

Imprint of evolution on protein structures

Guido Tiana*, Boris E. Shakhnovich†, Nikolay V. Dokholyan‡, and Eugene I. Shakhnovich§¶

*Department of Physics and Istituto Nazionale di Fisica Nucleare, University of Milano, Via Celoria 16, 20133 Milan, Italy; †Bioinformatics Program, Boston University, 44 Cummington Street, Boston, MA 02215; ‡Department of Biochemistry and Biophysics, School of Medicine, University of North Carolina, Chapel Hill, NC 27599; and §Department of Chemistry and Chemical Biology, Harvard University, 12 Oxford Street, Cambridge, MA 02138

Edited by Peter G. Wolynes, University of California at San Diego, La Jolla, CA, and approved December 22, 2003 (received for review October 16, 2003)

We attempt to understand the evolutionary origin of protein folds by simulating their divergent evolution with a three-dimensional lattice model. Starting from an initial seed lattice structure, evolution of model proteins progresses by sequence duplication and subsequent point mutations. A new gene's ability to fold into a stable and unique structure is tested each time through direct kinetic folding simulations. Where possible, the algorithm accepts the new sequence and structure and thus a "new protein structure" is born. During the course of each run, this model evolutionary algorithm provides several thousand new proteins with diverse structures. Analysis of evolved structures shows that later evolved structures are more designable than seed structures as judged by recently developed structural determinant of protein designability, as well as direct estimate of designability for selected structures by thermodynamic sampling of their sequence space. We test the significance of this trend predicted on lattice models on real proteins and show that protein domains that are found in eukaryotic organisms only feature statistically significant higher designability than their prokaryotic counterparts. These results present a fundamental view on protein evolution highlighting the relative roles of structural selection and evolutionary dynamics on genesis of modern proteins.

The wealth of data emerging from fully sequenced genomes and structural proteomics provide major insight into reconstruction of evolutionary history of protein domains (1). In particular, it was found that distributions of many properties observed in protein universe can be well fit by power law (2–4). We showed in our recent work that the observed power-law distribution stemming from domain structure comparison can be explained by evolutionary dynamics that models all proteins as diverging from one or few precursors (4). Our model succeeded in the quantitative description of power-law distribution in the degree similarity of protein domains (4). We make use of this recent success as a starting point for thinking about more concrete models describing the origins of modern protein domains. Many current models describing divergent evolution are formulated in abstract terms of protein domains or sequences as nodes of dynamically evolving graphs; as such, they tend to assign the observed inequalities in fold and sequence family size to pure evolutionary chance. It is therefore hard to evaluate how realistic these models are because they do not take into account the physical constraints imposed by the thermodynamics of the sequence–structure relationship in real proteins. Other researchers motivated mostly by arguments from protein physics proposed that structures of existing proteins are highly nonrandom. It was suggested (5–9) that one of the possible factors determining evolutionary success of a structure in evolutionary selection is its *designability*, i.e., its ability to accommodate numerous sequences that can fold stably into that structure.

Some have argued that the designability hypothesis implicitly assumes convergence as a major mechanism by suggesting that various sequences may converge to the same highly designable structures irrespective of their evolutionary history. Although the two views (evolutionary dynamics by divergence and the designability hypothesis) make valid points, both lack the necessary detail needed for evaluation of their relative correspondence with real domains. Attempts to reconcile the two views

have been made in the past (9, 10). For example, Taverna and Goldstein (9), using a two-dimensional lattice model where all sequences and conformations of short chains can be exhaustively enumerated, found that for some evolutionary scenarios, where stability conditions were imposed, the ensemble of evolved lattice proteins was indeed enriched by more designable structures.

