

Stochastic simulations of the origins and implications of long-tailed distributions in gene expression

Sandeep Krishna*, Bidisha Banerjee*, T. V. Ramakrishnan^{†‡}, and G. V. Shivashankar*^{§¶}

*National Centre for Biological Sciences, Tata Institute for Fundamental Research, Bangalore 560065, India; [†]Banaras Hindu University, Varanasi 221005, India; [‡]Department of Physics, Indian Institute of Science, Bangalore 560012, India; and [§]Raman Research Institute, Bangalore 560080, India

Edited by Charles R. Cantor, Sequenom, Inc., San Diego, CA, and approved February 11, 2005 (received for review August 31, 2004)

Gene expression noise results in protein number distributions ranging from long-tailed to Gaussian. We show how long-tailed distributions arise from a stochastic model of the constituent chemical reactions and suggest that, in conjunction with cooperative switches, they lead to more sensitive selection of a subpopulation of cells with high protein number than is possible with Gaussian distributions. Single-cell-tracking experiments are presented to validate some of the assumptions of the stochastic simulations. We also examine the effect of DNA looping on the shape of protein distributions. We further show that when switches are incorporated in the regulation of a gene via a feedback loop, the distributions can become bimodal. This might explain the bimodal distribution of certain morphogens during early embryogenesis.

fluctuations | genetic switches | single cell

The inevitable noise in gene expression, manifested at the subcellular level as distributions in protein numbers, has been observed experimentally in both prokaryotes and eukaryotes (1–6). Recent studies have investigated how organisms tolerate this noise and the kinds of regulatory strategies they use to control or minimize it (7, 8). One example where it has been suggested that noise is exploited for the benefit of the organism is bacterial chemotaxis (9). Analyses of noise in gene expression have highlighted the analogy with quantum many-body systems (10), and the authors of refs. 11–15 have explored the contribution of intrinsic and extrinsic sources, as well as the relative contribution of transcription and translation, focusing on the standard deviation of protein fluctuations. If the protein distributions were Gaussian, the mean, μ , and standard deviation, σ , would provide a complete description of the noise characteristics. However, recent experiments have revealed that protein distributions are often non-Gaussian and also time-dependent, showing a crossover from long-tailed to Gaussian (3). It is important, therefore, to understand the origins and implications of the long-tailed nature of the protein distributions.

We implement a stochastic chemical model of gene expression and show that it leads to distributions that fit the experimental observations, presented in this paper and in ref. 3, of protein distributions at different stages of bacterial growth. The predictions of the simulations were experimentally tested by single-cell-tracking experiments. We suggest that long-tailed protein distributions filtered by appropriate switches can lead to selection at the subcellular level. For example, a switch with a sharp threshold can be used to select a subpopulation of cells with a large concentration of a particular protein from a population of cells with a long-tailed distribution of that protein. By contrast, we show that symmetric Gaussian distributions are not as sensitive to switches.

When long-tailed distributions are combined with switches via a positive feedback loop, we find that the system becomes bistable and can result in bimodal protein distributions. Recent observations of a bimodal protein distribution in an autoregulatory system (3) and engineered gene circuits that couple a switch to a GFP gene (16) confirm our predictions. We also find bimodality in the distribution of the hunchback protein in early-stage *Drosophila*

embryos and suggest that this could be the result of the bicoid switch acting on a long-tailed distribution. Thus, long-tailed distributions and their response to switches might also be of relevance to processes in early embryogenesis. A schematic of our studies is presented in Fig. 1.

