

Hierarchy and feedback in the evolution of the *Escherichia coli* transcription network

M. Cosentino Lagomarsino^{*†‡}, P. Jona[§], B. Bassetti^{†¶}, and H. Isambert^{*‡}

^{*}Unité Mixte de Recherche 168/Institut Curie, 26 rue d'Ulm, 75005 Paris, France; [†]Università degli Studi di Milano, Dipartimento di Fisica, Via Celoria 16, 20133 Milano, Italy; [§]Politecnico di Milano, Dipartimento di Fisica, Pza Leonardo Da Vinci 32, 20133 Milano, Italy; and [¶]Istituto Nazionale di Fisica Nucleare, 20133 Milano, Italy

Edited by Michael S. Waterman, University of Southern California, Los Angeles, CA, and approved January 22, 2007 (received for review October 12, 2006)

The *Escherichia coli* transcription network has an essentially feedforward structure, with abundant feedback at the level of self-regulations. Here, we investigate how these properties emerged during evolution. An assessment of the role of gene duplication based on protein domain architecture shows that (i) transcriptional autoregulators have mostly arisen through duplication, whereas (ii) the expected feedback loops stemming from their initial cross-regulation are strongly selected against. This requires a divergent coevolution of the transcription factor DNA-binding sites and their respective DNA *cis*-regulatory regions. Moreover, we find that the network tends to grow by expansion of the existing hierarchical layers of computation, rather than by addition of new layers. We also argue that rewiring of regulatory links due to mutation/selection of novel transcription factor/DNA binding interactions appears not to significantly affect the network global hierarchy, and that horizontally transferred genes are mainly added at the bottom, as new target nodes. These findings highlight the important evolutionary roles of both duplication and selective deletion of cross-talks between autoregulators in the emergence of the hierarchical transcription network of *E. coli*.

The successful adaptation of microorganisms to an environment or host is determined by the correct response to external and internal stimuli through the simultaneous expression of a large set of genes. The basal mechanism that performs this task is transcriptional regulation, so that it becomes important to characterize this regulatory process from a global, or “network,” viewpoint. Transcriptional regulation networks are defined starting from the basic functional elements of transcription (1). To construct the associated graph, one usually represents each operon with a node, and each regulatory interaction with a directed link $A \rightarrow B$ between the target operon B and the operon A coding for a transcription factor (TF) that has at least one binding site in the *cis*-regulatory region of B. A transcription factor regulating its own expression is called an autoregulator (AR). With this definition, the interaction graph structure is accessible by large-scale and collections of small-scale experiments (2–5).

Some topological and evolutionary properties of transcription networks have been elucidated (6–8). In particular, they can be analyzed in terms of a hierarchy of inputs that produce output responses (9–11). Specifically, the *Escherichia coli* transcription network has an essentially feedforward layered structure, where feedback is mainly limited to autoregulations (9, 10). The abundance of the latter is, however, striking, as they concern more than half of the transcription factors (12). Here, after quantifying the marginality of these properties with respect to a null network ensemble, we investigate how they could have emerged during evolution. An assessment of the role of gene duplication based on protein domain architecture shows that (i) transcriptional autoregulators have mostly arisen through duplication, whereas (ii) the expected feedback loops stemming from their initial cross-regulation are strongly selected against. This requires a divergent coevolution of the autoregulator DNA binding sites and their respective DNA *cis*-regulatory regions. Moreover, we find that the network shows a tendency to grow by

expansion of the existing hierarchical layers of computation, rather than by addition of new layers. We also argue that *de novo* rewiring of regulatory links due to mutation/selection of novel transcription factor/DNA binding interactions does not affect the hierarchy, and that horizontally transferred genes are mainly added at the bottom, as new target nodes. Our findings are consistent with a view of prokaryote evolution based on ancient duplications and conservation of stable central parts despite widespread horizontal gene transfers (13, 14).

