

## Cell Fates as High-Dimensional Attractor States of a Complex Gene Regulatory Network

Sui Huang,<sup>1,\*</sup> Gabriel Eichler,<sup>1</sup> Yaneer Bar-Yam,<sup>2</sup> and Donald E. Ingber<sup>1</sup>

<sup>1</sup>*Vascular Biology Program, Departments of Pathology & Surgery, Children's Hospital and Harvard Medical School, Boston, Massachusetts 02115, USA*

<sup>2</sup>*New England Complex Systems Institute, Cambridge, Massachusetts 02138, USA*

(Received 13 September 2004; published 1 April 2005)

Cells in multicellular organisms switch between distinct cell fates, such as proliferation or differentiation into specialized cell types. Genome-wide gene regulatory networks govern this behavior. Theoretical studies of complex networks suggest that they can exhibit ordered (stable) dynamics, raising the possibility that cell fates may represent high-dimensional attractor states. We used gene expression profiling to show that trajectories of neutrophil differentiation converge to a common state from different directions of a 2773-dimensional gene expression state space, providing the first experimental evidence for a high-dimensional stable attractor that represents a distinct cellular phenotype.

DOI: 10.1103/PhysRevLett.94.128701

PACS numbers: 89.75.-k, 05.45.-a, 87.16.Yc, 87.18.-h

The maintenance of distinct phenotypic states of cells in multicellular organisms, such as the states of proliferation or differentiation into specialized cell types, as well as the switchlike transitions between these “cell fates,” are governed by a genome-scale network that consists of 10 000s of genes that regulate each other's activity (i.e., gene expression). This genome-wide regulatory network is identical in essentially every normal cell, and yet it establishes the distinct, stable gene expression profiles associated with the various cell types. Thus, the phenotypic state  $S(t)$  at time  $t$  of a cell may be represented by the activity  $x_i(t)$  of the  $N$  individual genes  $i$  in the genome, i.e.,  $S(t) = [x_1(t), x_2(t), \dots, x_N(t)]$ , a state vector, where  $x_i(t)$  depend on each other as determined by the network of gene regulatory interactions. It remains unclear how such a large and irregular (i.e., “complex”) network [1] can give rise to a globally coherent response that has macroscopic manifestations, such as the distinct, stable cell fates and the conditional cell fate transitions that require coordinated changes in the activity of thousands of genes [2,3]. Thus, if  $S(t)$  represents the dynamics of a complex network of a large number of interdependent variables  $x_i(t)$ , the challenge is to predict the long-term dynamic behavior of  $S(t)$ . Although data on the topology of large biological and nonbiological networks have become available in the past few years, analysis has been largely limited to the characterization of their static topology [4–7], and it has not been possible previously to model or experimentally measure the dynamic behavior of such complex networks due to the lack of appropriate information or tools. Thus, our understanding of the global dynamics is largely based on computer simulations of statistical ensembles of generic complex networks [8,9]. Such theoretical studies suggested that large, random gene networks can, given a particular network architecture, produce “ordered” dynamics, i.e., have relatively few stable attractor states in which a large fraction of the genes remain stationary despite their global interdependence. This theoretical result led to the proposal

that the attractor states represent the various differentiated cell types [8,9], or, more generally, different cell fates [2]. Although plausible theoretically, there is to date no experimental support for the existence of high-dimensional stability in a natural complex molecular network or for the concept that attractors represent stable cell fates.

Unlike in computer models, where one can use numerical simulations to evaluate the state space structure, in real cells we cannot access any arbitrary initial network state in order to monitor its time evolution. To obtain experimental evidence that a given stationary phenotypic state is a high-dimensional attractor state to which a volume of initial states contracts, we exploited the fact that human promyelocytic HL60 cells can reliably be triggered *in vitro* to assume a stable state of neutrophil differentiation by a variety of means. This allowed us to follow different state space trajectories to the differentiated state. Specifically, the solvent dimethylsulfoxide (DMSO) and the hormone all-trans-retinoic acid (atRA) both trigger the same cell fate switch into neutrophils in HL60 cells [10,11]. We used DNA microarrays [12] to monitor genome-wide mRNA steady state levels at various times. These gene expression profiles serve as surrogate measures for configurations of gene activation states, and hence, for  $S(t)$ . This allowed us to determine the two trajectories in the two different neutrophil differentiation processes,  $S^A(t)$  and  $S^D(t)$ , triggered by atRA and DMSO, respectively. By utilizing two biochemically distinct stimuli that are likely to target distinct sets of genes, we provoke the trajectories to initially *diverge* to different regions of the state space. If the final differentiated state is an attractor, it can be approached from different directions of the high-dimensional state space. Thus, the attractor hypothesis predicts that after an initial divergence, the two trajectories  $S^A(t)$  and  $S^D(t)$  will converge to the common end point. The convergence of trajectories from different directions across a large number of gene dimensions is a necessary condition for a high-dimensional attractor state and cannot be easily explained

