# How the first biopolymers could have evolved

V. I. ABKEVICH, A. M. GUTIN, AND E. I. SHAKHNOVICH

Harvard University, Department of Chemistry, 12 Oxford Street, Cambridge, MA 02138

ABSTRACT    In this work, we discuss a possible origin of the first biopolymers with stable unique structures. We suggest that at the prebiotic stage of evolution, long organic polymers had to be compact to avoid hydrolysis and had to be soluble and thus must not be exceedingly hydrophobic. We present an algorithm that generates such sequences for model proteins. The evolved sequences turn out to have a stable unique structure, into which they quickly fold. This result illustrates the idea that the unique three-dimensional native structures of first biopolymers could have evolved as a side effect of nonspecific physicochemical factors acting at the prebiotic stage of evolution.

The problem of how first biopolymers could have evolved prebiotically is one of the most fundamental in modern science. The first biopolymer-type macromolecules were likely to perform simple catalytic functions, possibly enhancing self-replication of prebiotic nucleic acids (1). The physical prerequisite for a macromolecule to serve even as a simplest enzyme is that it has a unique sufficiently rigid three-dimensional structure. This makes it possible to form a catalytic center by keeping functional groups in correct mutual positions and orientations.

The important question, from the evolutionary point of view, is whether it is likely that such sequences with unique three-dimensional structure appeared in the primordial soup. If initial polymerization of organic molecules at the prebiological stage resulted in the set of random sequences, then can it be formulated as what fraction of all sequences can fold to and be stable in a unique structure? This question concerns the very essence of the concept of prebiological evolution. Indeed, if the fraction of foldable sequences among all random sequences is large, then there may be no need in prebiotic selection, and protein-like sequences could have been selected randomly out of a multitude of randomly synthesized biopolymers. Alternatively, if the set of foldable sequences constitutes a vanishing fraction of all sequences, then some selection was likely to have been acting even at the prebiotic stage. However, the main problem here is that at the prebiotic stage a property such as a unique structure (and the ability to perform catalytic functions) was not likely to have been a factor in selection (that is why this stage is prebiotic). Therefore, in this case, some other factors could have been guiding synthesis of prebiological polymers, and sequences with unique structures could have evolved as a byproduct of other physicochemical processes. In this paper we suggest a principle possibility for a scenario of this type.

Significant progress has been achieved in understanding kinetic and thermodynamic requirements for the folding of heteropolymer sequences. Interest in this fundamental problem was motivated by development of the theory of protein folding (for reviews, see refs. 2–4), and the corresponding topics were discussed mainly in the context of this theory.

The high complexity of the protein folding problem motivated many investigators to develop simplified models in which a protein was represented as a simple self-avoiding chain often positioned on a lattice (5–9). The key advantage of such models is that they allow one to simulate a complete folding process—from a random coil to the native conformation. Another advantage, especially important for the purpose of the present study, is that such models are sufficiently general to be applicable to various heteropolymers, not only to polypeptides. The following analysis will be based on the lattice model developed and applied to the study of protein folding. Therefore, we will often use terminology accepted in protein folding but the reader should bear in mind that the model is sufficiently general to be applicable to different kinds of biological macromolecules, not just polypeptides.

For a polypeptide sequence to be protein-like, it must satisfy thermodynamic and kinetic criteria of folding. The thermodynamic requirement is that a sequence possesses a unique stable three-dimensional structure in which the chain spends most of the time at physiological temperature. It was shown in ref. 8 that the thermodynamic requirement is not very restrictive provided that the temperature does not exceed the critical temperature $T_c$, which depends on amino acid composition of polypeptides. In that case, a significant fraction of random sequences will have a unique structure. However, random sequences fold slowly at a temperature where their ground-state conformation is stable (10, 11).

Folding of model proteins was studied by using the lattice Monte Carlo simulations (12, 13). These simulations showed that the necessary condition for a polypeptide to have a kinetically reachable stable conformation is that it is a pronounced global energy minimum (11, 14, 15). It was estimated (16) that only an exponentially small fraction of random sequences [$\approx \exp(-\alpha N)$, where $N$ is the number of monomers in a polypeptide chain] satisfies this requirement for chains of realistic lengths. The obvious implication of this result is that it is unlikely to find a protein-like sequence as a result of random sampling of amino acid sequences.