Origins of Long-Tailed Protein Distributions

Stochastic Simulations of a Chemical Model of Gene Expression. To understand the microscopic processes producing the long-tailed distributions and the importance of different sources of noise, we have constructed a chemical model of the process of gene expression and analyzed it by using stochastic simulations. The model describes a single gene regulated by an operator site where a repressor molecule can bind and prevent transcription initiation. Expression of the gene is modeled as a series of chemical reactions, as in ref. 17. The Gillespie method (18) is used for stochastic simulation of the reactions. In the Gillespie method, the probability per unit time for the occurrence of a reaction is taken to be the product of a combinatorial factor, which is a function of the numbers of reactants, and the rate constant of the reaction. For second- and higher-order reactions, the volume of the cell also has to be taken into account and, here, is assumed to be linearly increasing through the cell cycle. Each cell cycle, of duration T , is implemented as follows. First, the Gillespie simulation is run for a time t_D , at which point the gene copy number, n , is doubled. The simulation is then run until time T when the cell volume and gene copy number are halved and other molecules are partitioned binomially. The process is then iterated for the next cell cycle, with one daughter cell followed after each partitioning (the model and the simulation algorithm are described in *Stochastic Simulation of a Chemical Model of Gene Expression* in Supporting Text, which is published as supporting information on the PNAS web site; they are based on models and algorithms used in refs. 13, 17, and 19). Fig. 2 shows the results of simulations with a fixed cell division time, $T = 1,800$ s, but different values of the total repressor number, R . The effects of different sources of noise on the protein time series are evident in these runs. When the protein number is large, the thermal noise (reflected in the stochastic occurrence of chemical reactions) is small but the noise from partitioning of molecules during cell division is visible. Thus, in the $R = 100$ and $R = 300$ runs, in the steady state, the protein number grows (on average) exponentially over each cell cycle: $n = N_0 e^{\alpha t}$, where $\alpha = \ln(2)/T$. In contrast, when the protein number is small (the $R = 10,000$ runs), the thermal noise becomes important and dominates the stochastic features of the time series. These aspects are clearly observed in our experiments discussed later.

The protein distributions corresponding to the runs of Fig. 2 are displayed in Fig. 3. The distributions are skewed and long-tailed for smaller mean protein numbers and closer to Gaussian for the run with the largest mean protein number. These distributions compare well with the experimentally observed distributions also shown in

This paper was submitted directly (Track II) to the PNAS office.

[¶]To whom correspondence should be addressed. E-mail: shiva@ncbs.res.in.

© 2005 by The National Academy of Sciences of the USA

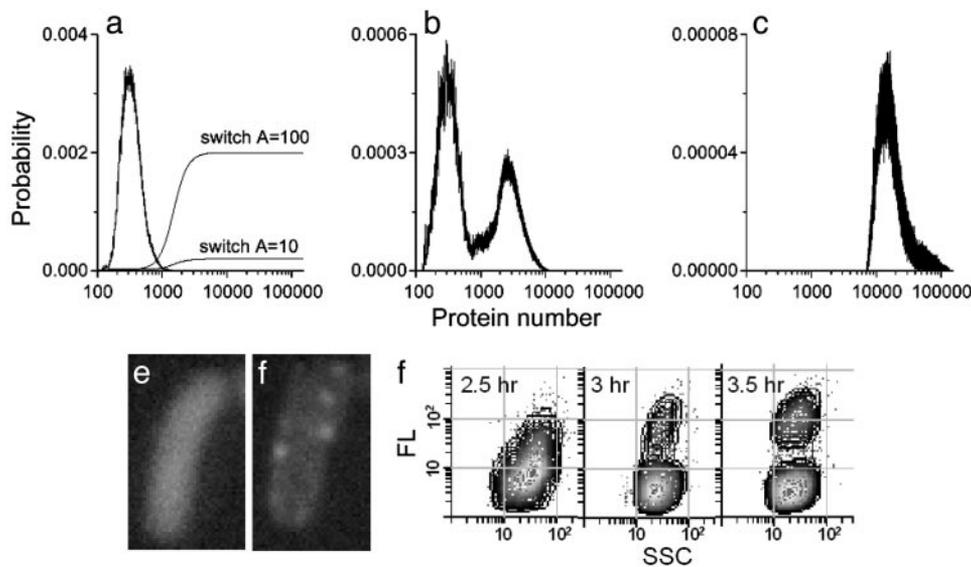


Fig. 7. Response of protein distributions to regulatory switches. (a) Long-tailed protein distribution from a simulation run, superimposed with the response functions of a switch with a Hill coefficient of 4, a threshold of 1,500, and amplitudes of 10 and 100, respectively (for clarity of display the functions have been scaled by a factor of 2×10^{-5}). (b) The modified protein distribution when the switch with $A = 10$ is added to the model (note the change in y scale). (c) The modified protein distribution when the switch with $A = 100$ is added to the model (note the change in y scale). (d and e) Fluorescence image of a cell under autoregulation at two different time points. d shows a uniform distribution of fluorescence, whereas e shows the same cell having bright spots at different places besides a uniform haze. (f) Scatter plot of fluorescence signal (FL) vs. side scatter (SSC) showing emergence of bimodality in the fluorescence distribution at 3.5 h, from a unimodal distribution at 2.5 h. The SSC distribution remains unimodal throughout.

served (see *Bistability due to Positive Feedback via a Switch in Supporting Text*). Bistability due to a switch added in a positive feedback loop has been observed in a number of systems, ranging from bacteriophages (22) to prokaryotes (23) to eukaryotes (24). Reviews of the requirements for constructing stable bistable switches can be found in refs. 25 and 26.