by the existing notion of a specific, unique “differentiation pathway” as the common target of the two drugs.

Stimulation of HL60 progenitor cells with either DMSO (1.25% v/v) or atRA ( $10^{-7}$  M) resulted in their differentiation into neutrophils within six days as previously reported [10,11] (see the supplementary material in [13] for details). Gene expression profiles across  $\sim 12\,600$  genes were measured for the differentiation processes induced by DMSO and atRA at 0, 2, 4, 8, 12, and 18 h and daily thereafter until day 7 using oligonucleotide DNA microarrays. The relative expression level with reference to that at 0 h was used for  $x_i(t)$ , expressed as the log-transformed ratio of the measured signals:  $x_i(t) = \log_2[\text{signal}_i^{A,D}(t)/\text{signal}_i(t=0)]$ , commonly referred to as “signal-log ratio” (SLR). A set of  $N = 3841$  genes remained after filtering out genes whose expression signal was too low in this cell type to be considered significant or that did not exhibit a

significant change in expression during the entire course of the experiment.

Unlike the use of DNA microarrays to identify specific genes, we treated genes as anonymous members of a single ensemble containing  $N$  genes and calculated the intertrajectory distance  $b(t)$  between  $S^A(t)$  and  $S^D(t)$  at corresponding time points. This ensemble property of the population of genes is a robust measure that is not biased by noise at the level of individual gene measurements.  $b(t)$  was quantified as the “inverse correlation,”  $b(t) = 1 - r(t)$ , where  $r(t)$  is the Pearson coefficient of correlation between the two state vectors  $S^A(t)$ ,  $S^D(t)$  at time  $t$ .

We first examined the ultimate convergence of the atRA- and DMSO-induced neutrophils, i.e., whether there is a negligible disparity at day 7 between  $S_A(7\text{ d})$  and  $S_D(7\text{ d})$  as expected based on functional comparison. The state vector disparity for microarray replicates was experimentally determined to be  $b_{\text{replicate}} < 0.01$  in three separate hybridizations. The disparity between different microarray samples within the same treatment group measured at different days after the cells had reached the stationary state (days 6 vs 7 and 11 vs 12; for both processes), and hence an upper bound estimate of intersample variability, was also low ( $b_{\text{stat}} = 0.14 \pm 0.02$ ). We found that at day 7 the final disparity was  $b(7\text{ d}) = 0.42$ . Because  $b(7\text{ d}) > b_{\text{stat}}$ , this final disparity cannot be explained by measurement noise alone. The lack of equivalence of the trajectories at the end of differentiation indicates that DMSO and atRA-induced neutrophils are not entirely identical, as previously suggested on biological grounds [14]. Thus, there are some genes whose expression may not be relevant to the macroscopic definition of neutrophils and are differentially regulated by DMSO and atRA. For the analysis of the trajectory we therefore focused on the subset of genes whose expression in atRA vs DMSO-induced neutrophils (at day 7) was considered to be not significantly different (see the supplementary material in [13]). With this filter we obtained a subset of  $N = 2773$  genes ( $=72\%$  of the initial set of  $N = 3841$  genes).

Since it is difficult to display thousands of dimensions (genes) for both processes simultaneously, we first reduced the dimensionality by using the gene expression dynamics inspector (GEDI) program [15] and principal component analysis. GEDI enables the visualization and comparison of multiple time series by mapping each expression profile, i.e., snapshot of  $S(t)$ , into a “mosaic” representation through dimension reduction and reordering of the genes into miniclusters of typically 1–10 similarly behaving genes using a self-organizing map. The GEDI mosaics [Fig. 1(a)] revealed that starting with an identical expression pattern, the two processes exhibit clearly distinct genome-wide gene expression patterns by 12–18 h after treatment with DMSO and atRA, as indicated by different color patterns. After this initial divergence the GEDI mosaics converged to a virtually identical pattern by day 6.