To overcome this difficulty *in machina*, a simple algorithm generating sequences with low energy for a given chain conformation was proposed in refs. 14 and 17. This is a stochastic optimization scheme that starts with random sequences and proceeds by making substitutions biased to minimize relative energy of a sequence in the native conformation. The bias is introduced explicitly via the rule according to which a substitution is accepted if it decreases relative energy. Substitutions increasing relative energy were also accepted if they satisfied the Metropolis criterion (14, 17).

The approach proved successful yielding sequences with a low relative energy in their native conformation. Folding simulations demonstrated that such sequences, which were designed to be more stable than random sequences, can also fold much faster into their native conformations (14, 18). The analysis made in the work (15) suggested that the opposite is also true: sequences designed to fold quickly into their native conformations are more stable as well. A simple algorithm of sequence selection in ref. 15 was based on the principles similar to biological evolution, i.e., random substitutions and selection pressure. The requirement of rapid folding into one specific

Abbreviation: MFPT, mean first passage time.

conformation modeled the selection pressure via acceptance of substitutions that result in faster folding and rejection of substitutions that slow down the folding rate. The initial sequence for the algorithm was quasirandom. As the algorithm proceeded, sequences with a low mean first passage time (MFPT) into their respective native conformations evolved. The most important finding of this study was that stability of such sequences in their native state increased too.

These results show that both key conditions for polypeptide sequences to be protein-like (i.e., thermodynamic stability and fast folding), can be achieved when one parameter, the relative energy, is minimized [we noted (15), though, that correlation between relative energy and folding rate may be somewhat limited, and other factors such as nucleus stability (19) can influence folding rate as well]. These findings imply that evolutionary pressure aimed at optimization of folding rate can optimize the stability and, the opposite is true as well: Optimization of stability can optimize the folding rate. This is due to the fact that one parameter, the relative energy, governs, to a large extent, both quantities.

Such schemes may model some aspects of protein evolution only at the biological stage where sequences, which reliably fold into their native conformations, have already evolved, and the evolutionary process presses for their improvement as more stable and fast-folding proteins. This fact manifests in both algorithms where the concept of unique native structure is introduced from the beginning, at the stage of formulation of selection rules.

This implies that such evolutionary improvement schemes must be "seeded" in the sense that first protein-like sequences that are able to fold into their unique native conformations must have evolved at the stage of prebiological evolution.

It is quite clear that at the prebiological stage of evolution, the notion of unique structure could hardly exist, and instead of "selection pressure," at which specific biological advantages were memorized and inherited, there could be only relatively nonspecific physicochemical factors that governed synthesis and hydrolysis of the first organic polymers. The question then arises as to how such nonspecific physicochemical factors could have selected the first protein-like sequences.

In this paper we give a possible example of such "prebiological" selection. While it is unlikely that it reproduces the actual process of prebiological evolution, it may illustrate an important idea of how biopolymers might have evolved as a result of nonspecific physicochemical factors.

We consider a model where physicochemical pressure consists of two factors: possibility of aggregation and hydrolysis (20). The latter could be catalyzed by other existing low and high molecular weight compounds. It is clear that a simple way to protect polypeptides from hydrolysis would be to make the chains compact. This biases equilibrium distribution over sequences toward more hydrophobic sequences, with prevailing attraction between residues. However, hydrophobic sequences tend to aggregate.

The real prebiotic selection of biopolymer precursors was likely to involve synthesis of a large number of sequences over a long period of time ($\approx 10^9$ years). The polymers that are less compact had a higher probability of getting hydrolyzed. This factor gradually shifts the sequence distribution toward the more compact and, as will be shown below, more stable sequences with unique structure.

In this work we propose a simple algorithm that biases sequence sampling toward compact and water-soluble sequences. While being unlikely to reproduce the details of the actual process of prebiological evolution, it generates in a computationally reasonable time the ensemble of sequences that meet these two criteria of compactness and solubility. Compactness can be measured directly in the simulation through the number of intramolecular contacts that the chain forms: the greater this number, the more compact the chain is.

With regard to solubility, it is computationally very costly to faithfully reproduce interactions of macromolecules with solvent while making a search in sequence space. The unrestricted condition of compactization would generate mainly hydrophobic sequences with a tendency to aggregate. To prevent this we will impose the simple condition of solubility as a requirement that the sequence search proceeds over sequences with a constant amino acid composition.

