

# Invariant Patterns in Crystal Lattices: Implications for Protein Folding Algorithms (Extended Abstract)

William E. Hart\* and Sorin Istrail\*\*

Sandia National Laboratories  
Massively Parallel Computing Research Laboratory  
P. O. Box 5800  
Albuquerque, NM 87185-1110

**Abstract.** Crystal lattices are infinite periodic graphs that occur naturally in a variety of geometries and which are of fundamental importance in polymer science. Discrete models of protein folding use crystal lattices to define the space of protein conformations. Because various crystal lattices provide discretizations of the same physical phenomenon, it is reasonable to expect that there will exist “invariants” across lattices that define fundamental properties of the protein folding process; an invariant defines a property that transcends particular lattice formulations. This paper identifies two classes of invariants, defined in terms of sublattices that are related to the design of algorithms for the structure prediction problem. The first class of invariants is used to define a master approximation algorithm for which provable performance guarantees exist. This algorithm can be applied to generalizations of the hydrophobic-hydrophilic model that have lattices other than the cubic lattice, including most of the crystal lattices commonly used in protein folding lattice models. The second class of invariants applies to a related lattice model. Using these invariants, we show that for this model the structure prediction problem is intractable across a variety of three-dimensional lattices. It turns out that these two classes of invariants are respectively sublattices of the two- and three-dimensional square lattice. As the square lattices are the standard lattices used in empirical protein folding studies, our results provide a rigorous confirmation of the ability of these lattices to provide insight into biological phenomenon. Our results are the first in the literature that identify algorithmic paradigms for the protein structure prediction problem that transcend particular lattice formulations.

## 1 Introduction

Crystal lattice models are vehicles for reasoning about the protein folding phenomenon through analogy. Crystal lattices are infinite periodic graphs that are

---

\* wehart@cs.sandia.gov, <http://www.cs.sandia.gov/~wehart/main.html>

\*\* scistra@cs.sandia.gov, <http://www.cs.sandia.gov/~scistra/main.html>

generated by translations of a “unit cell” that fill a two or three-dimensional space. In polymer science many important results have been obtained through the use of lattice models [4, 9]. In the context of protein folding, lattices provide a natural discretization of the space of protein conformations. The sequence of amino acids that defines a protein can be viewed as a path labeled with amino acids on vertices. A conformation of a protein is a self-avoiding embedding of this path into a lattice, where each vertex of the path is mapped to a vertex of the lattice and edges of the path are mapped to edges of the lattice. With every conformation we can associate an energy value using rules defined by the model, which take into account the neighborhood relationship of the amino acids. The central focus of this paper is the design of algorithms that construct a conformation of minimal or near-minimal energy for a given sequence.

Of particular interest here is the design of algorithms that can be applied to a variety of lattice models. Results that transcend particular lattice frameworks are of significant interest because they can say something about the general biological problem with a higher degree of confidence. In fact, it is reasonable to expect that there will exist invariants across lattices that fundamentally relate to the protein folding problem, because lattice models provide discretizations of the same physical phenomenon. However, the identification of such invariants has not been previously addressed.

This paper identifies invariants across lattice models that can be described in terms of sublattices. These invariants give the ability to address the following question. Given an energy formula for crystal lattices, does there exist an algorithm that takes a sequence and a lattice and produces a conformation with minimal energy? If such an algorithm exists, it may provide valuable insight into the protein folding process because it captures essential features of protein folding.

We address this question in two ways. First we design performance guaranteed approximation algorithms for protein folding in the hydrophobic-hydrophilic model. This model categorizes amino acids as hydrophobic (nonpolar) or hydrophilic (polar), and the energy of a conformation is equal to the number of hydrophobic-hydrophobic contacts. The invariant we use to design a “master” approximation algorithm employs special sublattices which we call laticoids. Laticoids impose a structure in which a skeleton of hydrophobic contacts can be constructed, thereby leading to folding algorithms whose performance can be analyzed. In the particular case of the square two-dimensional lattice, the laticoid describes the structure used in the approximation algorithms described by Hart and Istrail [8].

We prove that our master approximation algorithm has performance guarantees for a class of lattices that includes most of the lattices commonly used in simple exact protein folding models, e.g. two- and three-dimensional square lattice [4, 7, 12], the diamond (carbon) lattice [13], the face-centered-cubic lattice [2] and the 210 lattice used by Skolnick [14]. Furthermore, this class encompasses a large number of other lattices studied in crystallography. Our main theorems state that laticoids of the two-dimensional square lattice can be embedded

into all of these lattices, and therefore, every lattice in the class is approximable in linear time.

Second, we prove that lattice models related to those considered by Unger and Moulton [15] are NP-complete. The lattice model considered by Unger and Moulton uses a distance-related energy formula between an unbounded number of amino acid types. Our results extend their NP-completeness argument to any three-dimensional lattice into which a certain type of sublattice can be embedded. All of the three-dimensional lattices mentioned above fall into this class.

## 2 Lattice Models for Protein Folding

Lattice models for protein folding can be distinguished by at least five properties:

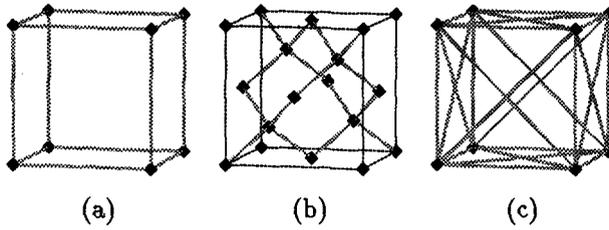
1. An alphabet of types of amino acids that the model considers;
2. The set of protein instances represented as sequences from this alphabet;
3. An energy formula specifying how the conformational energy is computed;
4. Parameters for the energy formula;
5. A crystal lattice that provides a discretization of the conformation space.