The Effect of DNA Looping. Regulation of transcription initiation by an operator site that is several hundred base pairs upstream of the promoter is common in eukaryotic systems (27). It has been suggested that this scheme of regulation, with DNA looping, can achieve repression levels far higher than regulation with a single operator site near the promoter (28). We have extended our stochastic simulations to include reaction schemes with DNA looping using an effective increase in the local concentration of the repressor molecule that is bound to the far upstream operator site (29) (see *DNA Looping in Supporting Text*). In these simulations, we find that the repression level is indeed enhanced, resulting in a much lower mean protein number for the same number of repressors, if the upstream site is sufficiently far from the main operator site. Comparing protein distributions with and without looping, having the same mean protein number, we find that looping results in a longer tailed protein distribution (see *DNA Looping in Supporting Text*).

We have found that combining DNA looping with a switch also results in sensitive selection of subpopulations of cells with higher protein concentrations. Because of the longer tails, the region of bimodality is much smaller for regulation with DNA looping as compared with operator regulation (see *Bistability due to Positive Feedback via a Switch in Supporting Text*). Thus, it is conceivable that different mechanisms for regulating transcription initiation, such as enhancer looping, might have been selected in various organisms for the kind of sensitivity they show in response to switches.

Bicoid and Hunchback in Early-Stage *Drosophila* Embryos. Switches with sharp thresholds are known to play an important role in embryogenesis. In early-stage *Drosophila* embryos, the maternal

morphogen bicoid acts as a switch that causes the hunchback gene to express in regions of the embryo where the bicoid concentration is above a critical threshold (30, 31). We have examined the expression of hunchback in a *Drosophila* embryo at cycle 14, taken from the FlyEx database (<http://urchin.spbcas.ru/flyex>). The database provides fluorescence images of the embryo showing the amounts of bicoid, hunchback, and other proteins in different parts of the embryo, as well as quantitative data of average fluorescence intensities for each protein, for each nucleus seen in the images (32, 33). We have used the quantitative data from the database to construct the histogram of hunchback concentration in the anterior portion of each embryo (see *Bicoid and Hunchback in Early-Stage *Drosophila* Embryos in Supporting Text*). Fig. 8 shows the histogram obtained for the embryo named hx21. Note that this histogram has

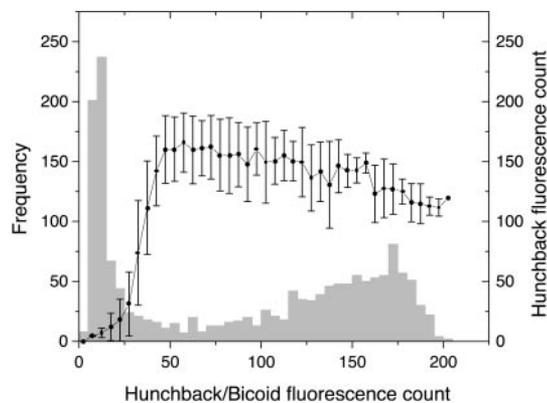


Fig. 8. Light gray bars, histogram of hunchback fluorescence for anterior portion of embryo hx21 from the FlyEx database. Only the 1,606 nuclei with anterior–posterior coordinate between 10% and 70% of the egg length have been considered to create this histogram. Filled circles, an estimate of the response function of the bicoid switch. Each data point shows the mean hunchback fluorescence for nuclei in the corresponding bin (width of five counts) of bicoid fluorescence. Error bars are one standard deviation.

been constructed from nuclei at various locations in the embryo and, hence, does not contain any spatial information. Also superposed on the figure is an estimate of the response function of the bicoid switch. We suggest that a long-tailed distribution of hunchback could be modified by the bicoid switch to produce the bimodal distribution of Fig. 8 in exactly the same way that the distribution of Fig. 7a is modified by the switch to produce the bimodal distribution of Fig. 7b.