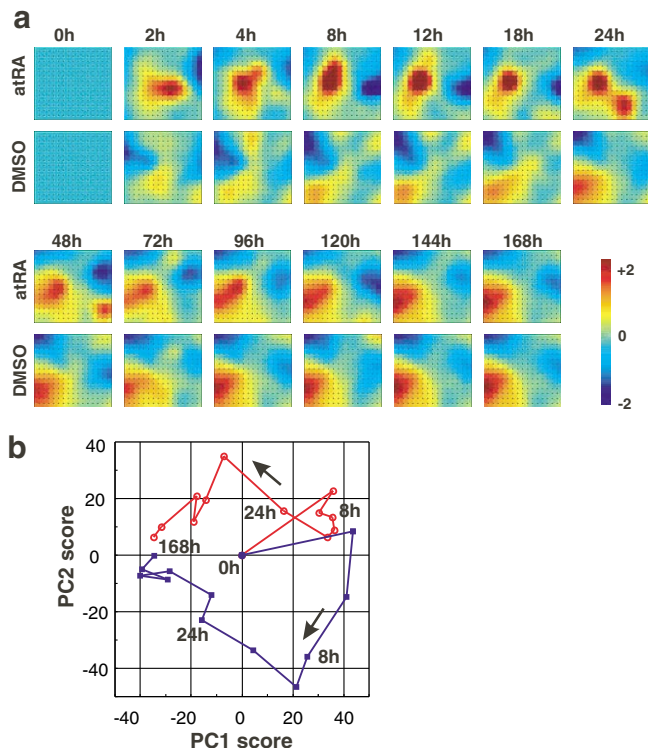


FIG. 1 (color). Comparison of the two gene expression trajectories for the subset of  $N = 2773$  genes during neutrophil differentiation. (a) The genes were clustered by a self-organizing map into  $15 \times 16$  “miniclusters” with regard to their temporal profiles across both differentiation processes using the GEDI program [15]. Each minicluster is mapped onto the same corresponding “tile” in all the “mosaics,” each of which represents a snapshot of  $S(t)$ . Tile colors indicate the expression level of the cluster centroid; numbers on color bar: gene expression levels in SLR units. (b) Principal component analysis. Each point represents an individual expression profile  $S(t)$  within one of the two differentiation processes (red circles: atRA; blue squares: DMSO) projected onto the first two principal components (PC1 and PC2).

Importantly, the genome-wide changes of the patterns in GEDI show that the convergence occurred with respect to a large portion of the genome, i.e., to a high number of state space dimensions.

To more formally display the trajectories, we next mapped the profiles into a reduced dimension space spanned by principal components (PCs) of the gene expression data with respect to variations in gene dimension. Figure 1(b) shows the trajectory in the reduced “state space” spanned by the two primary principal components, PC1 and PC2, which explain 31% and 14% of the variation in the expression data, respectively. In this projection, the two trajectories diverge and then converge at the end of differentiation. PC1 appears to account for the expression change due to differentiation, while PC2 reflects the difference between the two trajectories.

Finally, we quantified the time evolution of the intertrajectory disparity using the aggregate variable  $b(t)$ . Figure 2(a) shows that for the converging subset of  $N = 2773$  genes (solid symbols) the disparity increased initially to  $b(t) = 0.80$  (i.e., there is almost no correlation between the two trajectories) on day 1, and then decreased to 0.14, a level indistinguishable from experimental variability  $b_{\text{stat}}$  (see above). Significantly, even if  $b(t)$  for the raw data for all 3841 genes were plotted [empty symbols, Fig. 2(a)], the initial divergence and terminal convergence remained, except that the trajectories did not fully converge at day 7 as mentioned earlier. A nearly identical partial convergence of the trajectories was observed when gene expression profiles were monitored in a separate, independent repeat experiment using an entirely different (spotted cDNA filter) array technology [Fig. 2(b)]. Thus, the time course of

the aggregate variable  $b(t)$  behaves as if thousands of genes in the complex network exhibit a globally coherent dynamic pattern of attraction to a common stable state.