For example, the hydrophobic-hydrophilic (HP) model [3] can be described as follows. The alphabet used in an HP model is  $A = \{0, 1\}$ , and the set of protein instances is the set of binary sequences  $\sigma = \{0, 1\}^+$ . Each sequence  $s \in \sigma$  is the (hypothesized) hydrophobic-hydrophilic pattern of a protein sequence, where 1 represents a hydrophobic amino acid, and 0 represents a hydrophilic amino acid. We will refer to  $s$  as a protein instance. Contact energies are used in this model, so the energy formula is an energy matrix,  $\mathcal{E}$ . The energy matrix is indexed by the alphabet symbols,  $\mathcal{E} = (e(a, b))_{a, b \in A}$ . For HP models,  $e(a, b) = -1$  if  $a = b = 1$ , and  $e(a, b) = 0$  otherwise.

We consider protein folding models on a large class of crystal lattices, including the square lattice. Crystal lattices are infinite periodic graphs that are generated by translations of a "unit cell" that fill a two- or three-dimensional space. A unit cell contains a finite graph that is connected to neighboring unit cells. Examples of crystal lattices are shown in Figure 1.

One can interpret a protein sequence  $s = s_1 \dots s_m$  as an  $m$ -vertex node-labeled path, where for  $1 \leq i \leq m$ , node  $i$  is labeled with  $s_i$ . The path has  $m - 1$  edges that are called *bonds*. A *conformation*  $C$  of a protein sequence  $s$  in a lattice  $L$  is a path in the lattice in which the protein sequence is embedded, i.e., the protein vertices are mapped one-to-one to lattice points, and protein bonds are mapped to the corresponding lattice edges. The *energy* of a conformation of the protein sequence  $s$  in  $L$  can be computed using distances in the lattice. For example, in the HP model the energy is a function of the number of "contact edges." A contact edge is a lattice edge that is not a protein bond (in the embedding) but has both endpoints labeled. In HP models, contact edges with 1s at their endpoints have weight  $-1$  while all other contact edges have weight 0.

The *native* conformation of a protein is the conformation that has biological function. According to the Thermodynamic Hypothesis the native conformation



**Fig. 1.** Examples of crystal lattices: (a) cubic, (b) diamond, and (c) cubic with planar diagonals.

of a protein is the conformation with the minimum energy among the set of all conformations. Consequently, given a sequence  $s$  and a lattice model, the protein folding structure prediction problem is to find a native conformation of  $s$  in  $L$ . It is unknown whether this problem is NP-complete for HP models, but a few related models have been shown to be NP-complete [5, 10, 11, 15]. Furthermore, Hart and Istrail [8] have demonstrated that performance guaranteed approximation algorithms exist for HP models on square and cubic lattices.

Let  $\mathcal{Z}_L(s)$  be the energy of the conformation generated for protein instance  $s$  on lattice  $L$  with by algorithm  $\mathcal{Z}_L$ , and let  $OPT_L(s)$  be the energy of the optimal conformation of  $s$  on  $L$ . A standard performance guarantee used for approximation algorithms is the asymptotic performance ratio  $R^\infty(\mathcal{Z}_L)$  [6]. If  $R^\infty(\mathcal{Z}_L) = \tau$ , then as  $\mathcal{Z}_L$  is applied to larger protein instances, the value of solutions generated by  $\mathcal{Z}_L$  approaches a factor of  $\tau$  of the optimum. Here, “large” protein instances have low conformational energy at their native state, which may be independent of their length. Since  $\mathcal{Z}_L(s) \leq 0$  and  $OPT_L(s) \leq 0$ , both of these ratios are scaled between 0 and 1 such that a ratio closer to 1 indicates better performance.

### 3 Protein Sequence Structure in the HP Model

The protein folding models that we first analyze are HP models. HP models abstract the hydrophobic interaction process in protein folding by reducing a protein to a heteropolymer that represents a predetermined pattern of hydrophobicity in the protein. This is one of the most studied lattice models for protein folding, and despite its simplicity, the model is powerful enough to capture a variety of properties of actual proteins [4].

This section summarizes key definitions concerning the structure of protein instances from Hart and Istrail [8]. Let  $s = s_1, \dots, s_m$  be a protein instance,  $s_i \in \{0, 1\}$ . Let  $l(s)$  equal the length of the sequence  $s$ . Let  $M_{max}(s)$  equal the length of the longest sequence of zeros in  $s$ , and let  $M_{min}(s)$  equal the length of the shortest sequence of zeros in  $s$ . Finally, let  $E(s)$  equal the number of adjacent elements in the sequence,  $s_j$  and  $s_{j+1}$  for which  $s_j = 1$  and  $s_{j+1} = 1$ .

An instance  $s$  can be decomposed into a sequence of *blocks*. A block  $b_i$  has the form  $b_i = 1$  or  $b_i = 1Z_{i_1}1 \dots Z_{i_h}1$ , where the  $Z_{i_j}$  are odd-length sequences of 0's and  $h \geq 1$ . A *block separator*  $z_i$  is a sequence of 0's that separates two consecutive blocks, where  $l(z_i) \geq 0$  and  $l(z_i)$  is even for  $i = 1, \dots, h-1$ . Thus  $s$  is decomposed into  $z_0b_1z_1 \dots b_hz_h$ . Since  $l(z_i) \geq 0$ , this decomposition treats consecutive 1's as a sequence of blocks separated by zero-length block separators. Let  $N(b_i)$  equal the number 1's in  $b_i$ . Thus the sequence

$$0 \underbrace{10101}_{b_1} \underbrace{1}_{b_2} \underbrace{1}_{b_3} \underbrace{10101}_{b_4} 0000 \underbrace{1010101}_{b_5}$$

can be represented as  $l(z) = (1, 0, 0, 0, 4, 0)$  and  $N(b) = (3, 1, 1, 3, 4)$ .