Now the main idea of the algorithm can be formulated: beginning with random sequences, we make substitutions preserving the overall composition of amino acids. If the chain becomes more compact, or gets compact faster, the substitution is accepted, otherwise it is discarded. After a number of cycles the procedure will generate sequences that will be compact and still not exceedingly hydrophobic. The ensemble of such sequences can then be studied to see whether their properties deviate from those of random sequences.

The first example that illustrates the general idea is similar in spirit to the model that was used in our previous study (15). There we sought sequences of 27-monomer chains that fold quickly into a given native fully compact conformation. In the spirit of the discussion above, we now modify the requirement in such a way that sequences that fold quickly into any maximally compact conformation are sought. In other words, in the evolutionary algorithm, we substitute the specific condition of fast folding by a nonspecific requirement of fast compactization.

For this purpose, we used the previous model (15). A protein chain is modeled with a self-avoiding chain of 27 monomers on a cubic lattice. The energy of a conformation of the chain is the sum of the energies of pairwise contacts:

$$E = \sum_{1 \leq i < j \leq N} [B_0 + B(\xi_i, \xi_j)] \Delta^{ij}, \qquad [1]$$

where $\Delta_{ij} = 1$, if monomers $i$ and $j$ are lattice neighbors, and $\Delta_{ij} = 0$, otherwise. $\xi_i$ defines the type of amino acid residue in position $i$. $B(\xi, \eta)$ is a magnitude of a contact interaction between amino acids of types $\xi$ and $\eta$. The contact energies were determined from statistical distributions of contacts in real proteins (21). However, protein statistics may provide only relative energies of interactions (22). In addition, we also introduced the average contact energy $B_0$, which is independent of monomers forming a contact. The motion of the chain is modeled by the standard cubic lattice Monte Carlo algorithm (12). Simulations were performed at the temperature $T = 0.32$ and with $B_0 = -T = -0.32$.

As an estimate of stability of the native conformation, we used the relative value of its energy, or $Z$ score, introduced by Eisenberg and coworkers (23) as we did in our previous work (15):

$$Z = \frac{E_{nat} - E_{av}}{\sigma}, \qquad [2]$$

where $E_{nat}$ is the energy of the native conformation and $E_{av}$ is an average energy of compact nonnative conformations. The average energy of a nonnative compact conformation can be estimated as $E_{av} = K e_{av}$, where $K = 28$ is the number of contacts in a maximally compact conformation and $e_{av}$ is an average energy of all possible contacts with a corresponding dispersion $\sigma$.

Our selection mechanism is close to the one described in ref. 15. In step 1, 50 folding runs are performed to get a reasonably good estimate for the MFPT in any out of $\approx 10^5$ maximally compact conformation. In step 2, an attempt of point substitution is made, and a new MFPT to any maximally compact conformation is roughly estimated with two folding runs. If it is longer than the original MFPT in any maximally compact conformation, then the substitution is rejected; if it is shorter,

then we estimate the new MFPT in any maximally compact conformation more precisely in step 3. The purpose of step 2 is to reject obviously poor substitutions that constitute the majority of all substitutions, thus saving central processing unit time. In step 3, we perform 10 folding simulations and so are getting a more precise estimate for MFPT in any maximally compact conformation. If it is shorter than the original MFPT by 20%, then the substitution is accepted; otherwise it is rejected.

We started with a random sequence, for which the MFPT to a maximally compact conformation was about $3 \times 10^6$ Monte Carlo steps. The selection algorithm accelerated compactization up to about $5 \times 10^4$ Monte Carlo steps after about 50 accepted substitutions. For each evolved sequence, we found in a long Monte Carlo run ($10^8$ steps) the conformation with the lowest energy, which can be considered native. Strong average attraction between the monomers made the native conformations for all sequences maximally compact. The native conformations for the starting random sequence and for the last sequences generated by the selection are shown on Fig. 1 *a* and *b*. Through the analysis of native conformations, we found that for the first 20 substitutions the native conformation changed significantly, but subsequently, the algorithm picked one conformation and kept it unchanged. Small fluctuations of the native structure from sequence to sequence were an exception; those conformations had more than 80% of common contacts. Moreover, the relative energy of the native conformation ($Z$) decreases with substitutions (Fig. 1c). It implies that as a result of selection aimed at fast compactization, the procedure, after some number of substitutions, finally selects one compact conformation as the native one and decreases its relative energy to a value as low as $-29$. As was shown in ref. 15, the distribution of $Z$ for random sequences can be well fit with a Gaussian. The estimate of probability to choose a random sequence with such a low $Z$ is about $10^{-9}$. Therefore, only one of a billion 27-monomer random se-

quences will be as stable in the native state as the one selected by our algorithm.

