

Local Rules for Protein Folding on a Triangular Lattice and Generalized Hydrophobicity in the HP Model

Richa Agarwala* Serafim Batzoglou† Vlado Dančik‡ Scott E. Decatur§
Martin Farach¶ Sridhar Hannenhalli|| Steven Skiena**

Abstract

We consider the problem of determining the three-dimensional folding of a protein given its one-dimensional amino acid sequence. We use the HP model for protein folding proposed by Dill [3], which models protein as a chain of amino acid residues that are either hydrophobic or polar, and hydrophobic interactions are the dominant initial driving force for the protein folding. Hart and Istrail [5] gave approximation algorithms for folding proteins on the cubic lattice under HP model. In this paper, we examine the choice of a lattice by considering its algorithmic and geometric implications and argue that triangular lattice is a more reasonable choice. We present a set of folding rules for a triangular lattice and analyze the approximation ratio which they achieve. In addition, we introduce a generalization of the HP model to account for residues having different levels of hydrophobicity. After describing the biological foundation for this generalization, we show that in the new model we are able to achieve similar constant factor approximation guarantees on the triangular lattice as were achieved in

the standard HP model. While the structures derived from our folding rules are probably still far from biological reality, we hope that having a set of folding rules with different properties will yield more interesting folds when combined.

1 Introduction

A long standing problem in molecular biology is to determine the three-dimensional structure of a protein when only given the sequence of amino acid residues which compose the protein chain. Due to the complexity of the protein folding problem, scientists have proposed a variety of models which attempt to simplify the problem by abstracting only the “essential physical properties” of real proteins. In these models, the three dimensional space is often represented by a *lattice*. Residues which are adjacent in the primary sequence (*i.e.* covalently linked) must be placed at adjacent points in the lattice. A *conformation* of a protein is simply a self-avoiding walk along the lattice. The protein folding problem *STRING-FOLD* is that of finding a conformation of the protein sequence on the lattice such that the overall *energy* is minimized, for some reasonable definition of energy. The lattice formulation described so far leaves open two important questions: what type of lattice should one use; and what energy function is appropriate. Once these two questions have been answered, one may then address the algorithmic complexity of optimizing the energy function for the lattice. For a variety of such simple models, this minimization problem is in fact NP-hard [9, 8, 11].

In this paper, we consider the *Hydrophobic-Polar (HP) Model* introduced by Dill [3]. The HP model abstracts the problem by grouping the 20 amino acids which compose proteins into two classes: hydrophobic (or non-polar) residues and hydrophilic (or polar) residues. For concreteness, we will take our input to be a string from $\{H, P\}^+$, where P represents polar residues, and H represents hydrophobic residues. Dill *et.al.* [2] survey the literature analyzing this model.

1.1 Selecting an energy function. Given a conformation, we say that a pair of residues form a *topological contact* (or simply contact) if the residues are not covalently linked and are placed on neighboring lattice

*National Center for Human Genome Research/National Institutes of Health, Bethesda, MD 20892, (richa@helix.nih.gov) under a contract with R.O.W. Sciences, Inc. Most of this work was done while this author was at DIMACS center, Rutgers University and was supported by Special Year National Science Foundation grant BIR-9412594.

†MIT Laboratory for Computer Science and Department of Mathematics, 545 Technology Square, Room 342, Cambridge, MA 02139. (serafim@theory.lcs.mit.edu) Supported by a DOE Contract.

‡Department of Mathematics, University of Southern California, Los Angeles, CA 90089-1113. (dancik@hto.usc.edu) Supported by NIH grant GM36230.

§MIT Laboratory for Computer Science and Department of Mathematics, 545 Technology Square, Room 313, Cambridge, MA 02139. (sed@theory.lcs.mit.edu) Supported by a grant from the Reed Foundation through the MIT School of Science.

¶Department of Computer Science, Rutgers University, Piscataway, NJ 08855. (farach@cs.rutgers.edu) Supported by NSF Career Development Award CCR-9501942, an Alfred P. Sloan Research Fellowship, and NATO Grant CRG 960215.

||Department of Mathematics, University of Southern California, Los Angeles, CA 90089-1113. (hannenha@hto.usc.edu) Supported by NSF Young Investigator Award, NIH grant 1R01 HG00987 and DOE grant DE-FG02-94ER61919.

**Department of Computer Science, State University of New York, Stony Brook, NY 11794-4400. (skiena@cs.sunysb.edu) Supported by ONR award 400x116yip01 and NSF Grant CCR-9625669.

points. A *bond* refers to a topological contact between a pair of H's. Define the *free energy* of a conformation as $(-1) \times (\# \text{ of bonds})$. The optimal conformation for the protein is the one which has the lowest free energy.

The biological foundation of this energy function is the belief that the first-order driving force of protein folding is due to a "hydrophobic collapse" in which those residues which prefer to be shielded from water (hydrophobic residues) are driven to the core of the protein, while those which interact more favorably with water (polar residues) remain on the outside of the protein. The protein is hypothesized to fold in such a way as to minimize the surface area of hydrophobic residues exposed to water or polar residues.

1.2 Selecting a lattice. HP simulations have typically followed Dill's original choice of square lattices. Unfortunately, one rather severe consequence of the structure of the square lattice is that no two amino acids can be in adjacent lattice points if the string between them is of odd length. We call this the *parity constraint* of square lattices. Thus, the string $(PH)^n$ has no bonds in a square lattice, despite the fact that such a protein string has many potential bonds in "real space".

The bizarreness of the parity constraint illustrates another possible pitfall of algorithmic analysis of this problem. An approximation ratio for a maximization algorithm is the ratio of a lower bound on the performance of the algorithm and an upper bound on the optimal solution. Thus, it is desirable to raise the lower bound of the algorithmic solution, or to lower the upper bound of the optimal solution. But in the case of the *STRING-FOLD* problem on a square lattice, the upper bound is artificially low, due to the parity constraint. To illustrate, consider once again the string $(PH)^n$, which has a trivial upper bound of 0 bonds, so any algorithm will achieve the optimum. Thus, any approximation ratio on a square lattice will have little meaning for moving towards a "realistic" solution of the problem.