It is useful to divide blocks into two categories:  $x$ -blocks and  $y$ -blocks. For example, let  $x_i = b_{2i}$  and let  $y_i = b_{2i-1}$ . Let  $B_x$  and  $B_y$  be the number of  $x$ -blocks and  $y$ -blocks respectively. Further, let  $X = X(s) = \sum_{i=1}^{B_x} N(x_i)$  and  $Y = Y(s) = \sum_{i=1}^{B_y} N(y_i)$ . Let  $T_x(s)$  equal the number of endpoints of  $s$  that are 1's in  $x$ -blocks, and let  $T_y(s)$  equal the number of endpoints of  $s$  that are 1's in  $y$ -blocks. We assume that the division into  $x$ - and  $y$ -blocks is such that  $X \leq Y$  and if  $X = Y$  then  $T_x(s) \geq T_y(s)$ . For example, the sequence

$$0 \underbrace{10101}_{y_0} \underbrace{1}_{x_0} \underbrace{1}_{y_1} \underbrace{10101}_{x_1} 0000 \underbrace{1010101}_{y_2}$$

can be represented as  $z_0y_0z_1x_0z_2y_1z_3x_1z_4y_2z_5$ , where  $l(z) = (1, 0, 0, 0, 4, 0)$ ,  $N(x) = (1, 3)$ , and  $N(y) = (3, 1, 4)$ .

A *superblock*  $B_i$  is comprised of sequences of blocks as follows:  $B_i = b_{i_1}z_{i_1} \dots z_{i_{h-1}}b_{i_h}$ . Let  $N_x(B_i)$  equal the sum of  $N(b_j)$ , where  $b_j$  are  $x$ -blocks in  $B_i$ . Let  $N_y(B_i)$  equal the sum of  $N(b_j)$ , where  $b_j$  are  $y$ -blocks in  $B_i$ . Finally, let  $N(B_i) = N_x(B_i) + N_y(B_i)$ .

## 4 Master Approximation Algorithms for the HP Model

We now describe two paradigms for master approximation algorithms for the HP model. These master approximation algorithms are distinguished by properties of the lattices to which they apply. The first paradigm captures two aspects of the protein folding algorithms described by Hart and Istrail [8]: (1) the selection of a folding point that balances hydrophobicity and (2) the skeleton of contact edges that forms the hydrophobic core. We call this the *bipartite master approximation algorithm* because it is applicable to crystal lattices that can be described as a bipartite graph. These crystal lattices have the property that two 1's can be endpoints of a contact edge only if there is an even number of elements between them [8]. The second paradigm describes the *nonbipartite master approximation algorithm*, which is applicable to lattices that cannot be described as a bipartite graph. These graphs have the property that they contain odd cycles.

#### 4.1 The Bipartite Master Approximation Algorithm

Consider the following definitions.

**Definition 1** Given a path  $p$  in a lattice  $L$  from  $a$  to  $b$ , let  $d_p(a, b)$  be the length of  $p$ . A path  $p$  from  $a$  to  $b$  is polynomial evenly extensible if there exist paths  $p_k$  for every  $k \in \mathbf{Z}^{>0}$  such that  $d_{p_k}(a, b) = d_p(a, b) + 2k$  and there exists a polynomial time algorithm that given  $p$  and  $k$  constructs  $p_k$ . The collection of the paths of an polynomial evenly extensible path  $p$  is called the extension of  $p$  in  $L$ .

**Definition 2** Given polynomial evenly extensible paths  $p$  from  $a$  to  $b$  and  $q$  from  $c$  to  $d$ , we say that  $p$  and  $q$  are extensibly disjoint if their extensions are vertex disjoint.

**Definition 3** A bipartite latticoid,  $\hat{L}$ , of  $L$  is an infinite graph that contains an infinite sequence of contact edges  $(a_i, b_i)$  with the following properties:

- There is an polynomial evenly extensible path  $p_i^a$  from  $a_i$  to  $a_{i+1}$  and polynomial evenly extensible path  $p_i^b$  from  $b_i$  to  $b_{i+1}$ ,
- There is a constant  $\kappa > 0$  such that for every  $i$  and  $j$ ,  $d_{p_i^a}(a_i, a_{i+1}) = d_{p_j^b}(b_j, b_{j+1}) = 2\kappa$ , and
- The set of paths  $\{p_i^a, p_i^b \mid i = 1, \dots\}$  are mutually extensibly disjoint.

The dilation of the bipartite latticoid is  $\Delta_{\hat{L}} = \kappa$ .

Figure 2 illustrates the structure of a bipartite latticoid. Because the paths  $A_i$  are evenly extensible, the paths  $B_i$  and  $C_i$  can be constructed in polynomial time. Furthermore, the vertices in  $\{A_i, B_i, C_i\}$  and  $\{A_j, B_j, C_j\}$  do not intersect.

The bipartite master approximation algorithm takes a bipartite latticoid  $\hat{L}$  and selects a single folding point (turning point) that divides a protein instance into a  $y$ -superblock  $B'$  and an  $x$ -superblock  $B''$ . The folding point is selected using "Subroutine 1" from Hart and Istrail [8]. Subroutine 1 selects a folding point that balances the hydrophobicity between the  $x$ -blocks and  $y$ -blocks on each half of the folding point. The following lemma describes the key property of the folding point that is selected.