This result suggests that the simplest way to provide rapid folding into any of a hundred thousand maximally compact conformations is to choose one compact conformation and make its energy low. In other words, the requirement of a rapid compactization can lead to a result different from the compactization itself, that is, to the native conformation with very low energy. This suggests an idea of how proteins with stable native structures might have evolved. The main point here is that such stable structures could be a side effect of a different requirement, for example, the requirement of compactization.

However, while the calculations presented illustrate in general how the idea of nonspecific selection works to provide sequences with a unique structure, the condition of fast compactization seems a bit artificial. Instead, one should impose the condition that the chain gets compact and maintains such a compact state for some time. It is more natural to expect that this will prevent hydrolysis.

To this end, we changed the selection procedure in the following way. Instead of optimization of the folding time, we optimized the average compactness over half a million Monte Carlo steps starting with a random coil conformation. In other words, the requirement now is not only to reach a compact state rapidly but also to keep the compactness over the course of simulation. This is a more faithful representation of the idea that compactness should protect the chains from hydrolysis. This combines two requirements: that compactization will be fast and that it will lead to persistent compact conformations.

As was pointed out above, the simplest way to model the requirement of solubility and, therefore, prevent convergence to hydrophobic homopolymers is to keep amino acid composition unchanged in the process of substitutions. To model the requirement of solubility further, we introduced slight overall repulsion between molecules by taking $B_0 = 0.05$. Simulations were performed at temperature $T = 0.15$.

We slightly changed our sequence search mechanism. In step 1, 50 folding runs were performed to get a reasonably good estimate of average compactness during the first $5 \times 10^5$ Monte Carlo steps. In step 2, an attempt of pairwise substitution is made. This corresponds to picking two sites randomly and "swapping" amino acids between them. This is the natural way of changing sequence without changing amino acid composition. If the new average compactness (number of intrachain contacts) during the first $5 \times 10^5$ Monte Carlo steps is smaller, then the substitution is rejected. If it is larger, then we perform 10 folding simulations and get a more precise estimate of the new average compactness during the first $5 \times 10^5$ Monte Carlo steps in step 3. If it is at least 1% higher than the original compactness, then the substitution is accepted; otherwise, the substitution is rejected.

Thirty-two substitutions were accepted. For all of the generated sequences, we determined conformations (not necessarily maximally compact) with the lowest energy. This was done by means of extensive Monte Carlo simulations. These conformations were identified as the native ones. Some of them are shown on Fig. 2. The evolution of average compactness of the native conformations, with substitutions is shown on Fig. 3. It is clear that the compactness of the native conformations grows rapidly and reaches its plateau value corresponding to conformations having 25 contacts. After analyzing native conformations, we found that most of them differ one from another. Although for the first 7 substitutions, native conformations changed dramatically, for the next substitutions, the number of common contacts between native conformations for different sequences became large (>80%) so that sequences that evolved at the later stage of the procedure had structurally similar native conformations. One can see on Fig. 2 that native conformations that evolved as a result of this selection procedure seemed to be different from
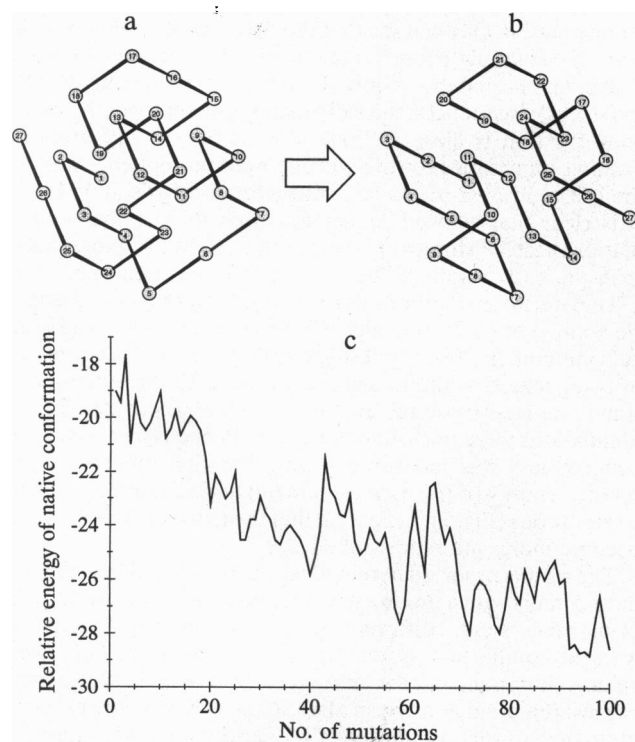