In this work, we address the question of the relative roles of chance and selection in protein evolution by simulating a more realistic divergent model of protein structure evolution. This version is based on a three-dimensional lattice model representation of protein structure. In the past, lattice models were instrumental in gaining fundamental insights into protein folding (9, 11–16). Despite their approximate character, they feature a unique sequence–structure relationship akin to that of real proteins (17). The major benefit of such models is that they are computationally tractable so that it is feasible to run a realistic evolutionary scenario that includes testing by direct kinetic simulations the ability of emerging proteins to fold and be stable. More details of the evolutionary algorithm are provided in *Methods*.

Methods

Lattice Model. We employ standard cubic lattice model where lattice amino acids occupy lattice sites and each site can be occupied by no more than one amino acid (15–17). Sequence neighbors occupy neighboring lattice sites. Only lattice amino acids that are in spatial contact, but are not sequence neighbors, can interact. Energy of each contact interaction is determined by the types of amino acids involved. We use the model with 20 types of amino acid and Miyazawa–Jernigan group potentials from table 6 of ref. 18. The Monte-Carlo folding algorithm is as described (17, 19) with move set including end moves, corner flips, and crankshaft moves. Every attempt to move a monomer is counted as a time step.

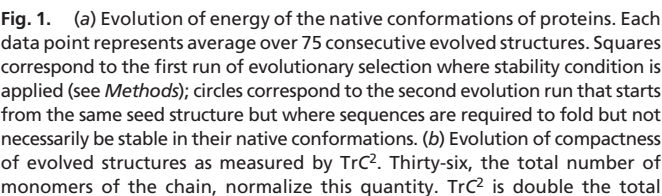
Evolutionary Algorithm. Our evolutionary model uses a cubic lattice of 36-mer as a basic model. It proceeds as follows:

1. Start from initial structure and design a sequence that stably folds into that structure with Monte-Carlo design in sequence space (17, 20). Check, with folding Monte-Carlo simulation in conformational space (17, 19), that the designed sequence does indeed stably fold into the target native structure.
2. Keeping target structure fixed, perform Monte-Carlo in sequence space (in the form of swaps as elementary move, to preserve amino acid composition). This step runs at a certain evolutionary temperature, T_{evol} , that has to be carefully selected (see below). This step creates sequence families providing divergence in sequence but not structural space.
3. Randomly select several evolved proteins and make gene duplication and point mutation attempts for each. Fold each of the new gene sequences several (10) times each, starting

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

¶To whom correspondence should be addressed. E-mail: shakhnovich@chemistry.harvard.edu.

© 2004 by The National Academy of Sciences of the USA



What is the reason for evolutionary pressure on somewhat esoteric structural parameters such as higher traces of model protein contact matrices? Comparison of two curves in Fig. 1c clearly rules out that structures with higher Tr^8 evolved in response to pressure of creating structures with higher stability. Furthermore, higher-order traces of contact matrix are correlated with contact density (Tr^2) in random ensemble (see Figs.

Tiana et al.

based phylogeny is not complete despite several fruitful efforts (29). Our approach here is based on the fact that eukaryotic cells evolved after prokaryotic ones and thus protein domains that exist only in eukaryotes (eukaryotic innovation domains) can serve as representatives of later-evolved protein structures, to be compared with domains that are exclusive to prokaryotes. Sequence analysis with the ELISA database (30) (<http://romi.bu.edu/elisa>) yielded 817 eukaryotic innovation Dali domains and 1,775 prokaryotic-only Dali domains. In accord with predictions from model evolution, we found that eukaryotic innovation domains are indeed statistically less compact than prokaryotic-only domains (Fig. 4a). Importantly, this trend is not a consequence of possible differences in length distribution between eukaryotic innovation and prokaryotic-only domains; it persists in all length windows (see Figs. 5–8).