Feedback and Hierarchy

A priori, one may expect that transcription networks contain abundant feedback loops involving two or more genes (15, 16). However, for the case of *E. coli*, the available data indicate that this is not the case (9–11). The Shen-Orr data set (2) (423 operons; 117 TFs, 578 interactions) does not contain any non-self-regulatory feedback loop for the *E. coli* transcription network. Such a tree-like directed graph is naturally organized in feedforward layers of computation, ending with target genes (TGs) as “leaves.” The layers and their numbering can be defined by the longest chain of (different) regulators upstream of each TF or TG in each layer (Fig. 1 *a* and *d*). Members of layer one are regulated by at most themselves, members of layer two are regulated by a chain of one transcription factor and possibly themselves, and so on. There are five hierarchical layers in the Shen-Orr data set (2), which is considerably lower than for randomized null networks (see Fig. 1*c*). Approximately 50% of the nodes (TFs and TGs) lay in layer two, with 69% of all TF nodes located in layer one. The notable exception to this general lack of feedback is the substantial presence of feedback loops involving a single node, or ARs (59 ARs of 117 TFs) (12, 17–19). The more recent publicly available database RegulonDB 5.5 (3) includes larger data sets (3, 9, 10) (648 operons; 147 TFs, 1,170 interactions, 85 ARs, excluding σ -factor interactions). By contrast with the Shen-Orr data set, it contains a few (4) non-self-regulatory feedback loops and a few more (a total of seven) hierarchical layers but still considerably less than in randomized null networks [see Note 5 in supporting information (SI) Appendix]. Hence, the same trend is seen for both Shen-Orr and RegulonDB 5.5 (3) data sets.

To quantify the significance of regulatory feedback and hierarchical properties of the *E. coli* transcription network, we compared it for each data set (Shen-Orr and RegulonDB 5.5) with randomized null networks with the same degree sequence, i.e., conserving the number of incoming and outgoing links for

Author contributions: M.C.L., B.B., and H.I. designed research; M.C.L. and P.J. performed research; and M.C.L. and H.I. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Abbreviations: TF, transcription factor; AR, autoregulator; TG, target gene.

[†]To whom correspondence may be addressed. E-mail: mcl@curie.fr or herve.isambert@curie.fr.

This article contains supporting information online at www.pnas.org/cgi/content/full/0609023104/DC1.