FIG. 1. Selection of sequences with 27 residues under the requirement of rapid folding into any maximally compact conformation. (*a*) Native conformations for the starting random sequence. (*b*) Native conformations for the last sequence generated by the selection. (*c*) Evolution of the relative energy of the native conformation.
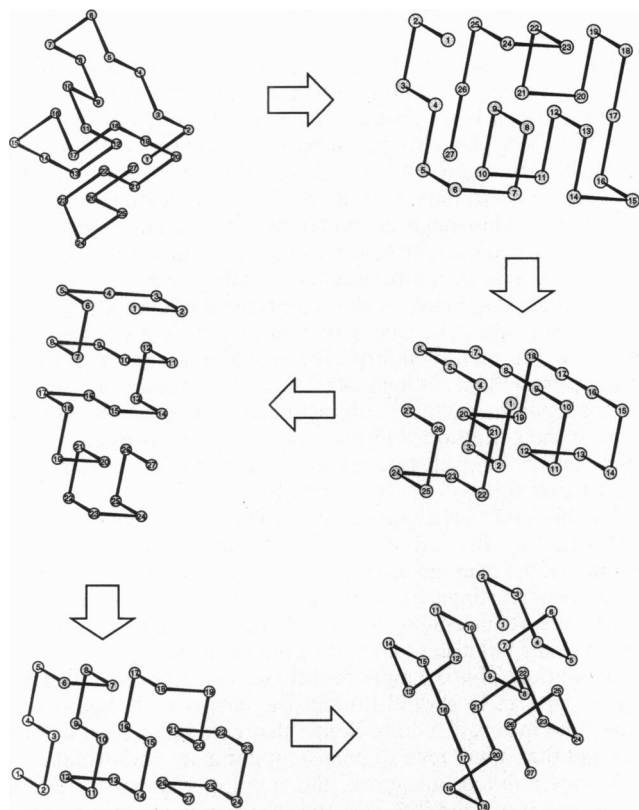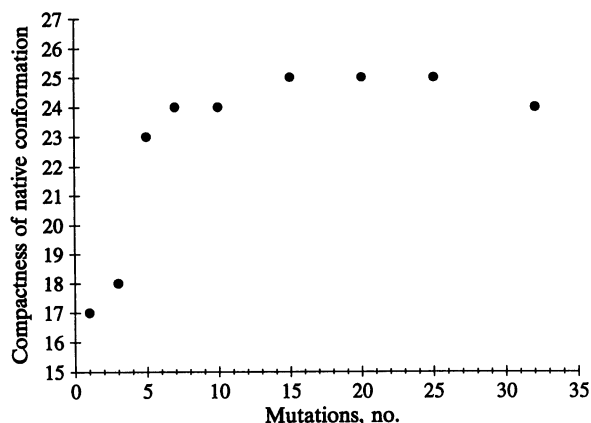
FIG. 2. Selection of sequences with 27 residues under the requirement of large average compactness over half a million Monte Carlo steps starting with a random coil conformation. Evolution of the native structure is shown. The native conformations of the original random sequence, sequence after the 5th substitution, sequence after the 10th substitution, sequence after the 20th substitution, sequence after the 25th substitution, and sequence after the 32nd (last) substitution are shown.



FIG. 4. (*a*) Evolution of the relative energy of the native conformation of 27-mers. (*b*) Evolution of the stability of the native conformation of 27-mers.

the randomly chosen compact conformations. They have a greater number of short-range contacts. However, when we repeated our selection algorithm, seeding it with another random sequence, this feature was not reproduced. Additional work has to be done to clarify what types of native conformations are more likely to evolve in this selection procedure.