To complete the analogy with our modeling, we turn to the analysis of higher-order contact traces of eukaryotic innovation domains and prokaryotic domains. In doing so one has to keep in mind that there is positive correlation between traces of second-order (CD) and higher-order contact traces. To this end, to make appropriate comparison between higher-order contact traces, we select only domains that fall into a narrow range of CDs, namely between 3 and 4. This range corresponds to domains of lower than average compactness (23) consistent with our expectations that selection mechanism based on higher-order contact traces is likely to work on domains of relatively low compactness. Comparison between eukaryotic innovation domains and prokaryotic-only domains (Fig. 4b) shows a statistically significant shift toward higher TrC^8 in eukaryotic innovation domains. The difference in distribution of TrC^8 between the two groups of protein domains is highly statistically significant: KS P value for the null hypothesis of no difference in distributions is 0.001. Comparison of eukaryotic innovation and prokaryotic-only domains shows that in this case possible loss of designability in eukaryotes due to their statistically lower contact density is compensated by contributions of the higher-order traces of contact matrices of their domains, similar to the effect observed in model evolution. This phenomenon is more complex than earlier observation that domains from thermophilic organisms are more designable (23). In the latter cases, greater designability of thermophilic proteins is observed at the level of contact density, i.e., in the lowest order of expansion in Eq. 1. Higher-order contact traces show similar trend in thermophiles vs. mesophiles. However, in the case of thermophilic organisms it is hard to say whether this is an independent trend or a consequence of correlation between traces of different order (see Figs. 5–8).

Our results clearly show that (i) increasing protein designability was certainly a factor in evolution of structures and (ii) that existing theory (25) correctly predicts the structural determinants of protein designability.

Why did model (and real) evolution optimize designability? More designable proteins were not selected through optimization of foldability (and stability) applied in model evolution. Nevertheless, divergent evolution came up with more designable structures. This is clearly an example of collective selection pressure where evolution is controlled not only by measures of fitness for an individual protein and its functionality inside an organism, but also by the accessibility of new structure for innovation. In this sense, designability represents a selective advantage for *ensembles* of evolving proteins.

Our simulations represent a rigorous albeit minimalist model of protein evolution where only the basic physical characteristics of proteins (their ability to fold and their stability) were chosen for selection. Whereas this requirement represents a bare-bones, necessary condition for proteins to survive, in reality, other requirements such as functional selection and ability to participate in protein–protein interactions were most likely factors in