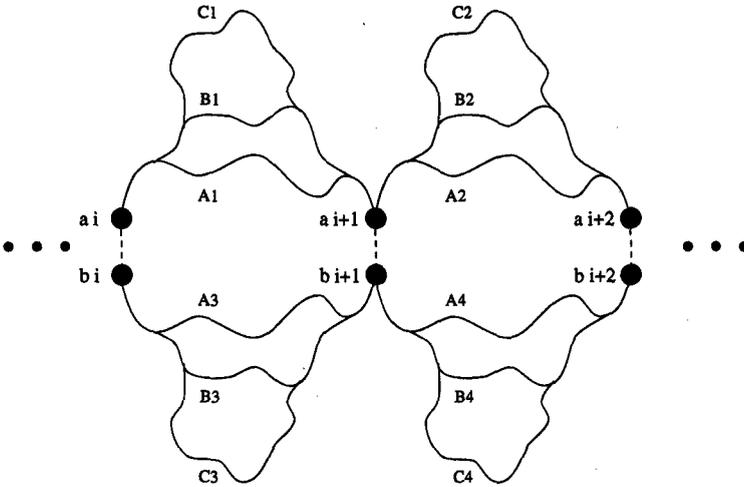
**Lemma 1** ([8], Lemma 1). *The folding point selected by Subroutine 1 partitions a protein instance  $s$  into two superblocks  $B'$  and  $B''$  such that either*

$$N_y(B') \geq \lceil (Y+1)/2 \rceil \quad \text{and} \quad N_x(B'') \geq \lceil X/2 \rceil$$

or

$$N_y(B') \geq \lceil Y/2 \rceil \quad \text{and} \quad N_x(B'') \geq \lceil (X+1)/2 \rceil .$$

After selecting the folding point, the conformation of the two superblocks is dictated by the bipartite latticoid  $\hat{L}$ . The bipartite latticoid specifies the placement of the contact edges between the superblocks, as well as the conformation of the loops within each superblock. This generalizes the notion of "normal form" that was used to describe the approximation algorithms in Hart and Istrail [8].



**Fig. 2.** A symbolic illustration of the structure of bipartite latticoids.

Decomposition into  $x$ - and  $y$ -blocks requires a single pass through the protein instance. Subroutine 1 requires a single pass through the sequence of blocks, which is no longer than the length of the protein instance. The construction of the final conformation requires polynomial time to create the paths for the zero-loops. Thus the computation required by Algorithm  $\mathcal{A}_{\hat{L}}$  is polynomial.

Let  $\mathcal{A}_{\hat{L}}(s)$  represent the energy of the final conformation generated by Algorithm  $\mathcal{A}_{\hat{L}}$ . The performance of Algorithm  $\mathcal{A}_{\hat{L}}$  can be bounded as follows.

**Lemma 2.**

$$\mathcal{A}_{\hat{L}}(s) \leq - \left\lceil \frac{X}{2\Delta_{\hat{L}}} \right\rceil + 1.$$

Let  $\delta(L)$  be the maximum degree of all vertices in  $L$ . Since  $L$  is a crystal lattice generated by a unit cell,  $\delta(L)$  is finite. It follows from the fact that  $L$  is bipartite that  $OPT_L(s) \leq -(\delta(L) - 2)X(s) - 2$ . Proposition 1 presents the asymptotic ratio for Algorithm  $\mathcal{A}_{\hat{L}}$ .

**Proposition 1**  $R^\infty(\mathcal{A}_{\hat{L}}) \geq 1/(2\Delta_{\hat{L}}(\delta(L) - 2))$ .

*Proof.* We know from Lemma 2 that

$$\mathcal{A}_{\hat{L}}(s) \leq - \left\lceil \frac{X(s)}{2\Delta_{\hat{L}}} \right\rceil + 1.$$

Now  $OPT_L(s) \leq -(\delta(L) - 2)X(s) - 2$ , so

$$R_{\mathcal{A}_L}(s) = \frac{\mathcal{A}_L(s)}{OPT_L(s)} \geq \frac{-\left\lceil \frac{X(s)}{2\Delta_L} \right\rceil + 1}{-(\delta(L) - 2)X(s) - 2} \tag{1}$$

$$\geq \frac{-\frac{X(s)}{2\Delta_L} + 1}{-(\delta(L) - 2)X(s) - 2} = \frac{X(s) - 2\Delta_L}{2\Delta_L(\delta(L) - 2)X(s) + 4\Delta_L} \tag{2}$$

For  $s \in S_N^L$ ,  $-(\delta(L) - 2)X(s) - 2 \leq N$ , so  $X(s) \geq -(N + 2)/(\delta(L) - 2)$ . Since Equation (1) is monotonically increasing for  $X(s) \geq 0$ , we have

$$R_{\mathcal{A}_L}(s) \geq \frac{-(N + 2)/(\delta(L) - 2) - 2\Delta_L}{-2\Delta_L(N + 2) + 4\Delta_L} = \frac{N + 2 - 4\Delta_L + 2\Delta_L\delta(L)}{2\Delta_L(\delta(L) - 2)N}$$

so

$$R^N(\mathcal{A}_L) \geq \frac{N + 2 - 4\Delta_L + 2\Delta_L\delta(L)}{2\Delta_L(\delta(L) - 2)N}$$