If the square lattice is so seriously flawed, what is a more interesting lattice choice? We propose using triangular lattices as a folding model. The two- and three-dimensional triangular lattices are shown in Figure 1 and Figure 8.¹ The triangular lattice does not exhibit the parity problem. That is, for every pair of sites x, y on any string, there exists a confirmation of the string on the triangular lattice such that x and y are neighboring sites on the lattice.

¹This three-dimensional lattice is based on the topology of the β form of Silicon Carbide. The nodes in this lattice correspond to the silicon atoms in the Silicon Carbide crystal. Two nodes in this lattice are connected if there exists a carbon atom that is bonded to both of the two silicon atoms corresponding to the nodes.

1.3 Past results. From a computational point of view, it is not known whether or not the problem of finding an optimal conformation using the above specified energy function (on either square or triangular lattices) is NP-hard. The problem is known to be NP-complete when the alphabet size is unbounded and the lattice is a two- or three-dimensional square lattice [10]. Hart and Istrail [5] presented approximation algorithms on the square lattice having approximation factor of $1/4$ on the two-dimensional square lattice and $3/8$ in three-dimensions but the optimal conformations on these lattices may be arbitrarily worse than the optimal on lattices without the parity problem. It is therefore interesting to examine the HP folding problem on triangular lattices, and to strive for conformations approaching the more natural optimal score found there.

1.4 Our results. We present a collection of local folding rules for the HP model on triangular lattices, in both two and three dimensions, and we prove approximation ratios for each of these rules. For all of our rules, these ratios are better than those achieved in by Hart and Istrail [5] for the square lattice. As pointed out above, an approximation ratio can often be a misleading. We provide these numbers as a rough guideline, since no more appropriate measure of goodness for a rule is available. Yet, since optimal solutions for triangular lattices are so much more densely packed than are optimal solutions for square lattices, we achieve our approximations by local rules which yield very dense structures. That is, since the upper bound of the optimal solution is much higher, the performance of the folding algorithms must increase correspondingly.

In Sections 2 and Sections 3, we describe the folding rules for the two- and three-dimensional lattices, respectively. All of our folding rules are implementable in linear time. The following table summarizes the proposed rules for protein folding under the HP model in the two- and three-dimensional triangular lattices.

Folding	2D	3D
backbone	1/4	3/10
improved backbone	1/2	2/5
arrow	1/2	14/30
improved arrow	6/11	n/a
star	n/a	16/30
combined backbone-star	n/a	44/75
improved star	n/a	3/5

In Section 4, we extend the HP model by considering a more general representation of hydrophobic residues. The new model is motivated by the fact that certain hydrophobic residues are more hydrophobic than others. While in the standard HP model all hydrophobic residues have identical hydrophobicities, our new model

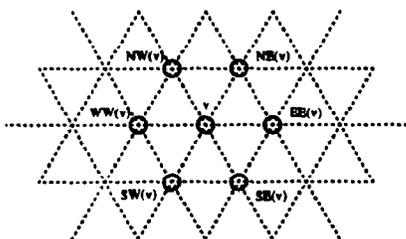


Figure 1: The two-dimensional triangular lattice

allows different residues to have different hydrophobicities and contacts between hydrophobic residues contribute to the energy function proportional to their combined hydrophobic strength. After describing the biological foundation for this generalization, we show that in the new model we are able to achieve similar constant factor approximation guarantees on the triangular lattice as were achieved in the standard HP model.

2 Rules for Two Dimensions

The two dimensional triangular lattice is shown in Figure 1. Every node in this lattice has 6 neighbors. For a node v , denote the six neighbors of v as shown in Figure 1.

For a binary sequence S , let s_1 be the number of H's in S . For simplicity, assume that S begins and ends with a P. For a conformation f of S , let $C(S, f)$ denote the number of bonds. We use s_1 to bound the number of bonds.

LEMMA 2.1. *For every conformation f of binary sequence S in the two dimensional triangular lattice*

$$C(S, f) \leq 2 \cdot s_1.$$

Proof. Since any internal residue (not corresponding to the beginning or the end of the sequence) can have at most 4 topological contacts in the two dimensional triangular lattice (out of six neighbors, two are connected along the sequence), $C(S, f) \leq 2 \cdot s_1$. ■

Notice that an H at the beginning or at the end of the sequence can have 5 contacts leading to $C(S, f) \leq 2 \cdot s_1 + 1$. This situation is avoided due to the simplifying assumption.

2.1 Backbone-folding: a 1/4-Approximation.

We describe a very simple folding strategy based on laying out two *backbones* anti-parallel to each other such that most of the bonds are achieved among the H's along a backbone and the H's across the backbones. The general scheme of backbone folding is shown in Figure 2. We refer to this conformation as *backbone-folding*. This is the folding rule used in [6].

Figure 3 describes the rules for backbone folding of

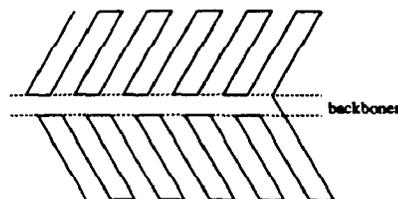


Figure 2: The backbone-folding

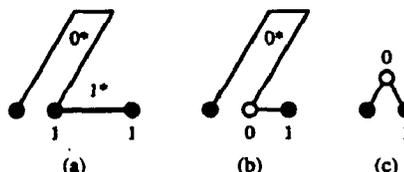


Figure 3: The conformation of $P^i H^j$, $i, j > 0$ when (a) $j > 1$, (b) $j = 1$, $i > 1$, (c) $j = 1$, $i = 1$. White circles represent P, black circles represent H, and gray circles represent H preceding P^i . Thick line represents sequence of $j - 2$ H's.

sequence $P^+ H^+$. It is easy to see that these rules create a conformation satisfying the following two conditions: (1) every H is placed on the backbone, (2) two P's on the backbone never form topological contacts.