An alternative explanation for the transient divergence of the two trajectories is that they move along the same path, but at different rates. However, we found more than 100 genes, including the neutrophil differentiation markers CD11b and G-CSFR1, that exhibited essentially equivalent dynamics in the atRA- and DMSO-triggered processes, hence excluding an overall temporal shift in differentiation. The convergence of the trajectories  $S^A(t)$  and  $S^D(t)$  defined by the 2773 genes is consistent with the transition into a high-dimensional attractor state with respect to this large portion (72%) of the gene network. However, this does not imply that the two observed trajectories are an unbiased sample of trajectories near the terminus since kinetics may play a role [16] and thus, state space properties near the attractor cannot be inferred directly. The number of time points monitored in the experiments also does not allow a detailed state space analysis as in other high-dimensional dynamic systems [17].

In the present study, we showed that even in the absence of knowledge of the specific network architecture, it is possible to use genome-wide gene expression profiling to probe the state space structure of a natural complex network and extract characteristic signatures of a stable high-dimensional attractor. It is not at all obvious that such stable behavior should arise from the interaction of a large number of irregularly connected elements [18,19]. But an important result from the analysis of statistical ensembles of discrete genetic networks is that given some global network architectural features, a complex network will spontaneously produce globally coherent patterns of gene activation, i.e., quickly settle down in one of a relatively small set of stable attractors instead of eventually visiting the entire state space [8,9]. The architectural features of the network known to enlarge the regime of ordered behavior include (i) sparseness of interactions [9], (ii) preferential use of a certain subset of functions for the regulatory interactions between genes [9,20], and (iii) a “scale-free” topology [21,22]. Sparseness, (near) scale-free architecture [5,23–25], and a bias in the occurrence of regulatory functions [26] have all been found to be predominant in real gene or protein regulatory networks.

The existence of discrete, mutually exclusive cell states has long been suggested to reflect multistability in small regulatory circuits comprised of two or a few elements that arise due to nonlinear regulatory relationships [27–29], and bistability in local signaling modules has recently been verified [30,31]. However, genomic and proteomic analysis of molecular networks, as well as graph theory principles, indicate that molecular interactions in the cell typically form a single, large connected network (giant component) that may span 90% of the genome [5,32,33]. Indeed, given the mutual exclusivity of cell fate programs

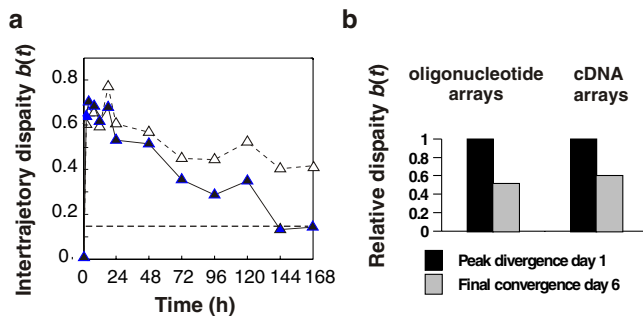


FIG. 2 (color online). Time evolution of the state space distance between the two trajectories for cells treated with atRA versus DMSO. (a) Time course of disparity  $b(t)$ . Open symbols: entire set of  $N = 3841$  genes; solid symbols: subset of  $N = 2773$  genes with common expression level in the atRA- and DMSO-differentiated neutrophils. (b) Comparison of two different methods for gene expression profiling. Oligonucleotide arrays: data obtained with Affymetrix microarrays ( $N = 3841$  genes); cDNA arrays: data from independent experiments using low-density cDNA arrays ( $N = 453$  genes).  $b(t)$  represents the disparity between the profiles in the two treatments at day 6 relative to the peak value in day 1.

[2,34], the action of postulated local “pathway modules” thought to serve individual cell functions must also somehow be globally coordinated. Moreover, a whole-genome view is also necessary because cells receive a broad range of simultaneous biological signals, as well as nonspecific (chemical or mechanical) perturbations [2,35] which influence genes across the entire genome—yet they reliably integrate all these inputs and select only one of a few possible cell fates [2].

Thus, formal network architecture considerations as well as experimental observation of cell fate behavior also support the idea that the genome-scale regulatory network can act as an integrated entity and give rise to coherent, higher-order dynamic patterns, such as stable high-dimensional attractors.