and

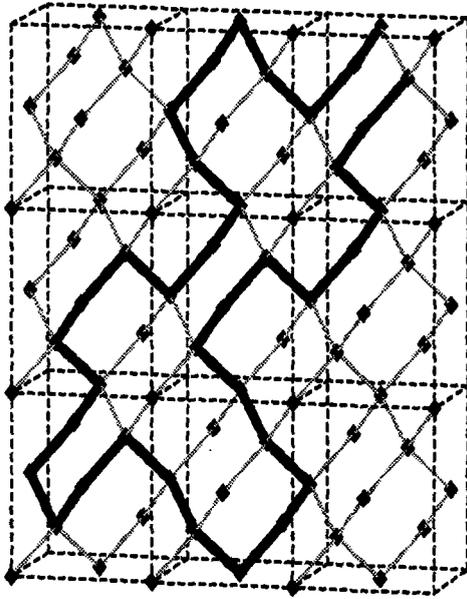
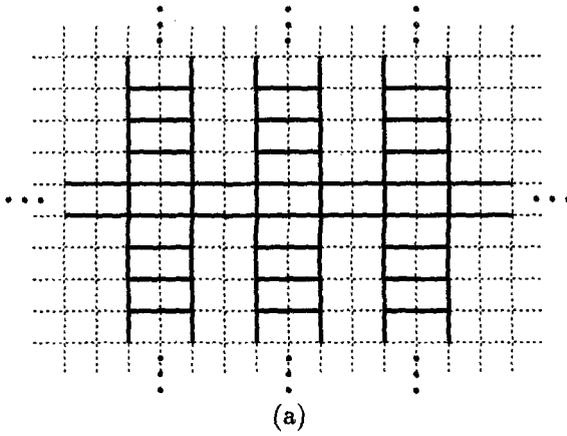
$$\begin{aligned} R^\infty(\mathcal{A}_L) = \sup\{r \mid R^N(\mathcal{A}_L) \geq r, N \in \mathbf{Z}\} &\geq \lim_{N \rightarrow \infty} \frac{N + 2 - 4\Delta_L + 2\Delta_L\delta(L)}{2\Delta_L(\delta(L) - 2)N} \\ &= 1/(2\Delta_L(\delta(L) - 2)). \end{aligned}$$

■

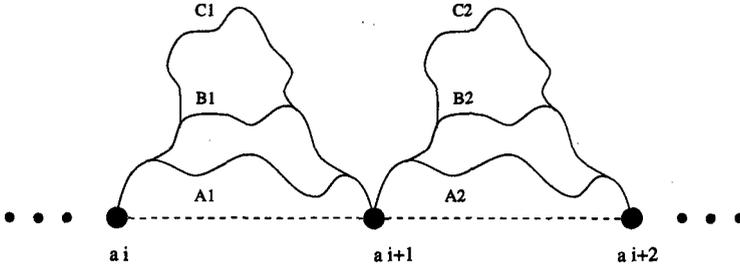
To illustrate the application of the bipartite master approximation algorithm, consider its application to the diamond lattice, which has previously been used in lattice models for protein folding [13]. Figure 3 shows the embedding of a “diluted” square lattice into a plane of unit cells for the diamond lattice. Figure 3a illustrates a bipartite latticoid of  $\hat{\mathbf{L}}_0$  that can be embedded into the diamond lattice. Figure 3b demonstrates this embedding. Dashed and solid lines between vertices in each unit cell indicate the edges of the diamond lattice that are used to embed a square lattice for which one dimension is dilated to length two. Edges not used for this embedding are omitted. The solid lines illustrate a conformation of a protein on this lattice that the bipartite master approximation algorithm would generate. Now  $\delta(L) = 4$  for the diamond lattice  $L$ . It follows from Proposition 1 that  $R^\infty(\mathcal{A}_L) = 1/8$ .

### 4.2 The Nonbipartite Master Approximation Algorithm

This section briefly sketches the details of the nonbipartite master approximation algorithm (full details will be given in the full paper). Figure 4 illustrates the structure of a *nonbipartite latticoid*. The hydrophobic residues in the protein sequence are placed along a path of  $a_i$  that are in contact. The paths  $A_i$  are *extensible*, which implies that in polynomial time they can be extended to any even or odd length beyond some minimal length. Furthermore, the paths  $\{A_i, B_i, C_i\}$



**Fig. 3.** Illustration of the embedding of a bipartite latticoid from  $\hat{L}_0$  into a diamond lattice: (a) the bipartite latticoid, and (b) the embedding into the diamond lattice.



**Fig. 4.** A symbolic illustration of the structure of nonbipartite latticoids.

and  $\{A_j, B_j, C_j\}$  are extensively disjoint. Note that because the hydrophobic-hydrophobic contacts are constructed along a path, the extensible paths may lie on either side of this path.

For a nonbipartite latticoid  $\hat{L}$ , the *dilation*  $\Delta_{\hat{L}}$  is half of the minimal length of a path from  $a_i$  to  $a_{i+1}$ . Given this, we can prove the following performance guarantee for a nonbipartite master approximation algorithm  $\mathcal{B}$  on lattice  $L$  with latticoid  $\hat{L}$ .

**Proposition 2**  $R^\infty(\mathcal{B}_{\hat{L}}) \geq 1/(2\Delta_{\hat{L}}(\delta(L) - 2))$ .

## 5 A Complexity Theory for Protein Folding on Bipartite Crystal Lattices

In this section we describe a framework for analyzing the design of efficient approximation algorithms with provable performance guarantees on bipartite lattices. The unifying theme is polynomial approximability asymptotic within a constant of optimal. This theory defines polynomial embedding reductions from one bipartite lattice to another, and relates the approximability on the first lattice to the approximability on the second. Further, this theory includes a notion of *completeness*, which defines the “hardest” members in the class. While we restrict our discussions to bipartite lattices, these notions naturally generalize to nonbipartite lattices.

*Definitions* A lattice  $L$  is *polynomial kernel-approximable* if there is a polynomial algorithm  $\mathcal{A}$  and constants  $\alpha_L, \beta_L \in \mathbf{Z}^{>0}$  such that for all protein instances  $s$ ,  $A(s) = -\alpha_L X(s) + \beta_L$ . A class of lattices  $\mathcal{L}$  is *polynomial kernel-approximable* if for every  $L \in \mathcal{L}$ ,  $L$  is polynomial kernel-approximable. Let **PKAL** be the class of polynomial kernel-approximable lattices. A lattice  $L$  is *polynomial approximable* if there is a polynomial algorithm  $\mathcal{A}$  and a constant  $\tau_L \in \mathbf{R}^{>0}$  such  $R^\infty(\mathcal{A}) \geq \tau_L$ . A class of lattices  $\mathcal{L}$  is *polynomial approximable* if for every  $L \in \mathcal{L}$ ,  $L$  is polynomial approximable. Let **PAL** be the class of polynomial approximable lattices. A

sublattice  $\hat{L}$  of  $L$  is a subgraph of  $L$  that is obtained by removing edges and vertices from  $L$ . A particular sublattice is the latticoid.