LEMMA 2.2. *Let f^b be the backbone-folding of binary sequence S . Then*

$$C(S, f^b) \geq \frac{1}{2} s_1.$$

Proof. Each residue of the backbone has two contacts with the other backbone and one of those has to be a bond. All H's are on the backbones, each creates at least one bond and therefore $C(S, f^b) \geq \frac{1}{2} s_1$. ■

Lemmas 2.1 and 2.2 together imply

COROLLARY 2.1. *Backbone-folding yields a 1/4-approximation for the STRING-FOLD problem in the HP model for two-dimensional triangular lattices for large values of s_1 .*

2.2 Improved Backbone-folding: a 1/2 Approximation. Using a more careful placement of H's (allowing H's to be placed outside the backbones), we can improve the backbone folding to achieve 1/2-approximation.

COROLLARY 2.2. *Backbone-folding yields a 1/2-approximation for STRING-FOLD in the HP model on two dimensional triangular lattice for large values of s_1 .*

2.3 Arrow-folding: a 1/2 Approximation. Define a *strip* as a maximal substring of S containing only H's. A k -*strip* is a strip of size k . For example, HHHH is a 4-strip starting at position 2 of $S =$

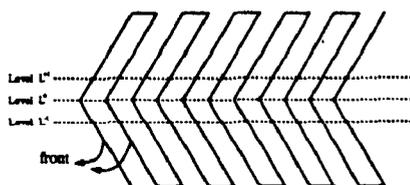


Figure 4: The arrow-folding

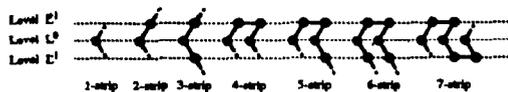


Figure 5: Placement of strips in the arrow-folding

PHHHHPHPPPPHHPHHHHHH. Define a *separator* to be a maximal substring of S containing only P's and a k -*separator* likewise. For example, PPPP is a 4-separator between 1-strip and 2-strip in the previous example. A k -strip (k -separator) is *long* if $k > 1$. In the following we describe a certain conformation which we call *arrow-folding* and show that for this conformation, we get at least s_1 bonds, leading to a $1/2$ -approximation.

Figure 4 shows the general scheme of arrow-folding. Define levels L^0 , L^{+1} and L^{-1} as shown in Figure 4. We fold a sequence along this pathway such that (i) all the H's are distributed only in L^0 , L^{+1} and L^{-1} and (ii) all the residues on level L^0 are H's. This is achieved by placing the strips as shown in Figure 5.

In the figure we show the case of k -strips, $1 \leq k \leq 7$, but it follows a general pattern as is evident in the figure. For $k \equiv 0 \pmod{3}$, both the first and the last residues of the strip are on level L^{+1} or L^{-1} . For $k \equiv 1 \pmod{3}$, both the first and the last residues of the strip are on level L^0 . For $k \equiv 2 \pmod{3}$, the first residue of the strip is on level L^0 and the last residue of the strip is on level L^{+1} or L^{-1} . Notice that a mirror image with respect to the L^0 line does not affect the general schema of folding. We connect two consecutive strips using the separators such that the entire sequence still follows the general schema dictated by arrow-folding. The precise description of arrow-folding is cumbersome and unnecessary. Figure 6 shows an example of arrow-folding.

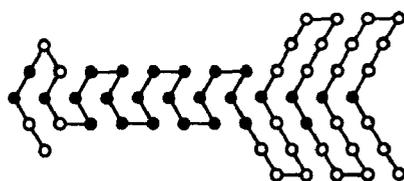


Figure 6: Example of arrow-folding for the sequence (PHHPHPPHPPHHHHHHHHHHHHHHHHHPPPPPHPPPPPHHPPPPHPPPPPHPP).

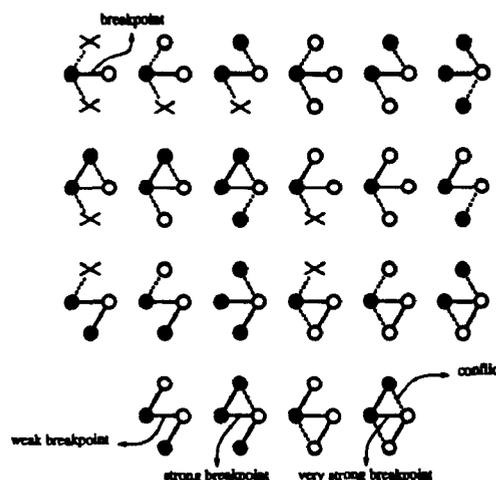


Figure 7: Conformation with breakpoints and conflicts; Solid lines represent the folding pathway; Dashed lines represent conflicts; Crosses represent absence of a residue of S .

LEMMA 2.3. Let f^a be the arrow-folding of binary sequence S . Then

$$C(S, f^a) \geq s_1 - 3.$$

Proof. To compute the number of bonds, notice that every H-residue v on level L^0 pairs exclusively to the H on L^0 to its right, i.e. $EE(v)$ (except for the terminal case). Also, every H-residue v at level L^{+1} (L^{-1}) pairs exclusively to the H-residue $SE(v)$ ($NE(v)$) at level L^0 . Simple counting gives the number of bonds as $s_1 - 3$. ■

Lemmas 2.1 and 2.3 together imply:

COROLLARY 2.3. Arrow-folding yields a $1/2$ -approximation for the STRING-FOLD problem in the HP model for two-dimensional triangular lattices for large values of s_1 .

