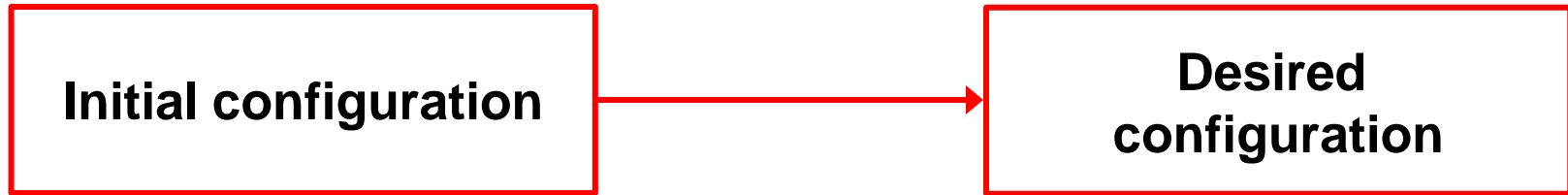
Our proposed DGC framework successfully identified known TFs for fibroblast to ESC and fibroblast to muscle cell reprogramming regimes. We use a biologically inspired dimension reduction via TADs, a natural partitioning of the genome. This comprehensive state representation was the foundation of our framework, and the success of our methods motivates further investigation of the importance of TADs as functional units to control the genome.

A dynamical systems view of the genome allows for analysis of timing, efficiency, and optimality in the context of reprogramming. Our framework is the first step toward this view. The successful implementation of time-varying reprogramming regimes would open unique avenues for direct reprogramming. This template can be used to develop regimes for changing any cell into any other cell, for applications that include reprogramming



# Controllability: An engineering perspective

---



Given an initial condition, target configuration, and the dynamics, are there input(s) that steer the system towards the desired configuration?

$$\frac{dx}{dt} = \mathbf{A}(t)x(t) + \mathbf{B}u(t)$$

Discrete model

$$x_{k+1} = \mathbf{A}_k x_k + \mathbf{B}u_k$$

## Controllability: A genome perspective

---

- Given any desired cell type, which TFs needed for conversion?
- Classical problem of control theory:

$$\left\{ \begin{array}{l} \text{Find } u \text{ such that} \\ \dot{x} = f(x, u) \\ x(0) = \text{initial cell type} \\ x(T) = \text{desired cell type} \end{array} \right.$$

Difficulty: what is  $x$ ? what is  $f$ ? what is  $u$ ?

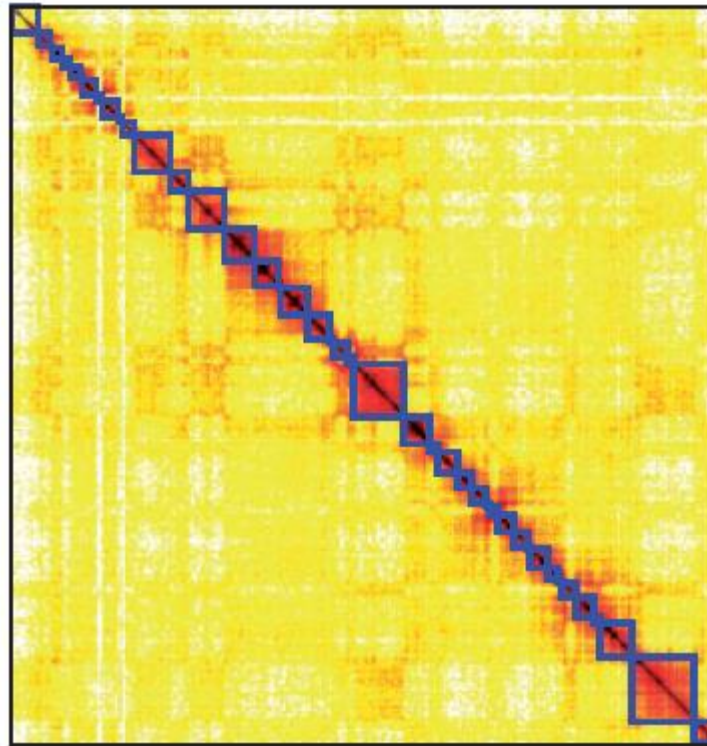
- We have data: RNA-seq, Hi-C, Transcription Factor binding