Conclusions

At the microscopic level, long-tailed distributions can occur in, broadly, two ways. One is if the microscopic processes are inherently non-Gaussian; for instance, if there is multiplicative noise in the system, or if the mean numbers are small and there is a lower cutoff. This is the case for the long-tailed distributions of conductance observed in one-dimensional wires, where it is a consequence of the localization of electrons due to disorder. A second way in which long-tailed distributions are produced (for instance, in relaxation spectra of glasses) is by a superposition of Gaussian distributions, each having slightly different means and variances.

Both these ways are observed in our chemical model of gene expression. We have observed two regimes: one having very small protein numbers and lognormal protein distributions that are primarily shaped by the thermal noise, and the other having larger protein numbers and protein distributions that are primarily shaped by the partitioning noise in cell division. In this regime our analysis indicates that the protein distribution is a superposition of Gaussians. Such a superposition can reproduce both the long-tailed distributions seen at early times in the growth curve of the bacterial population as well as the crossover to more symmetric distributions at later times.

Thus, the results of our detailed simulations of gene expression indicate that long-tailed protein distributions are a robust outcome of the constituent chemical processes, in combination with the different sources of noise that affect a cell. An interesting result of our simulations is that there exists a regime where the partitioning noise due to cell division dominates over the thermal noise. Many of the features seen in the simulation are validated by the single-cell experiment presented here.

Because long-tailed distributions are likely to occur, it is probable that cells would evolve to exploit this feature where beneficial and suppress it where harmful. We suggest that one way in which long-tailed distributions could be exploited is by combining them with cooperative switches. A simple, analytically tractable way of combining switches with protein distributions of different shapes reveals that long-tailed distributions are more sensitive than comparable Gaussian distributions. Thus, where sensitivity is useful, the cell might evolve to have a long-tailed distribution, and where a less

sensitive response is required, the cell might evolve to have a Gaussian protein distribution.

A more realistic way in which such switches could act on proteins is via a feedback mechanism that increases or decreases the level of expression of the corresponding gene. We have shown that when long-tailed distributions are combined with switches by using a positive feedback loop, this can make the system bistable and result in bimodal protein distributions. As discussed earlier in this article, such bimodal distributions are seen in engineered gene circuits which couple a switch to a GFP gene (16) as well as in an autoregulatory system (3). In ref. 23, there is a discussion of another engineered autoregulatory system with positive feedback that exhibits bimodal protein distributions. This mechanism also provides a way to increase the phenotypic diversity of a population of cells. This is of relevance in early embryogenesis, where producing phenotypic diversity is crucial for subsequent developmental processes, and could explain the observed bimodal distribution of hunchback in early-stage *Drosophila* embryos. In this context, morphogens could perform the role of the cooperative switches; a morphogen gradient would be a convenient way of exposing cells at different spatial locations to switches with different thresholds or amplitudes.

There are parameter regimes where a switch incorporated via positive feedback does not produce bimodality. In these regimes, the effect is rather to shift the peak of the distribution, thus selecting a subpopulation of cells from the original population. This might be relevant for later stages of embryogenesis, where it becomes necessary for the system to move from producing more phenotypic diversity to selecting specific cells that will trigger further developmental processes. Examples of such selection are well known from studies of neurogenesis, where the Notch signal plays the role of the selective switch. These observations suggest that cells might tune protein distributions to either be sensitive or robust to switches depending on the context and requirements of the cell.

Genome-scale microarray experiments have revealed that the overall distribution of gene expression, taking all of the genes into account, is also long-tailed (34). This overall distribution is a superposition of a number of distributions, with widely varying means, variances, and skews. Switches could act at this larger scale also to select out the portions of the tail of the overall distribution. In contrast to acting on individual genes, a switch acting on this scale would be selecting *modules* of interconnected genes that are present in the tail of the overall distribution. Thus, we can conjecture that the sensitivity of long-tailed distributions to switches is also exploited at genome-wide scales.

We thank G. Ananthakrishna, A. Sarin, and K. VijayaRaghavan for helpful discussions.