2.4 Improved Arrow-folding: a 6/11 Approximation. Even when folding sequence S optimally, it can happen that some internal residue H has less than 4 bonds, i.e it forms a contact with an internal residue P or does not form a contact at all. We call such absence of a potential bond a *conflict*. Let c_f be the number of conflicts in conformation f . Since each conflict excludes a bond, the number of bonds formed by conformation f is

$$C(S, f) = (4s_1 - c_f)/2.$$

Conflicts in conformations are tightly related to *breakpoints* in the folded sequence S . A breakpoint is a position in S at which a P is followed by a H or a H is followed by a P. Figure 7 shows all possibilities for two triangles adjacent to a breakpoint. Not all breakpoints in S force conflicts in conformations, *weak*

breakpoints HP-HP (or PH-PH) can be folded without creating conflicts, *Strong* breakpoints PP-HP (or HH-PH) in any conformation always create at least one conflict and *very strong* breakpoints PP-HH (or HH-PP) create two conflicts in any conformation.

Breakpoints P-HP..., ...PH-P (P-HH..., ...HH-P) at the beginning and at the end of S are strong (very strong). Let b_S be the number of strong breakpoints plus twice the number of very strong breakpoints. The same conflict can be forced by at most two different non-weak breakpoints, therefore $c_f \geq b_S/2$. Thus we have improved upper bound on $C(S, f)$:

LEMMA 2.4. For every conformation f of binary sequence S in the two dimensional triangular lattice we have

$$C(S, f) \leq 2s_1 - b_S/4.$$

Note that each long strip(or long separator) contributes 2 to the value of b_S (the first and the last strip of P's contribute only 1).

For certain parts of the sequence, a better analysis of arrow-folding leads to an approximation ratio better than 1/2. For other parts of the sequences we need to modify the "regular" arrow-folding to achieve a better approximation.

Suppose S is arbitrarily divided into k segments S_1, \dots, S_k , $S = S_1 \dots S_k$. Let C_1, \dots, C_k be the corresponding division of bonds (C_i is the number of bonds contributed exclusively by the \mathbb{H} in S_i), $C_1 + \dots + C_k = C$. We denote the number of \mathbb{H} 's in S_i by t_i and the number of breakpoints in S_i exclusively due to S_i by b_i . Let $r_i = \frac{C_i}{2t_i - b_i/4}$ be the ratio achieved by conformation f for sequence S_i . From the fact that $(a+c)/(b+d) \geq \min\{a/b, c/d\}$ we have

$$\begin{aligned} \frac{C(S, f)}{2s_1 - b_S/4} &= \frac{C_1 + \dots + C_k}{(2t_1 - b_1/4) + \dots + (2t_k - b_k/4)} \\ &\geq \min\{r_1, r_2, \dots, r_k\}. \end{aligned}$$

We scan the sequence left to right, cut the sequence at certain place and analyze the cut prefix. In any particular iteration we examine at most 4 strips and 4 separators ahead to decide on the place to make the next cut. Let the prefix of S containing 4 strips and 4 separators be $P^A \mathbb{H}^B P^C \mathbb{H}^D P^E \mathbb{H}^F P^G \mathbb{H}^H$. We decide on the position to cut depending on the values of A through H . An extensive case analysis of the different values of A through H (omitted here for brevity) yields worst ratio among all the cases of 6/11. We refer to the conformation thus obtained as *improved-arrow-folding*.

COROLLARY 2.4. *Improved-arrow-folding yields a 6/11-approximation for the STRING-FOLD in the HP model on a two dimensional triangular lattice for large values of s_1 .*

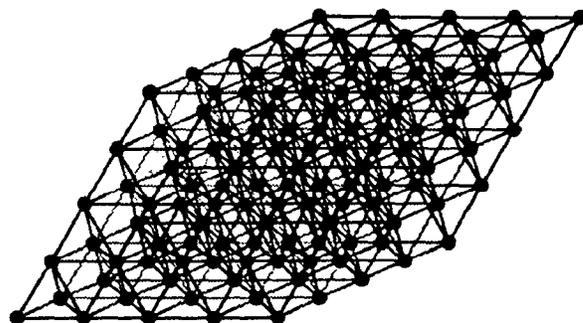


Figure 8: The Three-Dimensional Triangular Lattice.

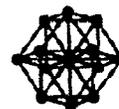


Figure 9: A single node and its twelve neighbors.

3 Rules for Three Dimensions

Analogous to the two dimensional case, the three dimensional triangular lattice mimics the configuration of equal sized solid spherical objects packed in three dimensional space. It can be described as a stack of two dimensional triangular lattices where every individual lattice is slightly offset with respect to the lattices above and below it (Figure 8). As in the two dimensional case, we assume that the sequence starts and ends with a \mathbb{H} .

LEMMA 3.1. For every conformation f of S on a three dimensional triangular lattice

$$C(S, f) \leq 5 \cdot s_1.$$

Proof. Every node has 12 neighbors in this lattice, as shown in Figure 9. Since any internal residue can have at most 10 topological contacts in the three dimensional triangular lattice (out of 12 neighbors, two are connected along the sequence), clearly, $C(S, f) \leq 5 \cdot s_1$. ■

3.1 Backbone-folding^{3D}: a 3/10 approximation. Since the three dimensional triangular lattice is a stack of two dimensional triangular lattices, an obvious extension of the backbone-folding for the three dimensional triangular lattice would look like a stack of backbone-foldings (*backbone-folding^{3D}*). We could fold the entire sequence into the two dimensional lattice, divide it into a few segments and stack up the conformations of individual segments on consecutive parallel planes. Instead of two parallel lines we now have two parallel backbone planes containing all \mathbb{H} 's of the sequence.