**TF Binding data from TRANSFAC and JASPAR to define a B matrix for each TF**



# Topologically Associating Domains (TADs)

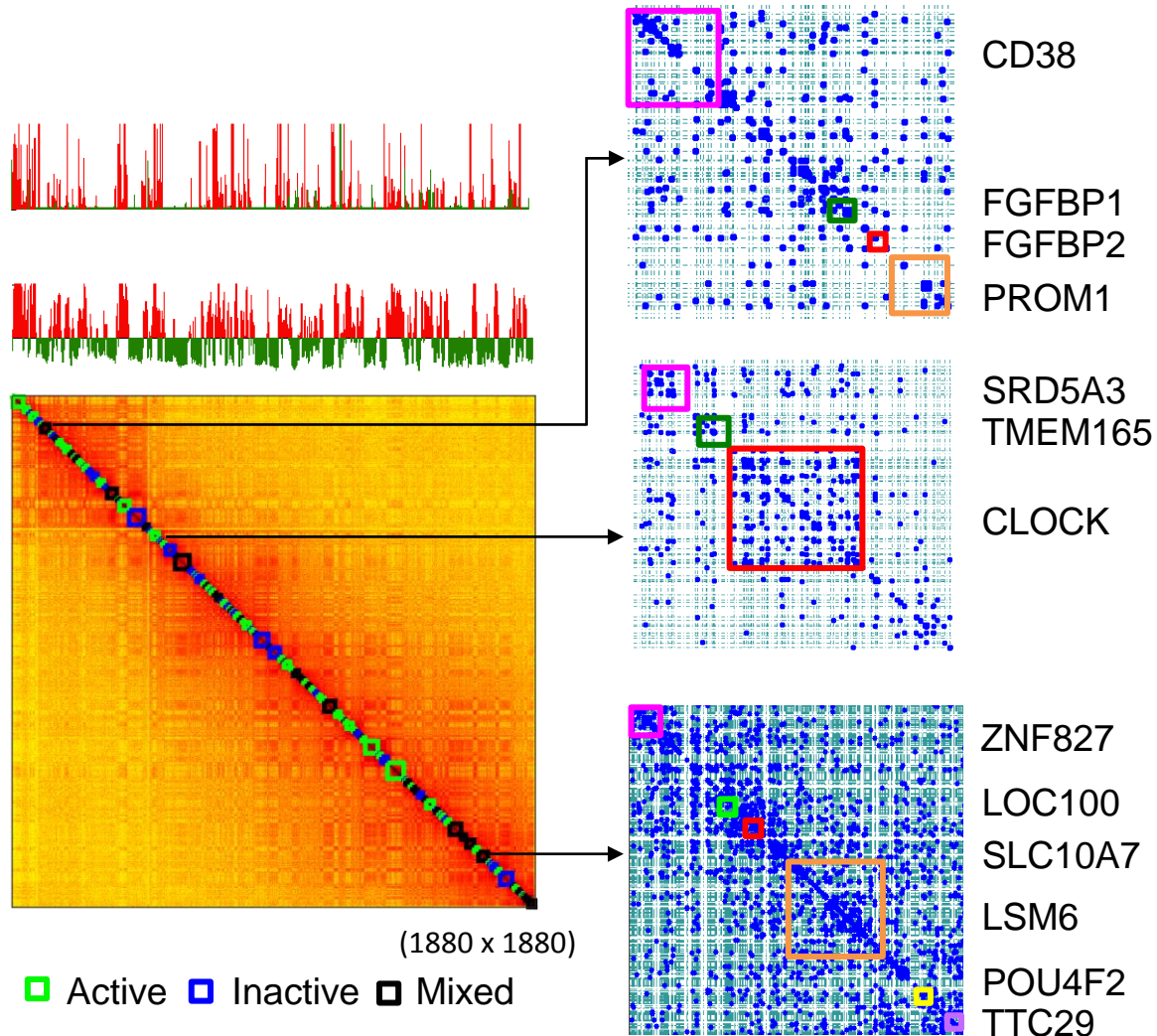
---



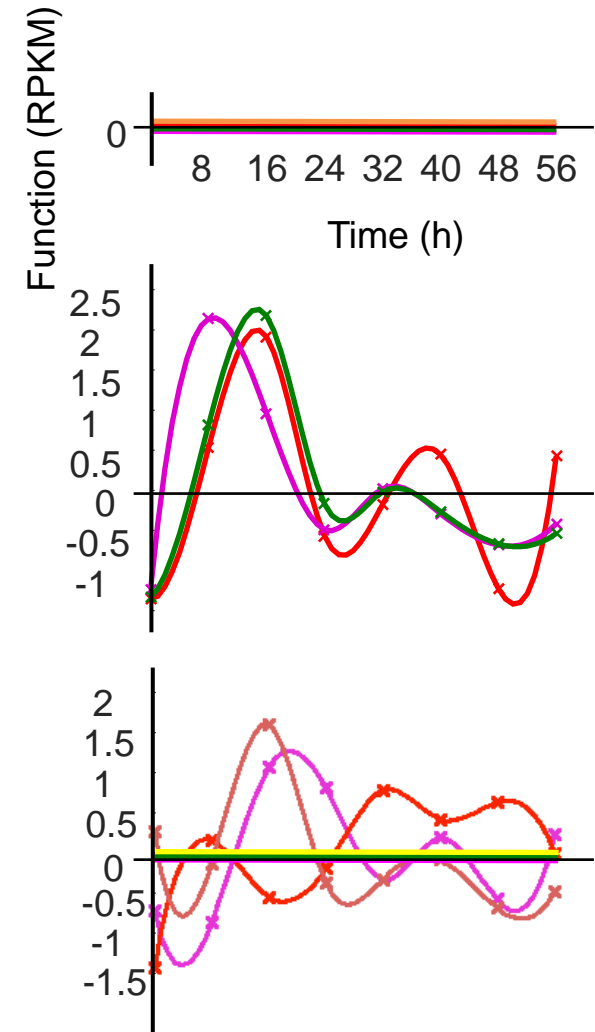
Chen J, Hero A, Rajapakse I. Spectral identification of topological domains, *Bioinformatics*, Published online May 5, 2016

# Data

Chromatin organization (Here only static information is shown)

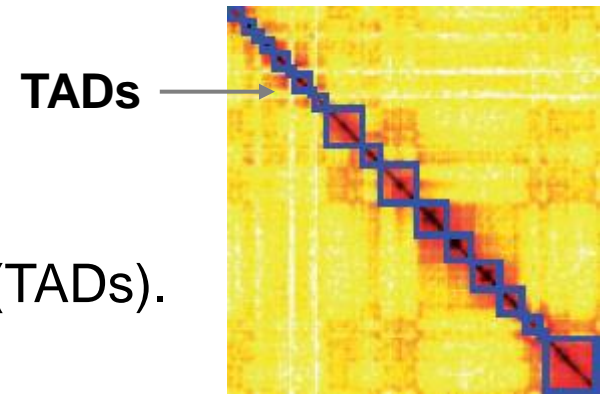


Gene expression over time



# Our methods and data for the algorithm

---



- Identified Topologically Associating Domains (TADs).  
 $\approx 2300 \times 2300$  matrix
- Measured 30,000 RNASeq counts at 8 hr intervals in synchronized cells
- Projected RNASeq counts onto TADs
- Fitted a dynamical model to match the time evolution of the data (2300 equations)
- Used bioinformatics data to infer a value of  **$B$**  for the equation  $A_k x_k + B u_k$  ( $u$  = transcription factors)