© 2007 by The National Academy of Sciences of the USA

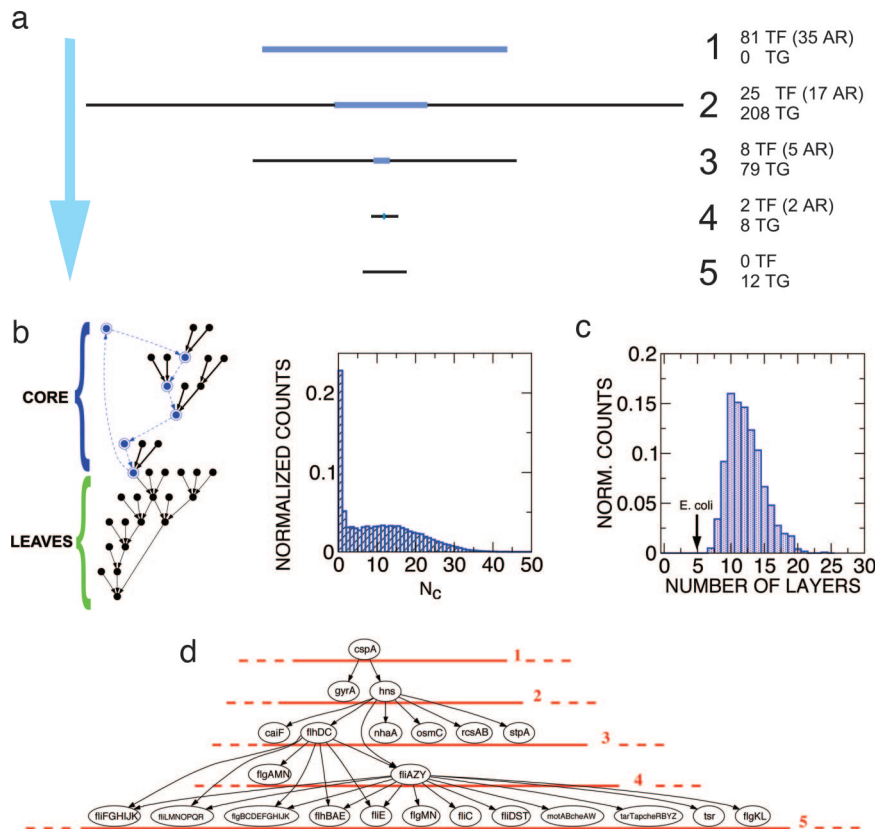


Fig. 1. Feedback and hierarchy in the *E. coli* transcription network. (a) Scheme of the layer structure of the network. Direction of regulatory links is from top to bottom. Each line represents a layer, populated by TFs (blue, thick line) and TGs (black, thin line). Members of layer i are regulated at most by $i - 1$ nodes plus themselves. By definition, layer one is constituted entirely by TFs. Annotations on the right side of the layers specify their population of TGs, TFs, and ARs. (b) Evaluation of feedback with the leaf-removal algorithm. (Right) Illustration of the leaf-removal algorithm. Leaves are nodes that do not regulate any other node. Removal of one leaf and its regulatory links may create a new leaf. Iterative removal of leaves has to stop at a core of nodes that contains loops (blue, circled nodes, dashed links). The core might contain tree-like components upstream of the loops (black). (Left) Histogram of the number of nodes in the core N_c for randomized counterparts of *E. coli* (16). The data refer to 1.1×10^6 accepted MCMC moves for randomization (see *Methods* and Note 1 in *SI Appendix*). (c) Histogram of the layer number in the randomized counterparts of the *E. coli* network. The average number of observed layers is ≈ 12 , to compare to the 5 of *E. coli*. The data correspond to a MCMC run where a total of 5.78×10^8 matrices were generated (of which $\approx 1.23 \times 10^8$ were tree-like). (d) The flagella-building subnetwork is the only example of functional subnetwork that spans all of the five layers. Here, this subnetwork is constructed arbitrarily starting from a member of layer one and following the tree downstream.

each node (Fig. 1 and Note 1 in *SI Appendix*). For both data sets, the number of ARs found in the empirical network greatly exceeds the same quantity for randomized counterparts, confirming previous observations on self-regulatory feedback (2, 12, 19). The importance of non-self-regulatory feedback was quantified by the size of the regulatory core obtained after pruning the tree-like input and output cascades using the leaf-removal algorithm (see Fig. 1b and Notes 1 and 5 in *SI Appendix*). From this analysis, we conclude that the importance of transcriptional, non-self-regulatory feedback is significantly lower in both empirical networks (Shen-Orr and RegulonDB 5.5) than in their randomized network counterparts (see Fig. 1b and SI Fig. 13 in *SI Appendix*).

The importance of hierarchy was also quantified. As there is no straightforward definition of hierarchy in general for networks including feedback, we have used the total number of layers in the tree-like input and output branches of the network as practical definition of hierarchy. This also corresponds to the number of iterations of the leaf-removal algorithm (however, see alternative definitions of hierarchy in Note 5 in *SI Appendix*). Note, in particular, that it correctly recovers the actual number of hierarchical layers for tree-like directed graphs (overlooking possible self-regulatory links as in the case of Shen-Orr data set). Comparisons with null models were restricted to randomized

networks with the same regulatory core size. Remarkably, the number of hierarchical layers was found to be considerably lower than in typical randomized network counterparts for both Shen-Orr and RegulonDB 5.5 data sets, see Fig. 1c and SI Fig. 14 and Note 5 in *SI Appendix*.

Evolutionary Drives

What is the evolutionary origin of this peculiar structure? There are three main mechanisms for the evolution of a transcription network: (i) gene duplication, (ii) rewiring of links by mutation/selection of TF/DNA interactions, and (iii) horizontal gene transfer. All three mechanisms, which we discuss below in the context of transcription network evolution, have been shown to play a substantial role in prokaryote evolution (1, 8, 14, 20–22). For clarity, the following discussion refers only to the Shen-Orr data set, which is still, to date, the most widely used data set. The same detailed analysis on the RegulonDB 5.5 data set is discussed in Note 5 in *SI Appendix*.

Duplication. Following previous analyses (8, 23), we define proteins that are likely to share a common ancestor through structural domain assignments of the SUPERFAMILY database (24). These domains allow for the definition of larger classes than sequence comparison alone (8). The database enables to asso-