While we aspire to a framework for general approximability for all lattices, our current framework applies to kernel-approximability on bipartite lattices.

**Lemma 3.** *If  $L$  is polynomial kernel-approximable, then there exists a polynomial algorithm  $A$  and constant  $C_L$  such that  $R^\infty(A) \geq C_L$ .*

**Corollary 1** *If  $\hat{L}$  is a sublattice of a lattice  $L$  and  $\hat{L}$  is polynomial kernel-approximable, then  $L$  is polynomial kernel-approximable.*

**Definition 4** *A core of a lattice  $L$  is a set of sublattices  $D(L) = \{\hat{L}^1, \hat{L}^2, \dots\}$ , where  $D(L)$  is finite or countably infinite.*

Folding algorithms in a lattice  $L_1$  can be transferred to folding algorithms in another lattice  $L_2$ , a folding “reduction”, if the sublattice used in  $L_1$  by the approximation algorithm can be embedded in  $L_2$ . This reduction can be polynomial in the sense that each unit cell is given by a finite description, and the symmetries in the crystal lattice are with respect to the neighboring cells (and thus also of finite description). This notion of reduction is formalized in the following definition.

**Definition 5** *A polynomial embedding reduction of  $L_1$  to  $L_2$  via core  $D(L_1)$  is a polynomial time function  $\psi : \hat{L}_1 \rightarrow \hat{L}_2$  such that: (1)  $\hat{L}_1$  is a sublattice in  $D(L_1)$ , (2)  $\hat{L}_2$  is a sublattice of  $L_2$ , and (3)  $\psi(\hat{L}_1)$  is lattice isomorphic to  $\hat{L}_2$  (i.e. graph isomorphic). We say that  $\hat{L}_1$  is embedded into  $\hat{L}_2$ . If there is a polynomial embedding reduction from  $L_1$  to  $L_2$  via core  $D(L_1)$ , we write  $L_1 \propto_{D(L_1)} L_2$ .*

**Definition 6** *A lattice  $L$  with core  $D(L)$  is polynomial core kernel-approximable if  $D(L) \subseteq \text{PKAL}$ .*

**Lemma 4.** *If a lattice  $L_1$  with core  $D(L_1)$  is polynomial core kernel-approximable and  $L_1 \propto_{D(L_1)} L_2$ , then  $L_2$  is polynomial kernel-approximable.*

The central concept of this theory is the notion of completeness defined as follows.

**Definition 7** *Let  $\mathcal{L}$  be a class of lattices. A lattice  $L$  is called  $\mathcal{L}$ -complete via core  $D(L)$  if (1)  $L \in \mathcal{L}$  and (2)  $\forall L' \in \mathcal{L}$ ,  $L \propto_{D(L)} L'$ .*

Similar to the theory of NP-completeness, if any member of the complete set is core-approximable then we can design polynomial approximation algorithms for all lattices in the class.

**Theorem 1** *Let  $L$  be a lattice with core  $D(L)$ . If  $L$  is  $\mathcal{L}$ -complete and polynomial core kernel-approximable then  $\mathcal{L} \subseteq \text{PKAL} \subset \text{PAL}$ .*

## 6 Approximable Lattices for the HP Model

In this section we describe a class of lattices  $\mathcal{L}$  for which performance guaranteed approximation algorithms exist.  $\mathcal{L}$  is a broad class of lattices that includes many of the lattices previously used in lattice models for protein folding. Further, it includes many other important crystallographic lattices. This result confirms that performance guaranteed approximability is not an artifact of the square and cubic lattices. Further, this lattice independence results suggests that the algorithmic mechanisms used to generate these approximate conformations may play a role in biological systems.

Our description of  $\mathcal{L}$  is split into the following sets of lattices:

1. Bravais lattices, which contain all points  $R$  of the form  $R = n_1a_1 + n_2a_2 + n_3a_3$ , where  $n_i$  are integers and  $a_i$  are linearly independent vectors in  $\mathbf{R}^n$  [1].
2. The planar triangular lattice, which tiles the plane with equilateral triangles, and the hexagonal close packed crystal structure.
3. The diamond lattice and the fluorite structure.
4. The hexagonal lattice, and bipartite lattices into which the hexagonal lattice can be embedded. This is significant since there are a large number of crystal lattices for which the hexagonal lattice can be embedded. The catalog of lattices in Wells [16] contains many three-dimensional lattices into which the hexagonal lattice can be embedded.
5. The “210 lattice” that Skolnick and Kolinkski [14] use to place  $\alpha$ -carbons. In this lattice, the  $\alpha$ -carbons are connected by the 3D generalization of the “knight’s walk” in chess.

The proof that these lattices are approximable uses the complexity theory outlined in the previous section. Although  $\mathcal{L} \subseteq \mathbf{PKAL}$ , it is unclear whether this relation is strict.  $\mathcal{L}$  certainly spans a broad class of crystal lattices. Furthermore, we believe that it contains many biologically relevant crystal lattices. For example, it contains most of the lattices previously used in protein folding lattice models [2, 4, 7, 12, 13, 14].

## 7 Hardness Results

In this section, we generalize the NP-hardness proof by Unger and Moutl [15] to show that it is applicable for a variety of lattices. Let  $L$  be a three-dimensional crystal lattice and let  $\mathbf{Z}$  be the set of integers. Suppose that  $S$  is a protein instance represented by a sequence of amino acids  $s_1, \dots, s_n$ . For a conformation of  $S$ , suppose the coordinate of  $s_i$  is  $(x_i, y_i, z_i)$ . Then  $d_{ij}^x = |x_i - x_j|$ ,  $d_{ij}^y = |y_i - y_j|$ , and  $d_{ij}^z = |z_i - z_j|$ . We can define a lattice-specific protein folding problem as follows.