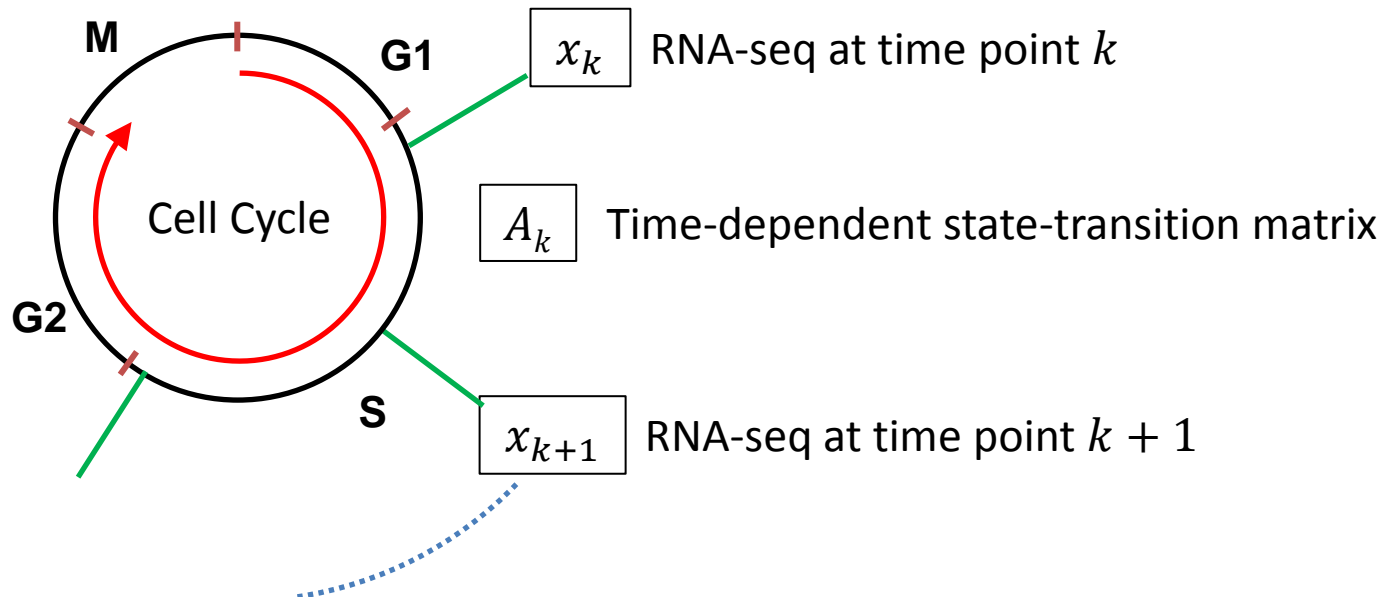
# A Matrix

$$x_{k+1} = A_k x_k + B u_k$$

Natural dynamics:  $x_{k+1} = A_k x_k$

Closest transition matrix to Identity subject to constraints from the data

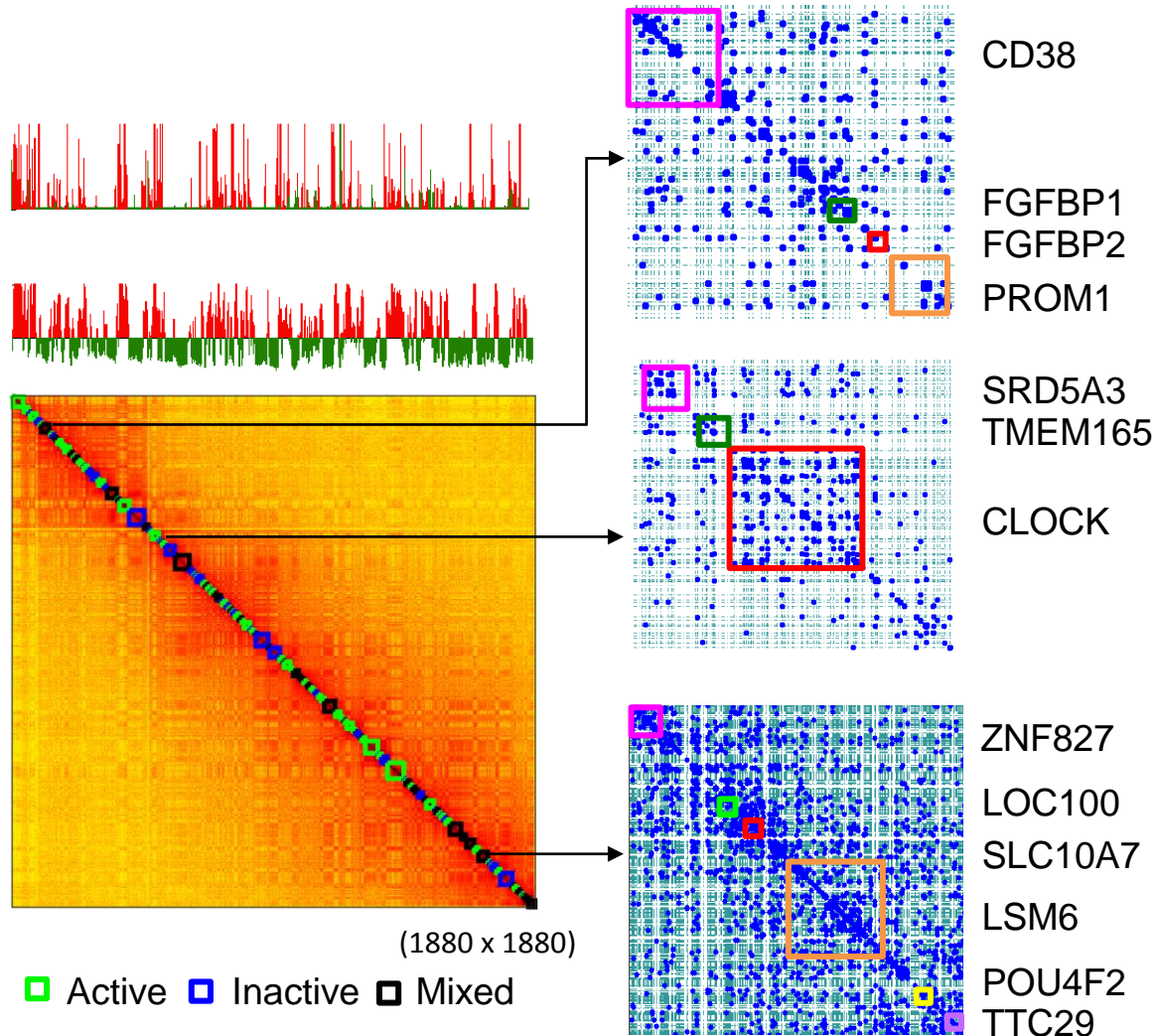
$$A_k = I + \frac{(x_{k+1} - x_k)x_k^T}{x_k^T x_k}$$