We thank T. Golub, K. Stegmaier, H. Radomska, and J. Wikswa for helpful discussions, L. Kunkel for access to the microarray facility, M. Han for technical assistance, and R. Shamberger and J. Folkman for their continued support. This work was supported by the AFOSR (S. H.), by NIH (D. E. I.), by NASA (D. E. I.), and by NSF (Y. B.).

---

\*Electronic address: sui.huang@childrens.harvard.edu

- [1] S. H. Strogatz, *Nature (London)* **410**, 268 (2001).
- [2] S. Huang, in *Gene Regulation and Metabolism: Post-Genomic Computational Approach* (MIT Press, Cambridge, MA, 2002), pp. 181–220.
- [3] Y. Bar-Yam and I. R. Epstein, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 4341 (2004).
- [4] A. L. Barabasi and R. Albert, *Science* **286**, 509 (1999).
- [5] H. Jeong *et al.*, *Nature (London)* **411**, 41 (2001).
- [6] S. Maslov and K. Sneppen, *Science* **296**, 910 (2002).
- [7] L. A. Amaral *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **97**, 11 149 (2000).
- [8] S. A. Kauffman, *J. Theor. Biol.* **22**, 437 (1969).
- [9] S. A. Kauffman, *The Origins of Order* (Oxford University Press, New York, 1993).
- [10] S. J. Collins *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **75**, 2458 (1978).
- [11] T. R. Breitman, S. E. Selonick, and S. J. Collins, *Proc. Natl. Acad. Sci. U.S.A.* **77**, 2936 (1980).
- [12] A. J. Holloway *et al.*, *Nat. Genet. Suppl.* **32**, 481 (2002).
- [13] See EPAPS Document No. E-PRLTAO-94-003513 for a description of the experimental methods and analysis. A direct link to this document may be found in the online article’s HTML reference section. The document may also be reached via the EPAPS homepage (<http://www.aip.org/pubservs/epaps.html>) or from <ftp.aip.org> in the directory/epaps/. See the EPAPS homepage for more information.
- [14] S. J. Collins, *Blood* **70**, 1233 (1987).
- [15] G. S. Eichler, S. Huang, and D. E. Ingber, *Bioinformatics* **19**, 2321 (2003).
- [16] G. E. Crooks, B. Ostrovsky, and Y. Bar-Yam, *Phys. Rev. E* **60**, 4559 (1999).
- [17] B. Gerstman and Y. Garbourg, *J. Polym. Sci., Part B: Polym. Phys.* **36**, 2761 (1998).
- [18] R. M. May, *Nature (London)* **238**, 413 (1972).
- [19] D. A. Meyer and T. A. Brown, *Phys. Rev. Lett.* **81**, 1718 (1998).
- [20] I. Shmulevich *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 10 734 (2003).
- [21] J. J. Fox and C. C. Hill, *Chaos* **11**, 809 (2001).
- [22] M. Aldana and P. Cluzel, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 8710 (2003).
- [23] H. W. Mewes *et al.*, *Nucleic Acids Res.* **30**, 31 (2002).
- [24] L. Giot *et al.*, *Science* **302**, 1727 (2003).
- [25] S. Li *et al.*, *Science* **303**, 540 (2004).
- [26] S. E. Harris *et al.*, *Complexity* **7**, 23 (2002).
- [27] M. Delbrück, in *Unités Biologiques Douées de Continuité Génétique Colloques Internationaux du Centre National de la Recherche Scientifique* (CNRS, Paris, 1949).
- [28] A. Novick and M. Weiner, *Proc. Natl. Acad. Sci. U.S.A.* **43**, 553 (1957).
- [29] J. Monod and F. Jacob, *Cold Spring Harb. Symp. Quant. Biol.* **26**, 389 (1961).
- [30] W. Xiong and J. E. Ferrell, Jr., *Nature (London)* **426**, 460 (2003).
- [31] E. M. Ozbudak *et al.*, *Nature (London)* **427**, 737 (2004).
- [32] E. M. Marcotte, *Nat. Biotechnol.* **19**, 626 (2001).
- [33] D. S. Callaway, J. E. Hopcroft, J. M. Kleinberg *et al.*, *Phys. Rev. E* **64**, 041902 (2001).
- [34] C. H. Waddington, *Principles of Embryology* (Allen and Unwin, London, 1956).
- [35] S. Huang and D. E. Ingber, *Exp. Cell Res.* **261**, 91 (2000).