### ***L*-PF**

**Instance:** A sequence  $S = (s_1, \dots, s_n)$ ,  $s_i \in A \subset \mathbf{Z}$ ; a positive function

$g : [0, n]^3 \rightarrow \mathbf{R}^+$ ; a matrix  $C \in \mathbf{Z}^{m \times m}$ ,  $m = |A|$ ;  $B \in \mathbf{Z}$ .

**Question:** Is there an embedding of  $S$  in  $L$  such that

$$\sum_{i=1}^n \sum_{j \neq i} C_{s_i, s_j} g(d_{ij}^x, d_{ij}^y, d_{ij}^z) \leq B?$$

Unger and Moulton [15] demonstrate that  $L$ -PF is NP-complete for the lattice  $L$  defined by the unit cell in Figure 1c. The NP-completeness of  $L$ -PF problems can be generalized to a variety of other lattices by noting a key property of the conformations used to construct their proof. The reduction from OLA used by Unger and Moulton requires that certain residues be placed along a line parallel to the  $x$ -axis in the optimal conformation. Further, it must be possible to construct vertex-independent paths between these residues for any permutation of their ordering along this line.

A second class of invariant patterns in lattices occurs in the context of this type of NP-completeness argument. We can abstract the type of structure needed for the reduction as a sublattice. Using ideas similar to the previous invariants, we can then construct NP-completeness reductions for a variety of crystal lattices. Figures 5 and 6 illustrate the concept of this class of invariants on two lattices: the cubic and diamond lattice. The numbers in these figures indicate the amino acids that are placed collinear parallel to the  $x$ -axis.

Our analysis uses a reduction from the Optimal Linear Arrangement Problem (OLA) [6]:

### OLA

**Instance:** A graph  $G = (V, E)$ ; a positive integer  $B$ .

**Question:** Is there a one-to-one function  $f : V \rightarrow \{1, 2, \dots, |V|\}$  such that

$$\sum_{\{u, v\} \in E} |f(u) - f(v)| \leq B?$$

**Theorem 2** *Let  $L$  be a Bravais, diamond, fluorite or hexagonal close packed lattice. Then  $L$ -PF is NP-complete.*

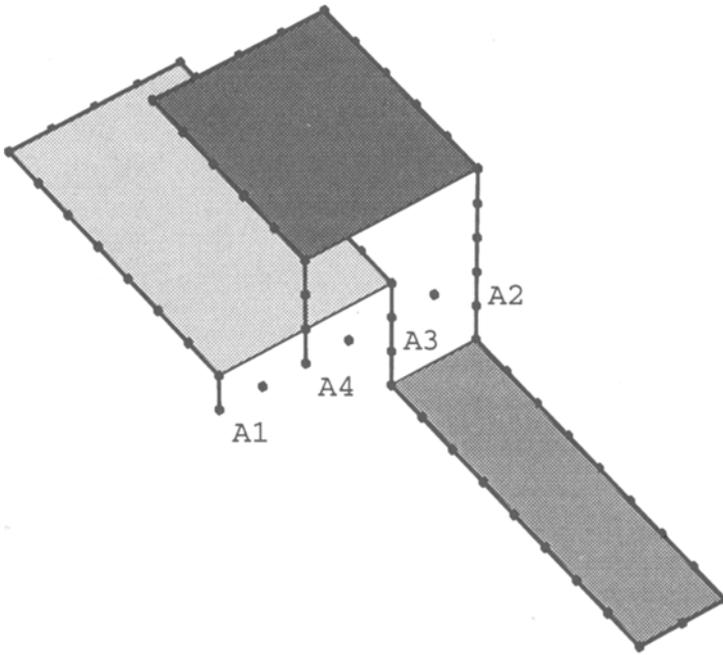
*Proof.* We show that if  $\mathcal{L}$  is the cubic lattice then  $\mathcal{L}$ -PF is NP-complete. The proof follows similarly for the other crystal lattices.

To transform an instance of OLA to  $\mathcal{L}$ -PF, we construct a protein instance as follows. Let  $A = V \cup \{x\}$  be a set of amino acids  $a_i$  that correspond to the vertices in  $V$  as well as a "dummy" amino acid  $x$ . Let  $\bar{f}(a_i) = f(v_i)$ , for  $a_i \in A$  and  $v_i \in V$ . Consider

$$S = a_1 \underbrace{xxx \dots xx}_{4n+3} a_2 \underbrace{xxx \dots xx}_{4n+3} \dots \underbrace{xxx \dots xx}_{4n+3} a_n.$$

The costs are

$$C_{s_i, s_j} = \begin{cases} |\bar{f}(s_i) - \bar{f}(s_j)| & \text{if } s_i, s_j \in A \\ 0 & \text{otherwise} \end{cases},$$