As we expected, the relative energy of the native conformation ($Z$) decreased with substitutions (Fig. 4*a*), indicating

that the native state became more and more stable. This conclusion is supported by the next plot on Fig. 4*b*, which shows evolution of the thermodynamic probability of the native conformation. One can see that the native conformation of the starting randomly chosen sequence is extremely unstable: its Boltzmann probability is only about 0.1%. At the same time for the selected sequences, the Boltzmann probability of the native conformation is close to 50% after about 10 substitutions. Folding trajectories for the starting random sequence and for the last evolved sequence are shown for comparison on Fig. 5. It is clear that evolved sequences reach their native conformations faster. Moreover, their native states are more stable than the native state of the starting random sequence.

To determine whether this result holds for longer sequences, we applied our selection algorithm to longer chains consisting of 36 monomers. The result suggests that the same requirement of compactness without aggregation leads to the evolution of stable native structures for the 36-mer shown on Fig. 6. Simulations were performed at $T = 0.26$ and $B_0 = 0$. Average compactness was measured during the first $10^6$ steps. The relative energy of the native conformation ($Z$) decreased with substitutions (Fig. 7), indicating that the native states of 36-mer become more and more stable.

The results reported herein illustrate the idea of how the first biopolymers with a unique structure could have evolved: In the primordial "soup," different species with different sequences were in equilibrium. Some fractions aggregated; some were soluble but noncompact and were hydrolyzed. Our results suggest that a small fraction of all sequences was able to satisfy mutually conflicting nonspecific requirements, and these sequences could have had a unique structure.

It is commonly believed that in the earliest cells, both the genetic and enzymatic components were RNA molecules (24–26). If this hypothesis is correct, proteins evolved during biological evolution. However, the question remains as to how



FIG. 3. Selection of sequences with 27 residues under the requirement of large average compactness over half a million Monte Carlo steps starting with a random coil conformation. Evolution of compactness (the number of contacts in the native structure) is shown.
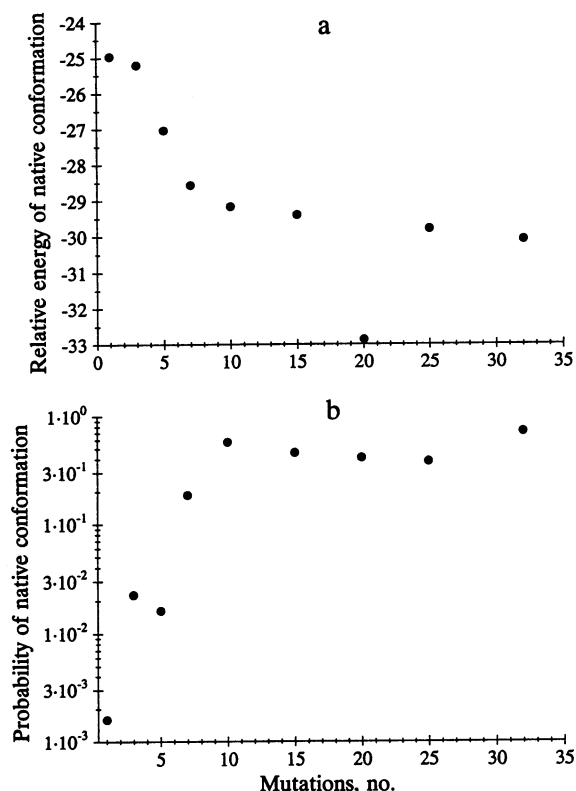
Chemistry: Abkevich *et al.*