LEMMA 3.2. Let f^b be the 3D backbone folding of sequence S on the three dimensional triangular lattice. Ignoring the boundary conditions

$$C(S, f^b) \geq \frac{3}{2} \cdot s_1.$$

Proof. The two parallel backbone planes consist of the parallel backbones of the segments. Each H on the i -th backbone creates at least 3 bonds, two with $(i-1)$ -th and $(i+1)$ -th backbones in the same plane and one with the i -th backbone in the opposite plane. ■

Lemmas 3.1 and 3.2 together imply

COROLLARY 3.1. Backbone-folding^{3D} yields a 3/10-approximation for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .

3.2 Improved Backbone-folding^{3D}: a 2/5 approximation. In the following we describe backbone-folding and show that for this conformation, we get at least $2 \cdot s_1$ bonds, leading to a 2/5-approximation. As in the two dimensional case, we can lay out any number of P's away from the backbone and construct a backbone of all the 1's in the protein. In three dimensions, we can construct a triangular "tube" from this backbone by folding it into three backbones such that one strand runs anti-parallel to the other two strands.

LEMMA 3.3. Let f^b denote backbone-folding of S with tube T on the three dimensional triangular lattice. Ignoring the boundary conditions

$$C(S, f^b) \geq 2 \cdot s_1.$$

Proof. In f^b , if x, y, z are neighbors and form a triangle along a cross section of T , then all the three pairs xy, yz, xz form bonds. Ignoring the boundary condition where we are folding the backbone to form T , each H on a triangular lattice i forms 2 bonds with backbone strands in both lattice $(i-1)$ and lattice $(i+1)$. Hence, each H forms at least 4 bonds. ■

COROLLARY 3.2. Improved backbone-folding^{3D} yields a 2/5-approximation for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .

3.3 Arrow-folding^{3D}: a 7/15 approximation. We extend two dimensional arrow-folding into three dimensions in the same manner as we did with backbone folding. We divide the entire sequence into a few segments and fold each segment according to the arrow-folding rule for the two dimensional lattice and stack up the conformations of individual segments on consecutive parallel planes.

LEMMA 3.4. Let f^b be the 3D arrow folding of sequence S on the three dimensional triangular lattice. Ignoring the boundary conditions

$$C(S, f^b) \geq \frac{7}{3} \cdot s_1.$$

Proof. All H's of sequence S are on one of level planes L^{-1}, L^0 or L^{+1} and there are only H's on the plane level L^0 . Each H on level L^0 forms three exclusive bonds, one with a H in the same segment and two with H's at level L^0 of the next segment. Each H on levels $L^{\pm 1}$ forms two exclusive bonds, one with a H at level L^0 of the same segment and one with a H at level L^0 of the next segment. Each front of arrow folding can contain one, two or three H's. In the first case we have three exclusive bonds per one H, in the second case we have five exclusive bonds per two H's and in the third case we have seven exclusive bonds per three H's. ■

Lemmas 3.1 and 3.4 together imply

COROLLARY 3.3. Arrow-folding^{3D} yields a 7/15-approximation for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .

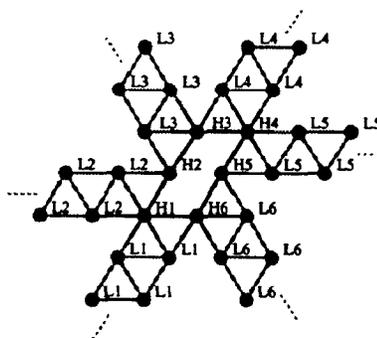


Figure 10: Individual plane in star-folding.

3.4 Star-folding: a 16/30 approximation. In the following we describe star-folding and show that for this conformation, we get at least $\frac{8}{3}s_1$ bonds, leading to a 16/30-approximation. We place all H's in a single, tightly-packed, hydrophobic core. The core is composed of six H's at each level and contains as many levels as is necessary to accommodate all H's in the sequence. The pattern for one level, containing six H's in the core and possibly ladders of P's is shown in Figure 10. Then we place these individual stars on consecutive planes in a slightly offset manner as shown in Figure 11 by connecting one end of a star with an end of the star at the level above and the other end to an end of the star at the level below.

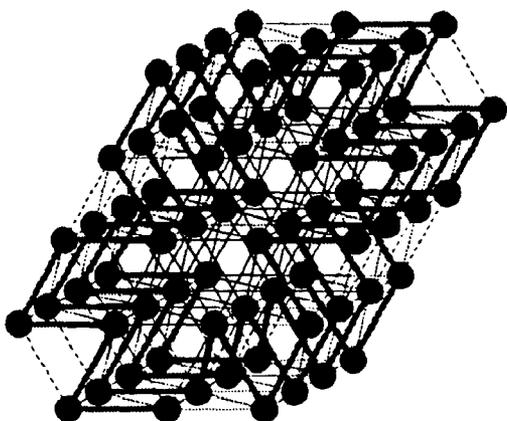


Figure 11: Four levels of the 16/30 factor layout.

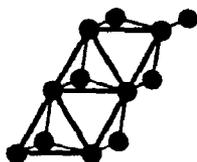


Figure 12: Hydrophobic contacts within one plane and to the next plane.

LEMMA 3.5. Let f' denote star-folding of S on the three dimensional triangular lattice. Ignoring the boundary conditions

$$C(S, f') \geq \frac{16}{6} \cdot s_1.$$

Proof. To compute the number of bonds, we will look at the six H's on a certain plane and count the number of bonds exclusive to these six H's. As shown in Figure 12, each group of six H's in the constructed core has 9 bonds within the group (thick lines) and 13 bonds to the next group (thin lines). At most 6 of these bonds could be covalent bonds. Thus, each group of six H's contributes at least $9 + 13 - 6 = 16$ exclusive bonds and $C(S, f') \geq \frac{16}{6} s_1$. ■

Lemmas 3.1 and 3.5 together imply

COROLLARY 3.4. Star-folding yields a 16/30-approximation algorithm for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .