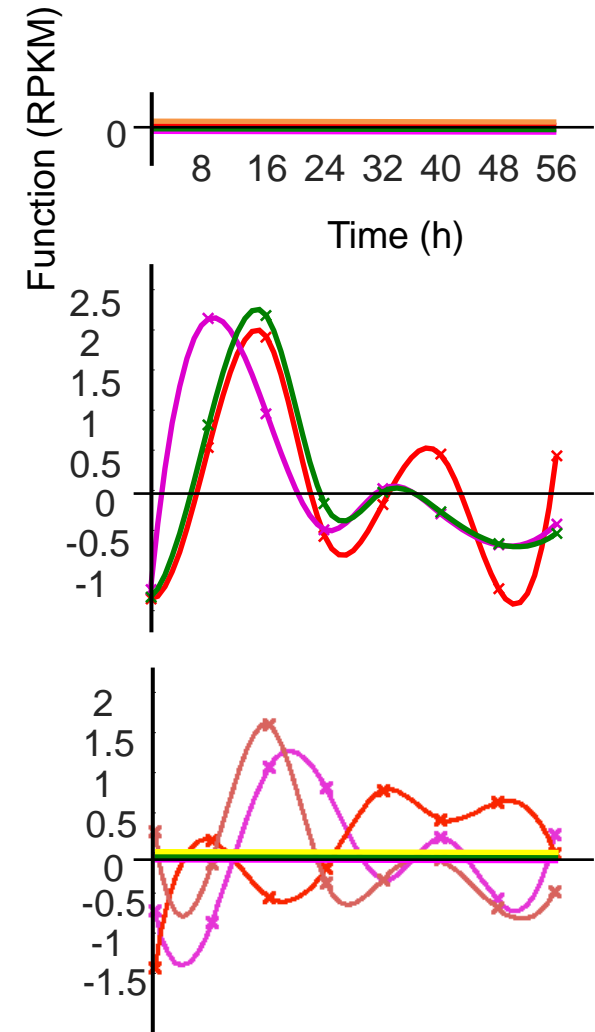


# Data

Chromatin organization (Here only static information is shown)