**Fig. 5.** The conformational invariant needed for the cubic lattice.

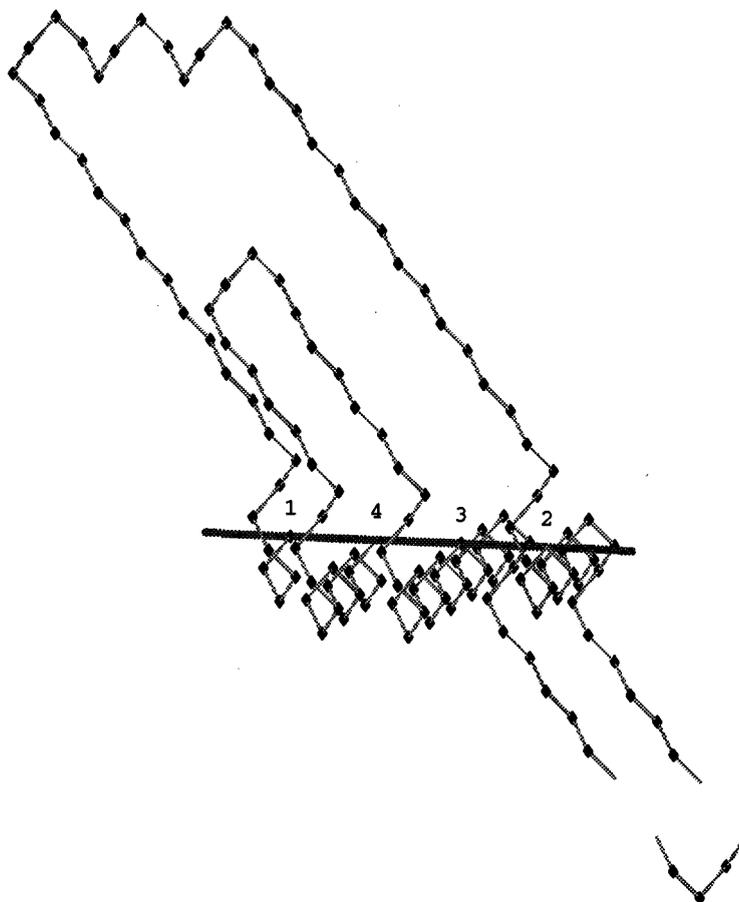
We use the same parameter  $B$  to bound the energy as in the OLA instance. The distance function  $g$  is given by

$$g(d_{ij}^x, d_{ij}^y, d_{ij}^z) = \begin{cases} d_{ij}^x/2 & \text{if } d_{ij}^y, d_{ij}^z = 0 \text{ and} \\ & d_{ij}^x \text{ is even} \\ (B+1)/C_{\min} & \text{otherwise} \end{cases},$$

where  $C_{\min}$  is the smallest nonzero cost in  $C$ .

As in Unger and Moulton's formulation, small energies are only possible if the  $a_i$  lie along a line parallel to the  $x$ -axis in the three dimensional lattice. The changes made to their reduction further restrict the optimal conformation to have the  $a_i$  lie at an even distance along the line. Figure 5 illustrates the structure of conformations that can assume low energy.

It follows that each of the  $a_i$  so configured can be connected by an even-length path of  $x$ s. Unger and Moulton's arguments suffice to demonstrate that the optimal conformation is found if and only if OLA is solved, with the observation that the additional  $x$ s added to the sequence  $S$  guarantee that the  $a_i$  can be connected when spaced apart in this fashion. ■



**Fig. 6.** The conformational invariant needed for the diamond crystal lattice. The break in the chain shortens the diagonally oriented loop.

## Acknowledgements

Our thanks to Ken Dill for suggesting the extension of our previous results to other lattice models and for discussions that inspired this work. We also thank Martin Karplus for his interest in our work and for his insight into the importance of performance guaranteed approximation algorithms for protein folding. This work was supported by the Applied Mathematical Sciences program, U.S. Department of Energy, Office of Energy Research, and was performed at Sandia National Laboratories, operated for the U.S. Department of Energy under contract No. DE-AC04-94AL85000.

## References

1. N. W. Ashcroft and N. D. Mermin. *Solid State Physics*. Holt, Rinehart and Winston, 1976.
2. D. G. Covell and R. L. Jernigan. *Biochemistry*, 29:3287, 1990.
3. K. A. Dill. *Biochemistry*, 24:1501, 1985.
4. K. A. Dill, S. Bromberg, K. Yue, K. M. Fiebig, D. P. Yee, P. D. Thomas, and H. S. Chan. Principles of protein folding: A perspective from simple exact models. *Prot. Sci.*, 4:561-602, 1995.
5. A. S. Fraenkel. Complexity of protein folding. *Bull. Math. Bio.*, 55(6):1199-1210, 1993.
6. M. R. Garey and D. S. Johnson. *Computers and Intractability - A guide to the theory of NP-completeness*. W.H. Freeman and Co., 1979.
7. A. M. Gutin and E. I. Shakhnovich. Ground state of random copolymers and the discrete random energy model. *J. Chem. Phys.*, 98:8174-8177, 1993.
8. W. Hart and S. Istrail. Fast protein folding in the hydrophobic-hydrophilic model within three-eighths of optimal. To appear in *Journal of Computational Biology*, Spring 1996. Extended abstract in *Proc. of 27th Annual ACM Symposium on Theory of Computation*, May 1995.
9. M. Karplus and E. Shakhnovich. *Protein folding: Theoretical studies of thermodynamics and dynamics*, chapter 4, pages 127-195. W. H. Freeman and Company, 1993.
10. J. T. Ngo and J. Marks. Computational complexity of a problem in molecular structure prediction. *Protein Engineering*, 5(4):313-321, 1992.
11. M. Paterson, March 1995. Personal communication.
12. E. I. Shakhnovich and A. M. Gutin. Engineering of stable and fast-folding sequences of model proteins. *Proc. Natl. Acad. Sci.*, 90:7195-7199, 1993.
13. A. Sikorski and J. Skolnick. Dynamic Monte Carlo simulations of globular protein folding/unfolding pathways. II.  $\alpha$ -helical motifs. *J. Molecular Biology*, 212:819-836, July 1990.
14. J. Skolnick and A. Kolinski. Simulations of the folding of a globular protein. *Science*, 250:1121-1125, 1990.
15. R. Unger and J. Moult. Finding the lowest free energy conformation of a protein is a NP-hard problem: Proof and implications. *Bull. Math. Bio.*, 55(6):1183-1198, 1993.
16. A. F. Wells. *Three-dimensional nets and polyhedra*. American Crystallographic Association, 1979.