- Ozbudak, E. M., Thattai, M., Kurtser, I., Grossman, A. D. & van Oudenaarden, A. (2002) *Nat. Genet.* **31**, 69–73.
- Elowitz, M., Levine, A. J., Siggia, E. D. & Swain, P. S. (2002) *Science* **297**, 1183–1186.
- Banerjee, B., Balasubramanian, S., Ananthakrishna, G., Ramakrishnan, T. V. & Shivashankar, G. V. (2004) *Biophys. J.* **86**, 3052–3059.
- Mihalcescu, I., Hsing, W. & Leibler, S. (2004) *Nature* **430**, 81–85.
- Blake, W. J., Kaern, M., Cantor, C. R. & Collins, J. J. (2003) *Nature* **422**, 633–637.
- Raser, J. M. & O’Shea, E. K. (2004) *Science* **304**, 1811–1814.
- Vilar, J. M., Kueh, H. Y., Barkai, N. & Leibler, S. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 5988–5992.
- Becskei, A. & Serrano, L. (2000) *Nature* **405**, 590–593.
- Korobkova, E., Emonet, T., Vilar, J. M., Shimizu, T. S. & Cluzel, P. (2004) *Nature* **428**, 574–578.
- Sasai, M. & Wolyne, P. G. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 2374–2379.
- Thattai, M. & van Oudenaarden, A. (2001) *Proc. Natl. Acad. Sci. USA* **98**, 8614–8619.
- Kepler, T. B. & Elston, T. C. (2001) *Biophys. J.* **81**, 3116–3136.
- Swain, P., Elowitz, M. & Siggia, E. D. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 12795–12800.
- Sato, K., Ito, Y., Yomo, T. & Kaneko, K. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 14086–14090.
- Poullson, J. (2004) *Nature* **427**, 415–418.
- Kobayashi, H., Kaern, M., Araki, M., Chung, K., Gardner, T. S., Cantor, C. R. & Collins, J. J. (2004) *Proc. Natl. Acad. Sci. USA* **101**, 8414–8419.
- Kierzek, A. M., Zaim, J. & Zielenkiewicz, P. (2001) *J. Biol. Chem.* **276**, 8165–8172.
- Gillespie, D. T. (1977) *J. Phys. Chem.* **81**, 2340–2361.
- Puchalka, J. & Kierzek, A. M. (2004) *Biophys. J.* **86**, 1357–1372.
- Shapiro, B. (1990) *Phys. Rev. Lett.* **65**, 1510–1513.
- Rossi, F. M. V., Kringstein, A. M., Spicher, A., Guicherit, O. M. & Blau, H. M. (2000) *Mol. Cell* **6**, 723–728.
- Hasty, J., Pradines, J., Dolnik, M. & Collins, J. J. (2000) *Proc. Natl. Acad. Sci. USA* **97**, 2075–2080.
- Isaacs, F. J., Hasty, J., Cantor, C. R. & Collins, J. J. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 7714–7719.
- Becskei, A., Seraphin, B. & Serrano, L. (2001) *EMBO J.* **20**, 2528–2535.
- Ferrell, J. E. & Machleder, E. M. (1998) *Science* **280**, 895–898.
- Ferrell, J. E. (2002) *Curr. Opin. Cell Biol.* **14**, 140–148.
- Schleif, R. (1992) *Annu. Rev. Biochem.* **61**, 199–223.
- Vilar, J. M. & Leibler, S. (2003) *J. Mol. Biol.* **331**, 981–989.
- Rippe, K. (2001) *Trends Biochem. Sci.* **26**, 733–740.
- Driever, W. & Nusslein-Volhard, C. (1989) *Nature* **337**, 138–143.
- Houchmandzadeh, B., Wieschaus, E. & Leibler, S. (2002) *Nature* **415**, 798–802.
- Kosman, D., Reinitz, J. & Sharp, D. H. (1999) in *Proceedings of the 1998 Pacific Symposium on Biocomputing*, eds Altman, R., Dunker, K., Hunter, L. & Klein, T. Available at www.smi.stanford.edu/projects/helix/psb98.
- Kosman, D., Small, S. & Reinitz, J. (1998) *Dev. Genes Evol.* **208**, 290–294.
- Naef, F., Hacker, C. R., Patil, N. & Magnasco, M. (2002) *Genome Biol.* **3**, 0018.1–0018.11.