Gene expression over time



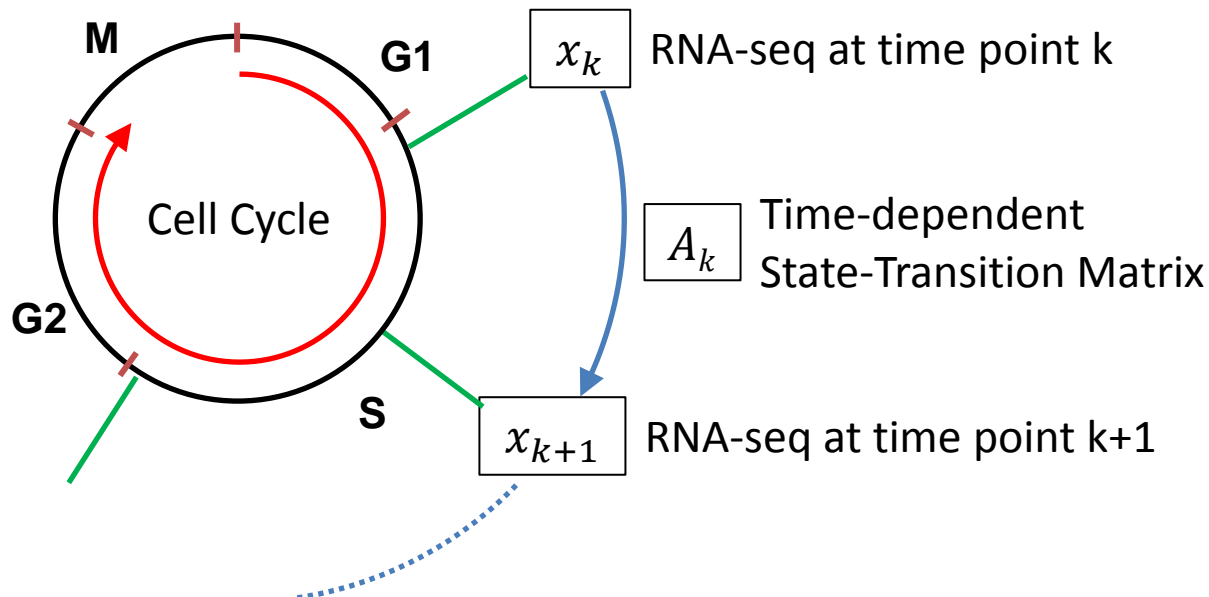
# Controllability – ‘A’ Matrix

$$x_{k+1} = A_k x_k + B u_k$$

No Control Input

$$x_{k+1} = A_k x_k$$

Natural Initial State Dynamics

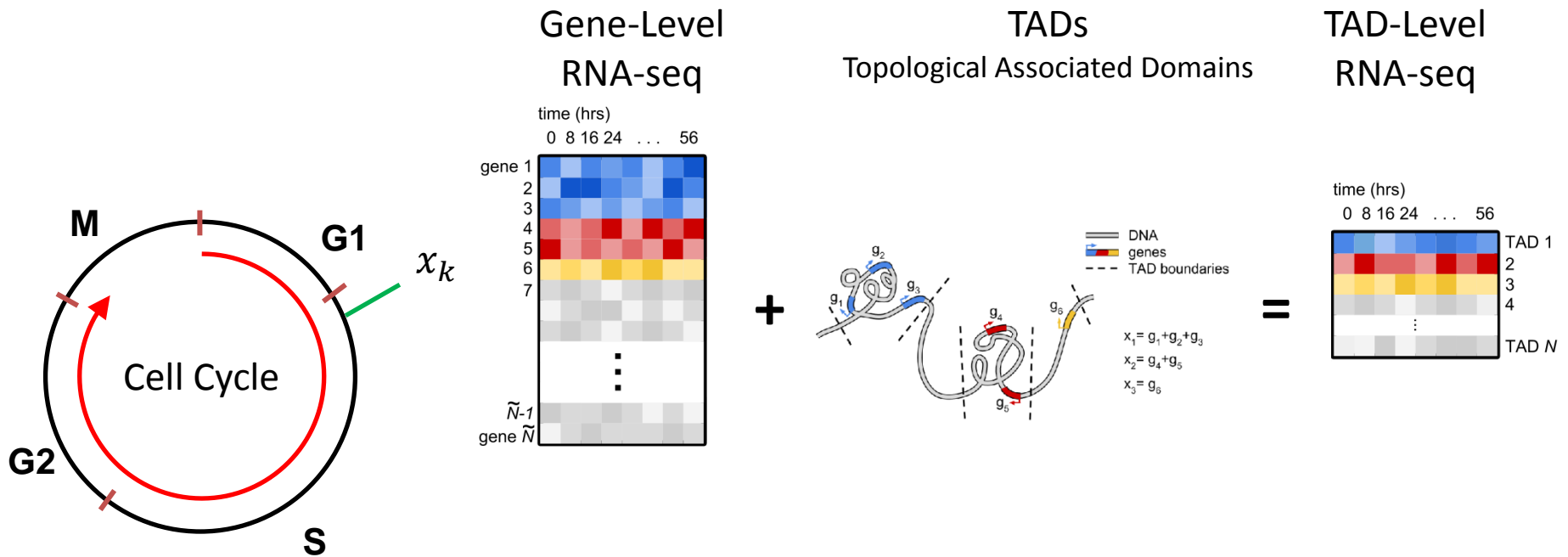


Closest transition matrix  
to Identity

$$A_k = I + \frac{(x_{k+1} - x_k)x_k^T}{x_k^T x_k}$$

# Natural Dimension Reduction

$$x_{k+1} = A_k x_k + B u_k$$

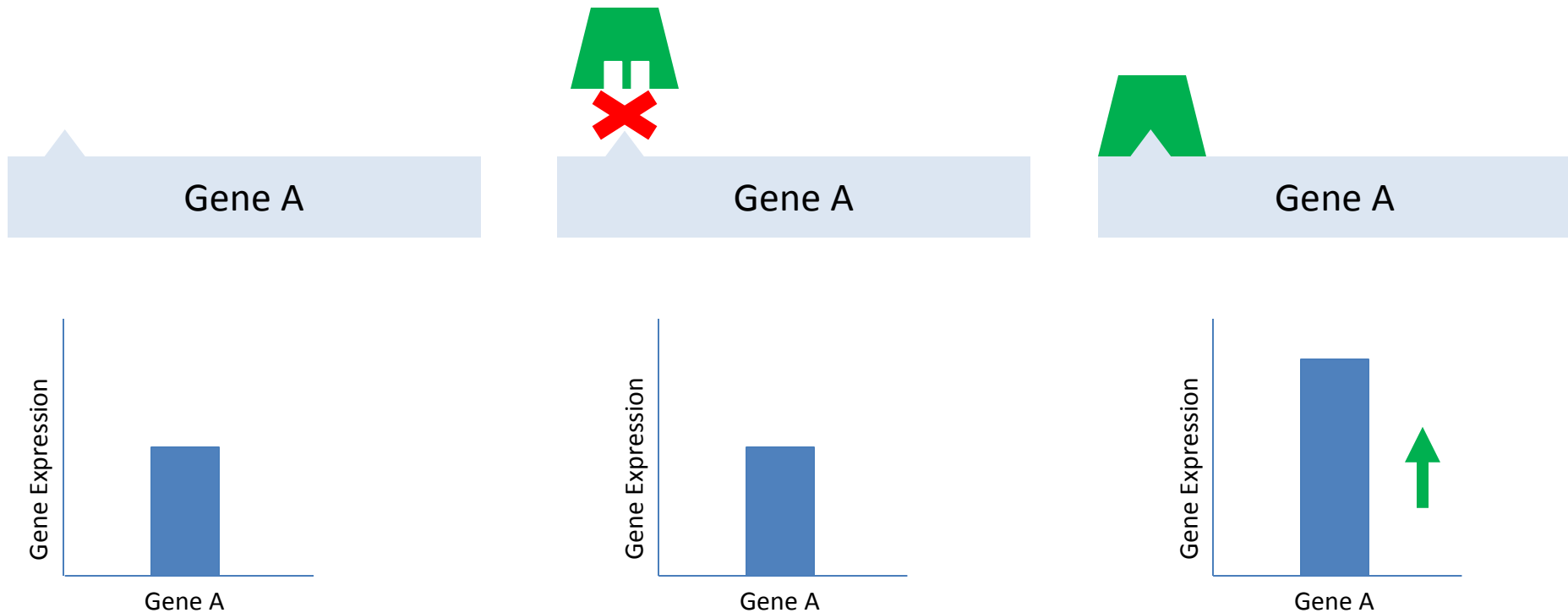


Data + Biology = Dimension-Reduced Data

# Controllability – Control Input

$$x_{k+1} = A_k x_k + Bu_k$$

Transcription factor = Protein that binds to DNA and changes local gene expression