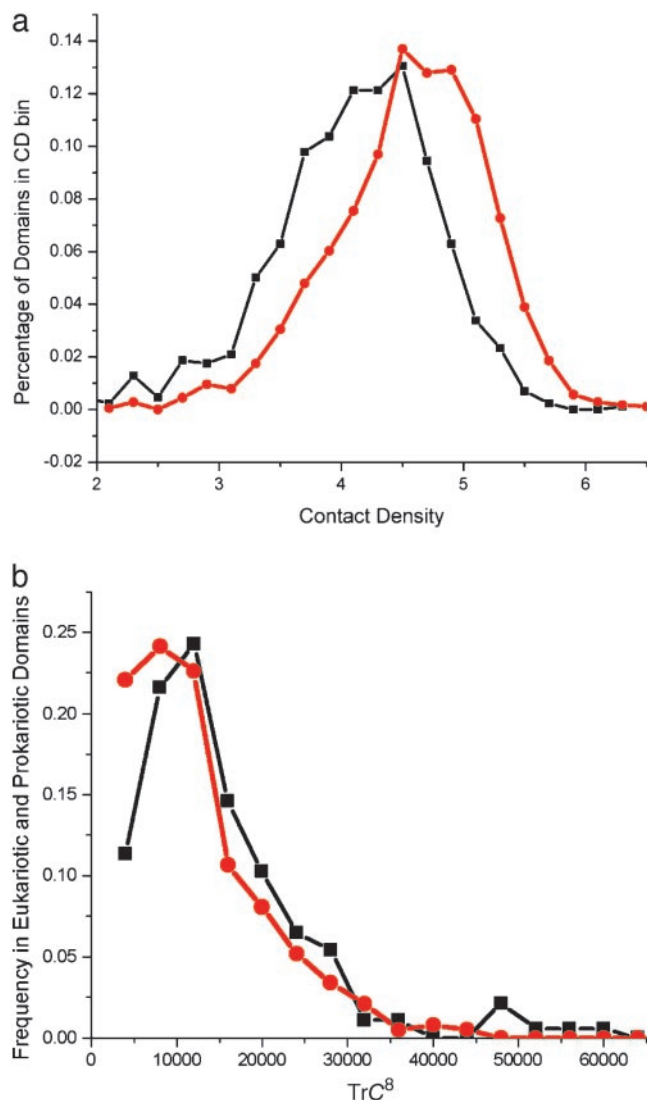


Fig. 4. (a) Distribution of CD in eukaryotic innovation domains (squares) and prokaryotic-only domains (circles). The data for these histograms are binned into bins of size 0.2. The CD is calculated as explained in *Methods*. (b) Distribution of TrC^8 normalized by domain length in prokaryotic-only Dali domains (circles) and eukaryotic innovation domains (squares). The data were binned with bin size 4,000. Only domains with CD in the range between 3 and 4 are taken; 262 eukaryotic innovation domains and 843 prokaryotic-only domains fall in this range. The Kolmogorov–Smirnov (KS) P value for the null hypothesis that these two datasets were drawn from the same distribution is 0.001. For control, we randomly split prokaryotic-only domains ensemble into two equal parts and compared their distributions of TrC^8 . In contrast to comparison between eukaryotic innovation and prokaryotic-only domains, these two sets appear to be identical: KS P value for the same null hypothesis is 0.647.

selection of specific protein structures (1, 31). Nevertheless, we see that even such minimal requirements on selection lead to significant consequences for the evolution of the protein structure universe.

Our analysis presents a rare example when a very specific prediction derived from theory appears to directly affect protein evolution, both in model and real proteins. Although earlier studies (23) suggested that this may be the case, they focused exclusively on comparison of CD of proteins from various proteomes. Although designability appeared to be the most plausible explanation for observed differences between meso-

philic and thermophilic proteomes in ref. 23, other, perhaps related, factors such as stability and/or aggregation could not be ruled out. Indeed as Fig. 1*b* shows, the less stable proteins evolve with lower compactness in general. The present study goes much further as it demonstrates that less obvious structural characteristics were evolutionarily selected. Furthermore, as model evolution suggests (Fig. 1*c*), such selection is unrelated to protein stability.

The relation between properties of contact matrices and protein topology points out to a possible reason for remarkable symmetry observed in proteins. Indeed, our analysis shows that availability of uninterrupted closed loops of intraprotein contacts may be beneficial for structure designability. One way to achieve this is to form regular structures, in particular with symmetric open interiors that are not interrupted by the crossing chain. This is one of the most common structural features of globular proteins, most of which have contiguous hydrophobic cores.

Finally, our findings highlight the interplay between selection and chance in protein evolution. While evolutionary pressure is

applied directly to select for proteins that can stably fold, it is countered by the difficulty in finding a proper structure–sequence match. The result of this balance between selection and effective entropic factors in structure space is the emergence of more designable structures, whereas in sequence space, evolution selects sequences that have pronounced, but not extreme, energy gaps in their native conformations (32, 33). This represents close analogy with statistical mechanics where temperature serves as a rough equivalent of the strength of selective pressure.

This study shows that a divergent model of protein structure morphogenesis is able not only to reproduce global power laws observed in protein universe (3, 4) but also to capture the unexpected selection of special structures that we observe to be predominant in the protein universe.

We are grateful to Eric Deeds and Richard Goldstein for useful discussions. This work was supported by the National Institutes of Health (E.I.S.). N.V.D. acknowledges support from a University of North Carolina Research Council grant.