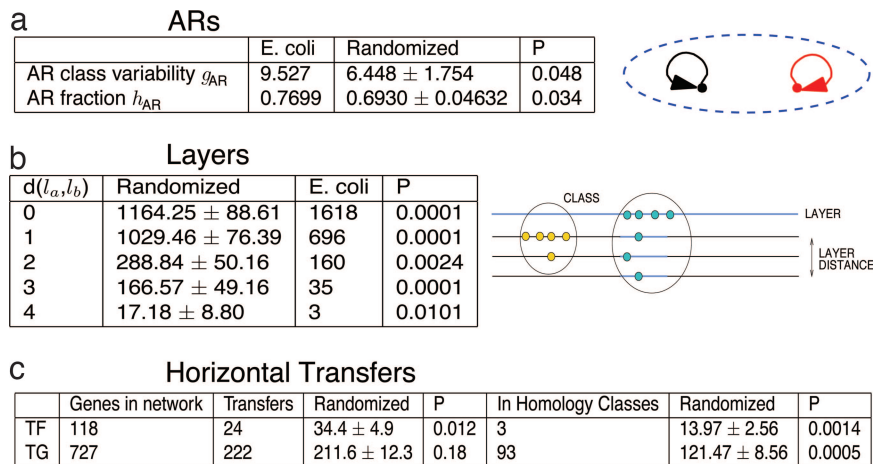


Fig. 2. Evaluation of different evolutionary drives (see also SI Table 1 in *SI Appendix*). (a) Duplicates of ARs tend to retain their self-links. This is quantified globally by the observables h_{AR} , the average fraction of ARs in classes with two or more ARs, and g_{AR} , measuring the spread in the AR population among classes that can be observed in Fig. 3a and SI Fig. 7 in *SI Appendix*. (b) Duplication and divergence preserve the layer structure. The first column indicates distance between layers (defined as the absolute difference in layer numbers), whereas the second and the third correspond to the population of duplicate genes (genes in the same homology class) at that distance, in 10^5 instances with randomized domain associations (average values) and the *E. coli* domain association data set respectively. For example, the first row (pairs of genes at distance zero) concerns the number of duplicate genes which occupy the same layer (see Fig. 3b and Note 2 in *SI Appendix*). The sketch in the right panel illustrates the distribution of nodes belonging to the same class of TFs (cyan) or TGs (yellow) among the layers, and the definition of distance between layers. (c) Fate of gene gains from horizontal transfer. TFs are underrepresented both in the class of gene gains (columns 2 and 3) and in the class of gene gains that have at least a paralog in the homology classes constructed with domain associations (columns 5 and 6).

ciate an ordered sequence of domains, or “domain architecture” to each protein. We define protein homologs as proteins whose domain architectures are identical neglecting domain repeats.^{||} We have analyzed the distribution of regulatory links between and within classes of likely duplicate genes. The statistical significance of the analysis in terms of homology classes is established (8) by comparison with random shufflings of genes (TFs and TGs separately) between classes.

The first result, summarized in Fig. 2a (see also SI Table 1 in *SI Appendix*), shows that duplicates of ARs tend to retain their self-links. We quantified this by using two global parameters, h_{AR} and g_{AR} . h_{AR} is the average fraction of ARs in classes with two or more ARs. It measures the tendency to have many ARs in one class if two are already present (the reason of the cutoff is to exclude from the count classes with two members and only one AR). g_{AR} is the variance across classes of the fraction of ARs within a class. This parameter measures the tendency to have classes that are more populated than average, and at the same time classes that are less populated than average, which can be observed in Fig. 3a (and SI Fig. 7 in *SI Appendix*). Despite this strong evidence for the proliferation of ARs through duplication events, we already mentioned the absence of any two-node feedback loops between homologous (or nonhomologous) ARs.^{**} This requires that the initial cross-regulation between duplicated ARs (reflecting the fact that binding sites are initially identical) is systematically suppressed even if self-regulation is conserved for both TF copies (Fig. 3a). We also find that single regulatory links between any kind of TFs in the same homology class are very scarce and always involve at least one AR (see Fig. S2.7). On average, 91% of the links within a homology class of TFs are self-links.

A simple duplication–divergence model (Fig. 3a and Note 3 in

SI Appendix) shows that the concomitant conservation of self-links and cancellation of cross-talks between duplicated ARs require a selective pressure for evolutionary decoupling.

This can be achieved through divergent coevolution (8, 25) of duplicate TF/DNA binding interactions. For instance, a straightforward analysis of the binding sites of CRP and FNR, two duplicate ARs regulating many TGs having no cross-regulation, shows that their own DNA *cis*-regulatory regions have higher specificities than the *cis*-regulatory regions of most of their TGs (see Note 4 in *SI Appendix*), which suggests decoupling of their self-regulatory links.