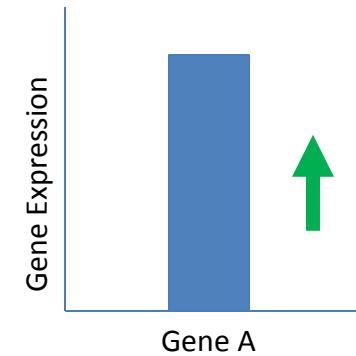
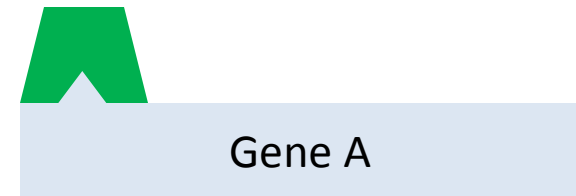
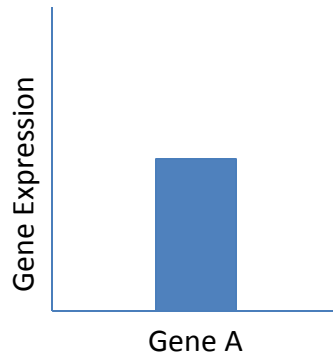
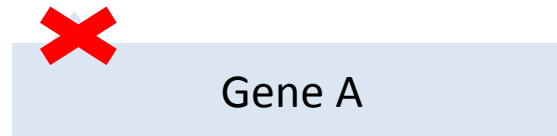
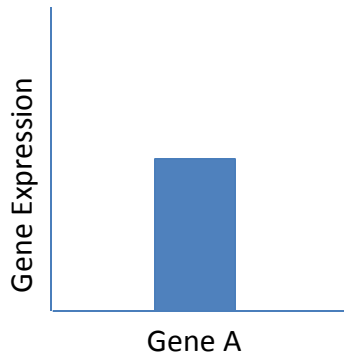
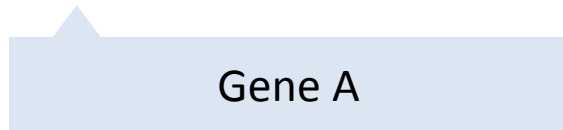




# Controllability – Control Input

$$x_{k+1} = A_k x_k + Bu_k$$

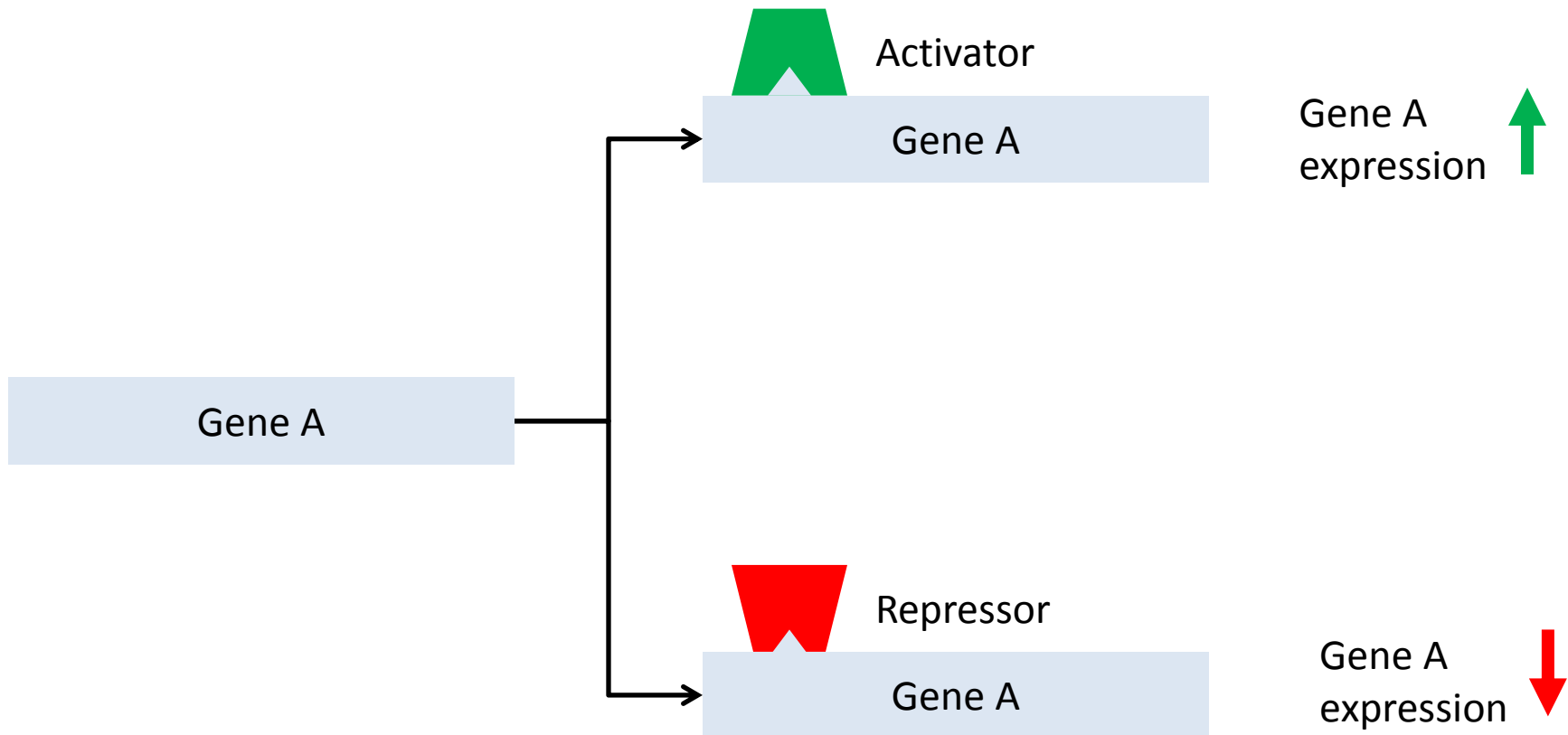
Cannot bind to protein inhibited regions