3.5 Combined (Backbone-Star) folding: a 44/75 approximation. To improve approximation ratio we derive better upper bound. Weak, strong and very strong breakpoints are defined in the same way as in section 2.4. Let b_S be twice the number of weak breakpoints plus three times the number of strong breakpoints plus four times the number of very strong

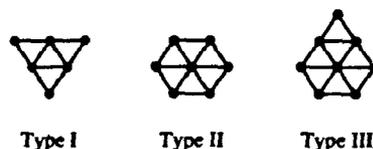


Figure 13: Cluster types.

breakpoints. Now we can bound the maximal number of bonds in the same manner as we did in Lemma 2.4.

LEMMA 3.6. For every conformation f of binary sequence S in the three dimensional triangular lattice we have

$$C(S, f) \leq 5s_1 - b_S/8.$$

Closer look on the backbone-folding reveals that a long separator of P's followed by a single H is the most unfavorable case. On the other side, single H's form extra bonds in the star-folding. We can decide the type of folding depending on the number of single H's. This yields 44/75 approximation.

Suppose the number of H's in the sequence S is s_1 and t_1 of them are single H's preceded by long separator of P's. Each such long separator contributes at least 6 to the b_S and for every conformation of S we have $C(S, f) \leq 5s_1 - 6t_1/8$. The backbone folding for S forms $C(S, f^b) = (7s_1 - 7t_1)/2$ bonds. In the star-folding f^s each single H forms two extra bonds, thus we have $C(S, f^s) = \frac{16}{6}s_1 + t_1$. For $t_1 < \frac{5}{27}s_1$ the backbone folding is more favorable, while for $t_1 > \frac{5}{27}s_1$ the star folding produces better approximation ratio. The worst case is when $t_1 = \frac{5}{27}s_1$ and we get

COROLLARY 3.5. Combined (backbone-star) folding yields a 44/75-approximation algorithm for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .

3.6 Improved Star folding: a 3/5 approximation. Note that the algorithm of Section 3.4 constructs a fold such that the position of the i -th hydrophobic residue in the sequence depends *only* on i . Specifically, the positions of the hydrophobic residues do not depend at all on the polar residues in the sequence. The algorithm presented in this section instead takes into account the polar residues in determining the fold, and is able to achieve an improved approximation factor. We first describe the algorithm and then analyze its performance.

As before, we construct a leveled hydrophobic core, but in the new algorithm the levels may differ as a function of the local makeup of hydrophobic versus polar residues. The three types of levels, or *clusters*, used are shown in Figure 13.

In Figure 14, we show the locations of the polar

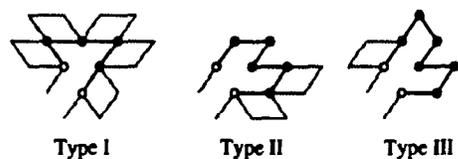


Figure 14: Ladders and pivots for the clusters.

ladders for each of the types as well as the entry and exit point for each type. The ladders are shown in light lines while the dark lines represent the sequence of H's in the cluster. Note that of the H's shown might not be covalently linked. The entry and exit H's (*pivot points*) are shown with empty circles. As is evident from Figure 14, any subsequence starting and ending with an H and containing 6 H's can always be folded into a Type I cluster. Consider a subsequence containing 7 H's, starting and ending with an H, and containing either one or two groups of consecutive P's between the H's. Such a subsequence must contain 3 consecutive H's, and can always be folded into a Type II cluster. Finally, a subsequence containing 8 H's and starting with 7 consecutive H-residues can always be folded into a Type III cluster. Following algorithm contains the rules for folding an arbitrary H/p protein sequence using the cluster types described above.

1. Let p be the number of intervening P groups in the subsequence that contains the next 7 H's.
2. If $p \geq 3$ fold the next 6 H's into a Type I cluster².
3. If $p = 1$ or $p = 2$ fold the next 7 H's into a Type II cluster.
4. If $p = 0$ (and therefore the next 7 H's are consecutive):
 - If the previous cluster is of Type I, fold the next 7 H-residues into a Type II cluster.
 - If the previous cluster is of Type II or III, fold the next 8 H's into a Type III cluster.
5. Continue again with Step 1.

Figure 15 displays all of the possible arrangements of one cluster type following another. The light shapes denote clusters at the previous level and the dark shapes represent clusters at the current level. A cluster of one level connects to a cluster of the next level using the pivot points shown in Figure 14. More specifically, a cluster is folded clockwise or counterclockwise depending on whether it started in the upper, or lower pivot point of Figure 14, respectively. Thus it reaches the other pivot point, which is then used to connect to the next level. The next level starts at the lower pivot point if the previous level ended at the lower pivot point, and similarly for the upper pivot point. A ladder, defined by the two pivot points where the two clusters connect, and

²The 7th H is left for the next level.

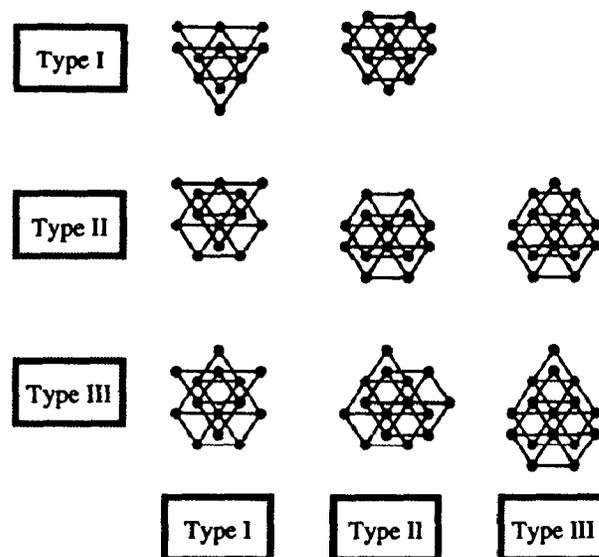


Figure 15: Arrangements of clusters from one level to the next.

the direction shown in Figure 14, is used to accommodate an arbitrary number of P's between the boundary H's of the two clusters.