Layer Hierarchy and Rewiring. As shown in ref. 8, a large fraction of the non-self-regulatory links of the *E. coli* transcription network likely originated from duplication events. Indeed, many pairs of TGs from the same homology class are regulated by a common TF; likewise, many homologous TFs regulate the same TGs, and many pairs of TFs from the same homology class regulate homologous pairs of TGs.

Clearly, the likely duplication events underlying this transcription network expansion conserve the number of TFs upstream of each target, hence leaving the layer hierarchy untouched. The only duplication event that can actually add a layer is the duplication of an AR, provided that a crosstalk is conserved. A comparison of the homology classes with the populations of the network layers (Figs. 2b and 3b, and SI Table 1 and Note 2 in *SI Appendix*), shows that globally genes of the same homology class tend to populate the same layer.

In fact, we find only five non-self-regulatory links within homology classes (see SI Fig. 9 in *SI Appendix*) and they all involve at least one AR, suggesting that they originated from duplication events of an AR. For example, the histone-like autoregulator H-NS, belonging to layer 2, regulates its homolog StpA, which belongs to layer 3 (SI Fig. 9 in *SI Appendix*). Yet, the coincidence between the number of non-self-regulatory links within homology classes and the number of hierarchical layers in *E. coli*, does not allow us to conclude that the layers were generated by AR duplication events. Evidence for some presumed rewiring of regulatory links also exists. For instance, the same AR H-NS (SI Fig. 9 in *SI Appendix*) is also regulated by the

^{||}This corresponds to a conservative view of homology where no domains are acquired or lost after duplication. More flexible and realistic definitions of homologs yield essentially the same results (Note 2 in *SI Appendix*).

^{**}This is not strictly true for the more recent RegulonDB 5.5 data set, where a few of these two-node feedback loops are observable, though the signature for negative selection remains (see Note 5 in *SI Appendix*).

al. (14). Finally, the binding sites for the clustering analysis FNR and CRP (see Note 4 in *SI Appendix*) were taken from the RegulonDB (3) data set.

Network Analysis. We used Fortran 77 implementations of different variants (see Note 1 in *SI Appendix*) of the leaf-removal algorithms on the Shen-Orr data set (including ARs) and its randomized counterparts, which were obtained by using a standard Markov Chain Monte Carlo (MCMC) algorithm that preserves the degree sequence (marginals of the adjacency matrix) (27). This algorithm is best formulated for the adjacency matrix of the graph, i.e., the matrix A such as $(A)_{ij} = 1$ if $i \rightarrow j$, and 0 otherwise. We considered unstructured counterparts of A . Randomizations with no self-links or structurally zero diagonal of A lead to different results. For all of the tree-like instances, the number of layers corresponds to the (whole-graph) iterations that are necessary for the leaf-removal algorithm to remove the entire graph. To consider a significant sample, the number of MCMC iterations was calibrated according to the number of accepted MCMC moves (27). Specifically, we stopped the algorithm after $T = K\tau$ accepted moves, where τ is the number of nonzero elements of A , and $K = 2,000$.

Evaluation of Duplications. We constructed classes of homologous genes by using similarity criteria of the SUPERFAMILY domain architecture. Our results refer to the case where two genes are considered homologs if they share the same domains in the same order, neglecting domain repeats. A gap is considered a domain. Different choices lead to very similar results (see Note 2 in *SI Appendix*). For this analysis, proteins coded by the same operon were considered as separate entities. Many classes generated this way, such as $\{CRP, FNR\}$, are supported by evidence based on protein sequence comparison. The classes of proteins obtained this way were compared with TF–TG links in the transcription network data set. Observations related to these classes were compared with randomizations that shuffle domain associations

to gene names, separately for TFs and TGs (8). The data correspond to 10^5 randomizations.

Graph Growth Model. A simple model of duplication-divergence was considered, where at each time step duplication of the graph is followed by cancellation of links with prescribed probabilities (Note 3 in *SI Appendix*). We analyzed the evolution equations for the fraction of ARs and of intraclass links, in the different scenarios of symmetric and asymmetric divergence, presence or absence of selective conservation of ARs, and presence or absence of constant inflow of ARs. The results were compared with the observed trends in the data.

Analysis of Horizontal Gene Transfers. We used lists of imported genes obtained by a phylogenetic tree reconstruction based on 51 bacterial species (14). We presented results obtained with a gain/loss penalty of two and the hypothesis of retarded transfer, or “DELTRANS” assumption. Different choices lead to similar results (data not shown). To evaluate the partition of transferred genes between TFs and TGs, we compared with a simple binomial model where the probability of import is given by the total fraction of imported genes. As a null model for the number of imported genes that appear in homology classes, we considered classes generated by shuffling associations of genes with domain architectures as above.