# Controllability – Control Input

$$x_{k+1} = A_k x_k + Bu_k$$

Transcription factors can be *activators* or *repressors*



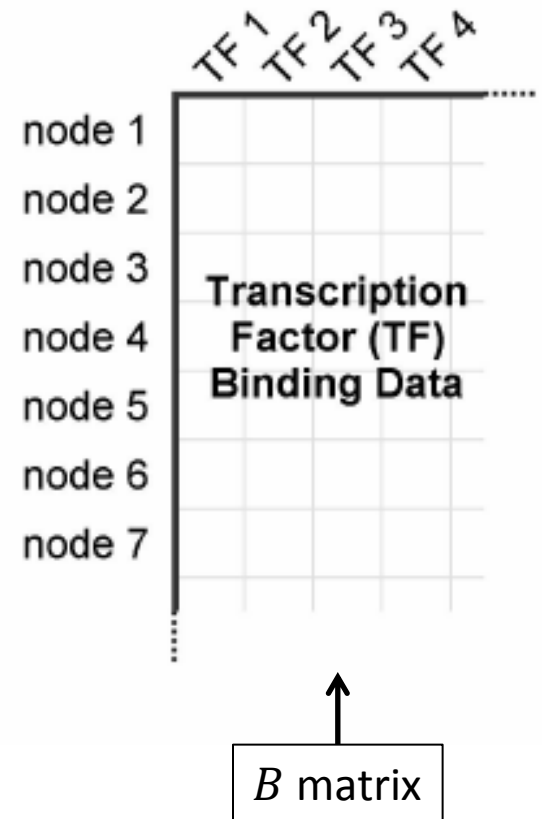
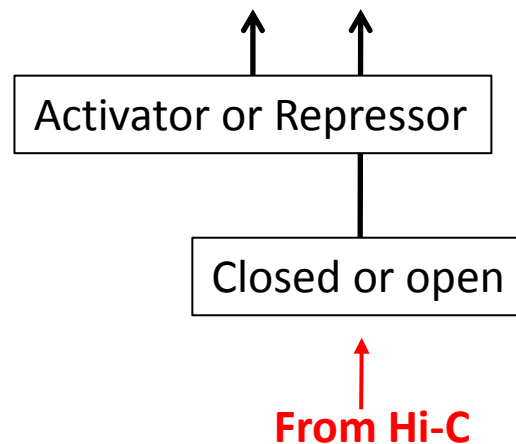
# Controllability – Control Input

$$x_{k+1} = A_k x_k + Bu_k$$

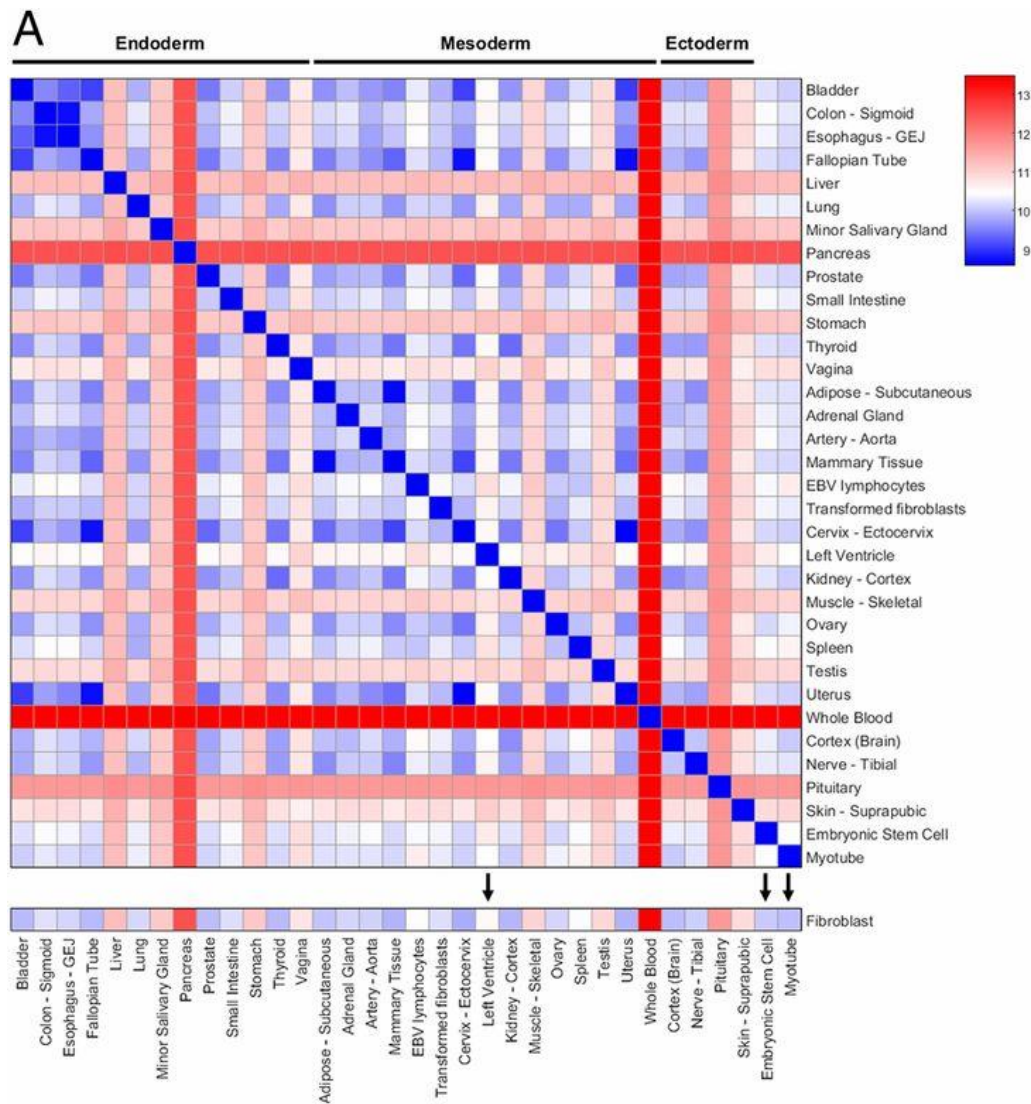
$B = \mathbb{R}^{N \times M}$ , fixed size