Let *inner* contacts of a cluster be the non-covalent bonds among the H's of the cluster, plus possibly the non-covalent bond between the exit H of this cluster and the entry H of the next cluster. Let *outer* contacts of a cluster be the contacts that the cluster achieves with the previous cluster, minus one (between the entry and exit H-residue). A Type I cluster contains 9 contacts within the cluster plus one contact between the pivot H's of this and the next level. Since there are at least 3 intervening polar groups in the next 7 H's subsequence starting with the entry H of this cluster, at most 3 of these 10 contacts may be covalent. Thus, a Type I cluster contributes at least 7 inner contacts. Similarly a Type II cluster achieves at least $7(= 12 + 1 - 6)$ inner contacts except when placed on top of a Type I cluster, when it can achieve as few as $6(= 12 + 1 - 7)$ inner contacts.³ A Type III cluster achieves at least $7(= 14 + 1 - 8)$ inner contacts. From the diagrams in Figure 15 we can enumerate the outer contacts each additional cluster is guaranteed to contribute. For example, placing a Type I cluster on top of a Type II cluster achieves 13 contacts, of which one may be a covalent bond between the two endpoint H-residues. Thus, it achieves at least 12 outer contacts. The enumeration of inner and outer contacts is summarized in Table 1.

Notice that any cluster on top of a cluster of the

³This exception is due to the rule in which we put a Type II cluster on top of a Type I cluster even if the next 7 H-residues contain no intervening polar groups.

New Residues	Next Cluster Type		
	Type I	Type II	Type III
	6	7	8
On Type I	18 = 11+7	20 = 14+6	—
On Type II	19 = 12+7	21 = 14+7	23 = 16+7
On Type III	20 = 13+7	22 = 15+7	24 = 17+7

Table 1: Number of contacts (outer+inner) achieved by each combination of types.

same type achieves at least 3 contacts per R. Say that we *ascend* a level whenever we move to a cluster of higher type (i.e. Type I to Type II, or Type II to Type III). Similarly say we *descend* one or two levels, whenever we move to a cluster of lower type. Notice in Table 1, that each time we ascend a level, we are one contact short of achieving 3 contacts per residue, whereas each time we descend one or two levels, we achieve 1 or 2 contacts more, respectively, than 3 contacts per residue. For the fold of one sequence, the total number of levels we ascend, has to equal the total number of levels we descend (\pm at most 2). Thus this algorithm achieves asymptotically 3 contacts per residue.

COROLLARY 3.6. *Improved Star-folding yields a 3/5-approximation algorithm for the STRING-FOLD problem in the HP model on a three dimensional triangular lattice for large values of s_1 .*

4 Generalized Hydrophobicity

The standard HP model makes the simplifying assumption that all hydrophobic residues have the same energy contribution to the hydrophobic collapse. Yet it is well known that various residues vary in their *hydrophobicity* (see for example Kyte and Doolittle [7] or Engelman, Steitz and Goldman [4].) We therefore propose an extension of the HP model which accounts for the varying hydrophobicities [1].⁴

In the new model, we allow each of the 20 amino acids to have a value from the set $\{0, 1, 2, \dots\}$. Zero represents polar residues and non-zero values represent proportional levels of hydrophobicity. We then consider the value of a contact to be the amount of hydrophobicity buried from polar residues and solvent. Thus, a contact between any residue and solvent shields no hydrophobicity and is given a value of 0. A contact between a hydrophobic residue (with value from $\{1, 2, \dots\}$) and

⁴Note that although surface area of the actual residue may contribute to the strength of hydrophobicity attributed to a residue, we do not model the actual size differences in the spatial layout of the protein. Furthermore, we do not address the fact that polar residues also have diversity when compared to one another as well as when compared to water. These extension not considered here would still further improve the modeling of the protein folding problem.

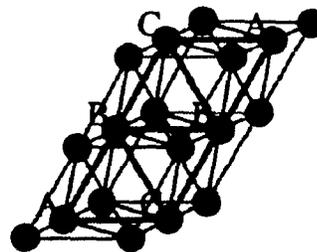


Figure 16: Contacts for positions A, B and C.

a polar residue (value 0) also shields no hydrophobicity and is given a value of 0. Finally, a contact between two hydrophobic residues (each with value from $\{1, 2, \dots\}$) shields both and is given value equal to their sum, the combined amount of hydrophobicity that is shielded.

If the ratio between the largest and smallest hydrophobic values assigned to the 20 amino acids is ρ , then blindly using an algorithm for the binary HP model could in the worst case result in a factor of ρ loss in the approximation factor. In the case of the algorithm of Section 3.4, as ρ grows arbitrarily the 16/30 approximation factor only falls by a constant to 2/5. These losses occur since some positions in a fold are involved in fewer contacts than others. If the strongly hydrophobic residues are placed in these positions, then the approximation suffers. Therefore, we must adapt our algorithms to adjust the construction based on the strength of the hydrophobic residues in the protein. Below, we adapt the algorithm of Section 3.4 for the generalized HP model and analyzed its performance. In addition, other folding algorithms such as the one in Section 3.6 should also be examined in this respect and one would hope that their performance could also be kept comparable to their performance in the binary HP model.

In adapting the algorithm of Section 3.4 we are able to achieve the same 16/30 asymptotic approximation that was possible in the binary HP model. The new algorithm will once again place all hydrophobic (non-zero) residues in the core using the layout shown in Figure 11. In this layout, the number of contacts that a hydrophobic residue participates in depends on which of three different types of positions within the core it occupies (See Figure 16). If we number the hydrophobic residues starting at 0, then the number of hydrophobic contacts (2 of which may be covalent) in which the r -th hydrophobic residue participates is

$$\begin{aligned} 6 &\Rightarrow \text{if } r \equiv 0 \pmod{3} && \text{(Position A)} \\ 9 &\Rightarrow \text{if } r \equiv 1 \pmod{3} && \text{(Position B)} \\ 7 &\Rightarrow \text{if } r \equiv 2 \pmod{3} && \text{(Position C)} \end{aligned}$$

out of a possible 12 contacts. Thus, there are three equivalence classes of residues, depending on their hydrophobic position in the primary sequence modulo 3.