Specificity of TF Binding Sites. Binding sites of two duplicate TFs were scored against their logos (28), obtained with the list of all available binding sites from RegulonDB. The specificity was defined as the difference between the scores of the same binding sites on two different logos. To improve the sensitivity, logos were computed keeping into account reverse-complement sequences and the entropy of mixing of the sets of binding sites of the two TFs under examination (see Note 4 in *SI Appendix*).

We thank M. Lercher for generously providing and illustrating data from ref. 14; H. Salgado for help with the regulonDB data set; and U. Alon, F. Poelwijk, P. ten Wolde, J. Widom, M. Vergassola, and F. Kepes for stimulating discussions.

1. Babu MM, Luscombe NM, Aravind L, Gerstein M, Teichmann SA (2004) *Curr Opin Struct Biol* 14:283–291.
2. Shen-Orr SS, Milo R, Mangan S, Alon U (2002) *Nat Genet* 31:64–68.
3. Salgado H, Santos-Zavaleta A, Gama-Castro S, Peralta-Gil M, Penaloza-Spinola MI, Martinez-Antonio A, Karp PD, Collado-Vides J (2006) *BMC Bioinformatics* 7:5.
4. Lee TI, Rinaldi NJ, Robert F, Odom DT, Bar-Joseph Z, Gerber GK, Hannett NM, Harbison CT, Thompson CM, Simon I, et al. (2002) *Science* 298:799–804.
5. Harbison CT, Gordon DB, Lee TI, Rinaldi NJ, Macisaac KD, Danford TW, Hannett NM, Tagne JB, Reynolds DB, Yoo J, et al. (2004) *Nature* 431:99–104.
6. Milo R, Itzkovitz S, Kashtan N, Levitt R, Shen-Orr S, Ayzenshtat I, Sheffer M, Alon U (2004) *Science* 303:1538–1542.
7. Warren PB, tenWolde PR (2004) *J Mol Biol* 342:1379–1390.
8. Teichmann SA, Babu MM (2004) *Nat Genet* 36:492–496.
9. Ma HW, Buer J, Zeng AP (2004) *BMC Bioinformatics* 5:199.
10. Ma HW, Kumar B, Ditzges U, Gunzer F, Buer J, Zeng AP (2004) *Nucleic Acids Res* 32:6643–6649.
11. Yu H, Gerstein M (2006) *Proc Natl Acad Sci USA* 103:14724–14731.
12. Thieffry D, Huerta AM, Perez-Rueda E, Collado-Vides J (1998) *BioEssays* 20:433–440.
13. Charlebois RL, Doolittle WF (2004) *Genome Res* 14:2469–2477.
14. Pal C, Papp B, Lercher MJ (2005) *Nat Genet* 37:1372–1375.
15. Thomas R (1973) *J Theor Biol* 42:563–585.
16. Cosentino Lagomarsino M, Jona P, Bassetti B (2005) *Phys Rev Lett* 95:158701.
17. Wall ME, Hlavacek WS, Savageau MA (2004) *Nat Rev Genet* 5:34–42.
18. Becskei A, Serrano L (2000) *Nature* 405:590–593.
19. Atkinson MR, Savageau MA, Myers JT, Ninfa AJ (2003) *Cell* 113:597–607.
20. Conant GC, Wagner A (2003) *Nat Genet* 34:264–266.
21. Dekel E, Mangan S, Alon U (2005) *Phys Biol* 2:81–88.
22. Mazurie A, Bottani S, Vergassola M (2005) *Genome Biol* 6:R35.
23. Madan Babu M, Teichmann SA (2003) *Nucleic Acids Res* 31:1234–1244.
24. Gough J, Karplus K, Hughey R, Chothia C (2001) *J Mol Biol* 313:903–919.
25. Poelwijk FJ, Kiviet DJ, Tans S (2006) *PLoS Comput Biol* 2:0467.
26. Rosenfeld N, Elowitz MB, Alon U (2002) *J Mol Biol* 323:785–793.
27. Rao AR, Jana R, Bandyopadhyay S (1996) *Indian J Stat* 58:225–242.
28. Schneider TD (2002) *Appl Bioinformatics* 1:111–119.