$B_{j,m}$  = transcription factor  $m$  affect on TAD  $j$

$B_{j,m} = (\#of\ binding\ sites)(\pm)([0,1])$

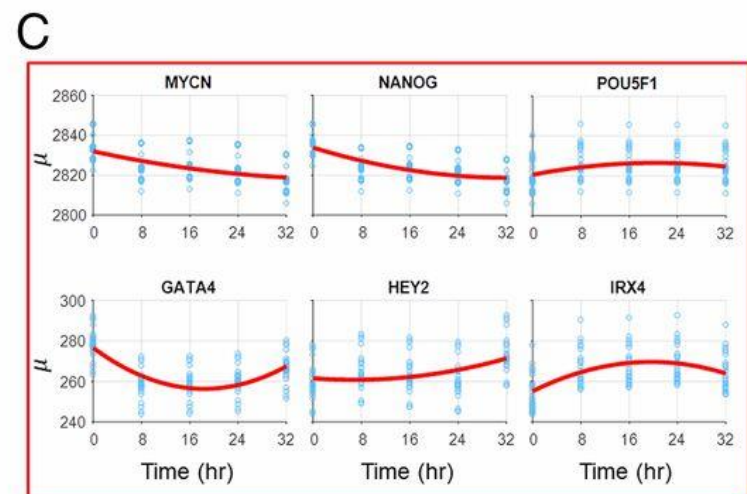


# Quantitative measure between cell types and TF scores



**B**

	ESC	$\mu$	Myotube	$\mu$	Heart - Left Ventricle	$\mu$
1	MYCN	2634	MYOG	1753	HEY2	78
	ZFP42	2453	MYOD1	1065	IRX4	57
	NANOG	2013	PKNOX2	797	HIST2H2BE	44
	OTX2	1628	SIX2	745	HAND1	17
	SOX2	1025	PITX2	744	SOX7	1
2	MYCN, ZFP42	2746	MYOG, PKNOX2	1804	GATA4, HEY2	141
	MYCN, NANOG	2704	MYOG, SIX2	1766	SOX7, IRX4	120
	MYCN, OTX2	2672	MEF2C, MYOG	1753	GATA4, IRX4	117
	MYCN, POU5F1	2634	MYF6, MYOG	1753	HEY2, SOX7	114
	MYCN, SOX2	2634	MYOD1, MYOG	1753	HEY2, IRX4	93
3	MYCN, NANOG, POU5F1	2840	MYOG, SIX2, PKNOX2	1804	GATA4, HEY2, IRX4	263
	MYCN, ZFP42, POU5F1	2801	MEF2C, MYOG, PKNOX2	1804	GATA4, HEY2, HIST2H2BE	218
	MYCN, NANOG, ZFP42	2747	MYF6, MYOG, PKNOX2	1804	HEY2, SOX7, IRX4	190
	MYCN, OTX2, ZFP42	2746	MYOD1, MYOG, PKNOX2	1804	GATA4, HAND1, HEY2	166
	MYCN, SOX2, ZFP42	2746	MYOG, PITX2, PKNOX2	1804	GATA4, HIST2H2BE, IRX4	159





# Controllability – Single Input

---

$$x_{k+1} = A_k x_k + B u_k$$

Assume  $u$  is constant for all  $k$

Assume reprogramming over 1 cell cycle (5 time points)

$$C = A_5 A_4 A_3 A_2 B + A_5 A_4 A_3 B + A_5 A_4 B + A_5 B + B$$

linear non-negative least squares fit

$$\min_u \|x^{target} - A_5 A_4 A_3 A_2 A_1 x_1 - C u\|$$

What value of  $u$  minimizes this equation?

# Controllability – TF Scoring

---

$$x_{k+1} = A_k x_k + B u_k$$

TFs can be scored relative to each other based defined parameter  $\mu$   
(single input, time point 1 addition case shown below)

$$\begin{aligned} \min_u \quad & \|x^{target} - A_5 A_4 A_3 A_2 A_1 x_1 - C u\| \\ \text{subject to} \quad & \begin{cases} u_{m,k} \geq 0, & \text{if } m \in p \\ u_{m,k} = 0, & \text{if } m \notin p \end{cases} \end{aligned}$$

Where  $\hat{p}$  is a set that refers to transcription factor indices we are analyzing

Note: The above equation can be solved through MATLAB *lsqnonneg* function [1].

$$d = \|x^{target} - A_5 A_4 A_3 A_2 A_1 x_1 - C u\|$$

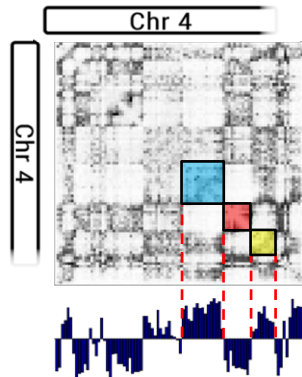
$$d_0 = \|x^{target} - x^{initial}\|$$

$$\mu = d_0 - d$$

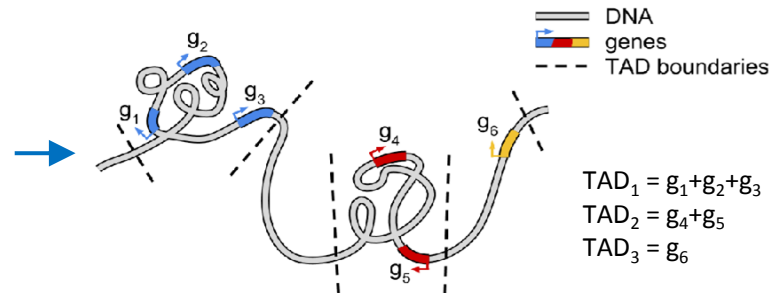
$\mu$  can be thought of as *distance progressed towards target*

# An over view of “Algorithm for cellular reprogramming”

Genome-wide contact map (Hi-C)



Initial partitioning using  
Fiedler vector<sup>1</sup>



TAD: Topologically associating domain, an inherent unit of chromatin organization

Dimension reduction using intrinsic delineations in the genome (TADs)

+ Gene expression

2000 functional units inferred

Mathematics + Bioinformatics  
(TRANSFAC and FIMO databases)

## Skin to Muscle

Rank	TF
1	MyoD
2	PitX2
3	PKNOX2+
4	Six2

Weintraub, 1989

Our algorithm predicted transcription factors (TF) with known reprogramming capability<sup>2</sup>

Prediction: any cell to any cell!

## Skin to Stem Cell

Rank	TF
1	NANOG+OCT4+SOX2
2-11	NANOG+OCT4+other
12	SOX2+OCT4+MYCN

Yamanaka, 2007