For $i \in \{0, 1, 2\}$, let H_i be the set of residues at position r such that $r \equiv i \pmod{3}$, let W_i be the sum of the hydrophobic values of all of the residues in set H_i and let indices i_1, i_2, i_3 be such that $W_{i_1} \leq W_{i_2} \leq W_{i_3}$. If we define $W = W_{i_1} + W_{i_2} + W_{i_3}$, then the optimal score for a protein is at most $(12 - 2) * W$. Note that a single scan of the residues is sufficient to determine the values W_0, W_1 and W_2 .

By placing the first hydrophobic residue at a different starting position in the construction, we effectively replace r in the above equations with $r + 1$ or $r + 2$. Furthermore, by wrapping the core in the opposite direction, we effectively replace r by $-r$. Therefore, we adapt the construction to place the residues of H_{i_3} in position B, the residues of H_{i_2} in position C, and the residues of H_{i_1} in position A. If we let $x = W_{i_1}$, $y = W_{i_2} - W_{i_1}$ and $z = W_{i_3} - W_{i_2}$, then we achieve:

$$\begin{aligned}
 \text{Score} &= (9 - 2) * W_{i_3} + (7 - 2) * W_{i_2} + (6 - 2) * W_{i_1} \\
 &= 7 * (x + y + z) + 5 * (x + y) + 4 * (x) \\
 &= 16x + 12y + 7z \\
 &\geq 16x + \frac{32}{3}y + \frac{16}{3}z & (4.1) \\
 &= \frac{16}{3} * (3x + 2y + z) \\
 &= \frac{16}{3} * W
 \end{aligned}$$

The approximation factor is therefore at least 16/30. Inequality (4.1) is equality (yielding an approximation factor lower bound of 16/30) when $y = z = 0$, i.e. $W_0 = W_1 = W_2$. When the weights of the three classes differ, the approximation factor bound increases and asymptotically approaches 7/10 as all of the weight becomes concentrated in one of the three classes.

In some proteins, the weight of residues in one class (e.g. W_0) may be largest in the first half of the protein sequence, while the weight of residues in another class (e.g. W_2) may be largest in the second half of the protein sequence. In such cases, we would prefer one wrapping of the core in the first half and another wrapping in the second half. In order to change the wrapping in the middle of the sequence, there is a loss in the score due to the break in the core. Yet, in some cases this loss would be more than made up for by placing the heavier weighted residues at the B position in both halves of the sequence. Furthermore, this change in wrapping can be done as many times as is productive. Using dynamic programming, one could construct the best score possible by this strategy in time $O(nb)$ where n is the number of residues in the protein and b is the number of breaks permitted.

Note that these arguments for improving the approximation do not yield better *worst case* asymptotic

approximations since here the worst case is when all classes have equal weight. But, since real proteins undoubtedly have diverse hydrophobic makeup, the ability of an algorithm to leverage this diversity when present would be preferable. As this diversity is not guaranteed to exist in a worst-case analysis, it would also be of interest to determine if there are properties in biological protein sequences which HP folding algorithm *could* make use of to improve their performance guarantees, possibly using an average case analysis.

Acknowledgments

Thanks to Bonnie Berger, Gunnar Hoest, Jonathan King, Jon Kleinberg, S. Muthukrishnan and Lior Pachter for discussion on protein folding and lattices. Images generated with the help of RasMol Molecular Renderer, Version 2.5.1 (Roger Sayle, October 1994).

References

- [1] S. Decatur. Protein folding in the generalized hydrophobic-polar model on the triangular lattice. Technical Memo MIT-LCS-TM-559, Massachusetts Institute of Technology, May 1996.
- [2] K. A. Dill, Sarina Bromberg, Kaizi Yue, Klaus M. Fiebig, David B. Yee, Paul D. Thomas, and hue Sun Chan. Principles of protein folding: A perspective from simple exact models. *Prot. Sci*, 4:561-602, 1995.
- [3] K.A. Dill. Theory for the folding and stability of globular-proteins. *Biochemistry*, 24(6):1501-1509, 1985.
- [4] D. M. Engelman, T. A. Steitz, and A. Goldman. Identifying nonpolar transbilayer helices in amino acid sequences of membrane proteins. *Annu. Rev. Biophys. Chem.*, 1986.
- [5] W. Hart and S. Istrail. Fast protein folding in the hydrophobic-hydrophilic model within three-eighths of optimal. *Journal of Computational Biology*, 3(1):53-96, 1996.
- [6] W. Hart and S. Istrail. Invariant patterns in crystal lattices: Implications for protein folding algorithm. In *Combinatorial Pattern Matching, Proc. 7th Annual Symposium (CPM'96)*, 1996. 288-303.
- [7] J. Kyte and R. F. Doolittle. A simple method for displaying the hydropathic character of protein. *J. Mol. Biol.*, 1982.
- [8] R. Lathrop. The protein folding threading problem with sequence amino acid interaction preferences is NP-complete. *Protein Engineering*, 7(9):1059-1068, 1994.
- [9] J. T. Ngo and J. Marks. Computational complexity of a problem in molecular structure prediction. *Protein Engineering*, 5(4):313-321, 1992.
- [10] M. Paterson and T. Przytycka. On the complexity of string folding. In *ICALP'96*, volume 1099 of *LNCS*, pages 658-669, 1996.
- [11] R. Unger and J. Moulton. Finding the lowest free energy conformation of a protein is a NP-hard problem: Proof and implications. *Bulletin of Mathematical Biology*, 55(6):1183-1198, 1993.