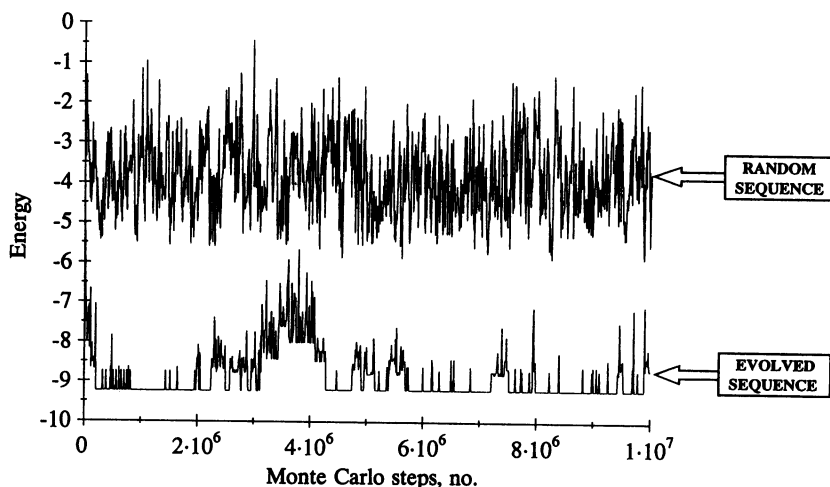*Proc. Natl. Acad. Sci. USA* 93 (1996)     843



FIG. 5.  Folding trajectories of the starting random sequence and for the last evolved sequence of 27-mer ($B_0 = 0.05$ and $T = 0.15$).

ribozymes with unique native structure have evolved. The model described above can be used to understand better the driving forces of prebiological creation of any biopolymers with stable native structure.

The selection algorithm that we used herein may resemble biological selection rather than prebiological selection since some favorable substitutions were accepted and passed by to further "generations" of sequences. This resemblance is superficial. The possible feedback in prebiological evolutionary scenario could be that more compact sequences were more likely to avoid hydrolysis. In the situation when sequences got synthesized randomly, this factor gradually shifted the sequence distribution to increase probability of stable sequences. The possible physical mechanism of it is very simple: in spite

of the fact that probability to randomly synthesize stable sequence is very small, these sequences were not hydrolyzed and accumulated over many cycles of hydrolysis and synthesis; whereas more frequently occurring unstable sequences underwent many turns of recycling. At each recycling step, there existed a (small) probability to synthesize a stable sequence.

Synthesis of random sequences represents a random walk in sequence space. Within the framework of this analogy, the stability of folded structures, which make sequences nonhydrolyzable, is equivalent to absorbing walls; it is clear that the concentration of species at or near adsorbing boundaries is increased. This gives an equivalent explanation of the observed phenomenon and suggests the way to construct an analytical model of prebiological evolution.

Another striking result is that, though the fraction of sequences with a stable unique structure is very small, the algorithm generated them after a small number of mutations. This is the common feature of evolution-like algorithms: selection pressure generates a directed drift in sequence space that makes generation of desirable (in our case compact) chains much more feasible than if such sequences were searched for randomly.

Our model does not consider such important properties of proteins as chirality and enzymatic activity. But we found that the requirement of compactness without aggregation can cause a stable three-dimensional native conformation, without which enzymatic activity is unlikely to be possible. The model presented in this work, while being unlikely to reproduce the
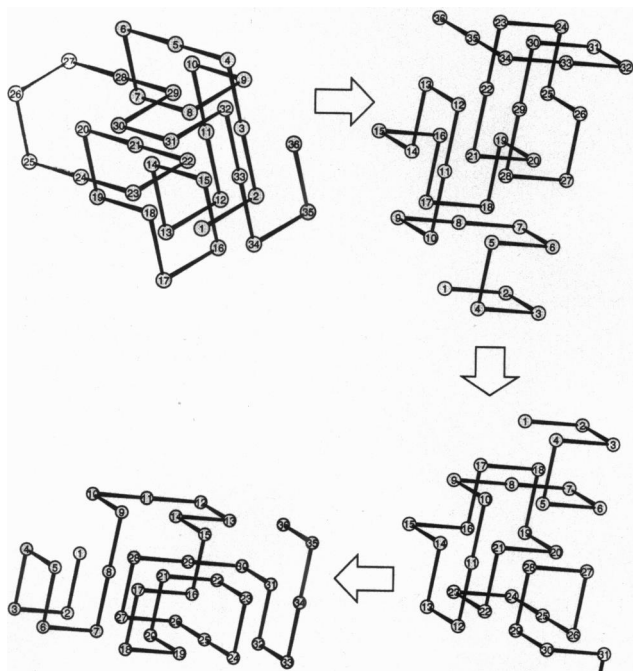