- Ponting, C. P. & Russell, R. R. (2002) *Annu. Rev. Biophys. Biomol. Struct.* **31**, 45–71.
- Qian, J., Luscombe, N. M. & Gerstein, M. (2001) *J. Mol. Biol.* **313**, 673–681.
- Koonin, E. V., Wolf, Y. I. & Karev, G. P. (2002) *Nature* **420**, 218–223.
- Dokholyan, N. V., Shakhnovich, B. & Shakhnovich, E. I. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 14132–14136.
- Finkelstein, A. V. & Ptitsyn, O. B. (1987) *Prog. Biophys. Mol. Biol.* **50**, 171–190.
- Finkelstein, A. V., Gutin, A. M. & Badretdinov, A. (1995) *Subcell. Biochem.* **24**, 1–26.
- Govindarajan, S. & Goldstein, R. A. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3341–3345.
- Li, H., Helling, R., Tang, C. & Wingreen, N. (1996) *Science* **273**, 666–669.
- Taverna, D. M. & Goldstein, R. A. (2000) *Biopolymers* **53**, 1–8.
- Xia, Y. & Levitt, M. (2002) *Proc. Natl. Acad. Sci. USA* **99**, 10382–10387.
- Mirny, L. & Shakhnovich, E. (2001) *Annu. Rev. Biophys. Biomol. Struct.* **30**, 361–396.
- Klimov, D. K. & Thirumalai, D. (2001) *Proteins* **43**, 465–475.
- Socci, N. D., Onuchic, J. N. & Wolynes, P. G. (1998) *Proteins* **32**, 136–158.
- Pande, V. S., Grosberg, A. Y. & Tanaka, T. (1997) *Biophys. J.* **73**, 3192–3210.
- Sali, A., Shakhnovich, E. & Karplus, M. (1994) *Nature* **369**, 248–251.
- Chan, H. S. & Dill, K. A. (1996) *Proteins* **24**, 335–344.
- Shakhnovich, E. I. (1994) *Phys. Rev. Lett.* **72**, 3907–3910.
- Miyazawa, S. & Jernigan, R. L. (1985) *Macromolecules* **18**, 534–552.
- Sali, A., Shakhnovich, E. & Karplus, M. (1994) *J. Mol. Biol.* **235**, 1614–1636.
- Shakhnovich, E. I. & Gutin, A. M. (1993) *Protein Eng.* **6**, 793–800.
- Miyazawa, S. & Jernigan, R. L. (1996) *J. Mol. Biol.* **256**, 623–644.
- Shakhnovich, E. I. (1998) *Fold. Des.* **3**, R45–R58.
- England, J. L., Shakhnovich, B. E. & Shakhnovich, E. I. (2003) *Proc. Natl. Acad. Sci. USA* **100**, 8727–8731.
- Abkevich, V. I., Gutin, A. M. & Shakhnovich, E. I. (1994) *Biochemistry* **33**, 10026–10036.
- England, J. L. & Shakhnovich, E. I. (2003) *Phys. Rev. Lett.* **90**, 218101.
- Wolynes, P. G. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14249–14255.
- Landau, L. D., Lifshitz, E. M. & Pitasevskii, L. P. (1978) *Statistical Physics* (Pergamon, Oxford).
- Grzybowski, B. A., Ishchenko, A. V., Shimada, J. & Shakhnovich, E. I. (2002) *Acc. Chem. Res.* **35**, 261–269.
- Mirkin, B. G., Fenner, T. I., Galperin, M. Y. & Koonin, E. V. (2003) *BMC Evol. Biol.* **3**, 2.
- Shakhnovich, B. E., Harvey, J. M., Comeau, S., Lorenz, D., DeLisi, C. & Shakhnovich, E. (2003) *BMC Bioinformatics* **4**, 34.
- Teichmann, S. A., Chothia, C. & Gerstein, M. (1999) *Curr. Opin. Struct. Biol.* **9**, 390–399.
- Mirny, L. A., Abkevich, V. I. & Shakhnovich, E. I. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 4976–4981.
- Gutin, A. M., Abkevich, V. I. & Shakhnovich, E. I. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1282–1286.