FIG. 6.  Selection of sequences with 36 residues under the requirement of large average compactness over a million Monte Carlo steps starting with a random coil conformation. Evolution of the native structure is shown. The native conformations of the original random sequence, sequence after the 15th substitution, sequence after the 30th substitution, and sequence after the last (47th) substitution are shown.
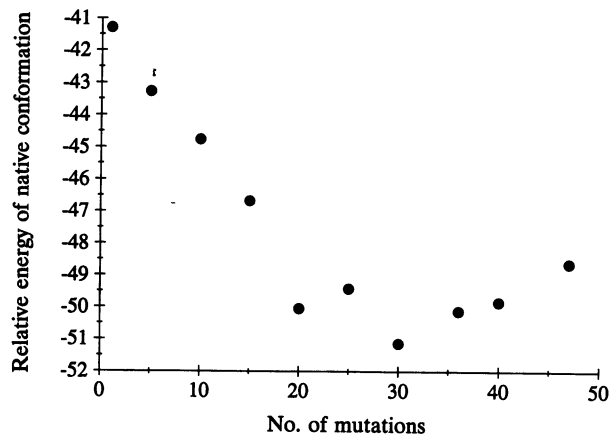


FIG. 7.  Evolution of the relative energy of the native conformation of 36-mers.

details of prebiotic evolution, illustrates a general idea that a stable native structure could have evolved as a side effect of equilibrium conditions that selected sequences satisfying simple physicochemical requirements.

1. Eigen, M. & Schuster, P. (1979) *The Hypercycle: A Principle of Natural Self-Organization* (Springer, Berlin).
2. Karplus, M. & Shakhnovich, E. (1992) in *Protein Folding*, ed. Creighton, T. E. (Freeman, New York), pp. 127–195.
3. Chan, H. S. & Dill, K. A. (1993) *Phys. Today* **46**, 24–32.
4. Frauenfelder, H. & Wolynes, P. G. (1994) *Phys. Today* **47**, 58–64.
5. Taketomi, H., Ueda, Y. & Go, N. (1975) *Int. J. Pept. Protein Res.* **7**, 445–459.
6. Lau, K. F. & Dill, K. A. (1989) *Macromolecules* **22**, 3986–3997.
7. Covell, D. & Jernigan, R. (1990) *Biochemistry* **29**, 3287–3294.
8. Shakhnovich, E. I. & Gutin, A. M. (1990) *Nature (London)* **346**, 773–775.
9. Skolnick, J. & Kolinski, A. (1990) *Science* **250**, 1121–1125.
10. Sali, A., Shakhnovich, E. I. & Karplus, M. (1994) *J. Mol. Biol.* **235**, 1614–1636.
11. Goldstein, R., Luthey-Schulten, Z. A. & Wolynes, P. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9029–9033.
12. Verdier, P. H. (1973) *J. Chem. Phys.* **59**, 6119–6126.
13. Hilhorst, H. J. & Deutch, J. M. (1975) *J. Chem. Phys.* **63**, 5153–5161.
14. Shakhnovich, E. I. & Gutin, A. M. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7195–7199.
15. Gutin, A. M., Abkevich, V. I. & Shakhnovich, E. I. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1282–1286.
16. Gutin, A. M. & Shakhnovich, E. I. (1993) *J. Chem. Phys.* **98**, 8174–8177.
17. Shakhnovich, E. I. & Gutin, A. M. (1993) *Protein Eng.* **6**, 793–800.
18. Shakhnovich, E. I. (1994) *Phys. Rev. Lett.* **72**, 3907–3909.
19. Abkevich, V. I., Gutin, A. M. & Shakhnovich, E. I. (1994) *Biochemistry* **33**, 10026–10036.
20. Grosberg, A. & Khohlov, A. (1985) *Unsolved Problems in Physics of Polymers*, preprint.
21. Miyazawa, S. & Jernigan, R. (1985) *Macromolecules* **18**, 534–552.
22. Finkelstein, A. V., Gutin, A. M. & Badretdinov, A. Ya. (1993) *FEBS Lett.* **325**, 23–28.
23. Bowie, J. U., Luthy, R. & Eisenberg, D. (1991) *Science* **253**, 164–169.
24. Joyce, G. F. (1989) *Nature (London)* **338**, 217–224.
25. Benner, S. A., Ellington, A. D. & Tauer, A. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 7054–7058.
26. Wilson, C. & Szostak, J. W. (1995) *Nature (London)* **374**, 777–782.