

# A Predictive Method for the Evaluation of Peptide Binding in Pocket 1 of HLA-DRB1 via Global Minimization of Energy Interactions

I.P. Androulakis,<sup>1</sup> N.N. Nayak,<sup>1</sup> M.G. Ierapetritou,<sup>1</sup> D.S. Monos,<sup>2</sup> and C.A. Floudas<sup>1\*</sup>

<sup>1</sup>Department of Chemical Engineering, Princeton University, Princeton, New Jersey

<sup>2</sup>Department of Pediatrics, University of Pennsylvania, The Children's Hospital of Philadelphia, Philadelphia, Pennsylvania

**ABSTRACT** Human leukocyte antigens (HLA) or histocompatibility molecules are glycoproteins that play a pivotal role in the development of an effective immune response. An important function of the HLA molecules is the ability to bind and present antigen peptides to T lymphocytes. Presently there is no comprehensive way of predicting and energetically evaluating peptide binding on HLA molecules. To quantitatively determine the binding specificity of a class II HLA molecule interacting with peptides, a novel decomposition approach based on deterministic global optimization is proposed that takes advantage of the topography of HLA binding groove, and examined the interactions of the bound peptide with the five different pockets. In particular, the main focus of this paper is the study of pocket 1 of HLA DR1 (DRB1\*0101 allele). The determination of the minimum energy conformation is based on the ECEPP/3 potential energy model that describes the energetics of the atomic interactions. The minimization of the total potential energy is formulated on the set of peptide dihedral angles, Euler angles, and translation variables to describe the relative position. The deterministic global optimization algorithm,  $\alpha$ BB, which has been shown to be  $\epsilon$ -convergent to the global minimum potential energy through the solution of a series of nonlinear convex optimization problems, is utilized. The PACK conformational energy model that utilizes the ECEPP/3 model but also allows the consideration of protein chain interactions is interfaced with  $\alpha$ BB. MSEED, a program used to calculate the solvation contribution via the area accessible to the solvent, is also interfaced with  $\alpha$ BB. Results are presented for the entire array of naturally occurring amino acids binding to pocket 1 of the HLA DR1 molecule and very good agreement with experimental binding assays is obtained. *Proteins* 29:87–102, 1997. © 1997 Wiley-Liss, Inc.

**Key words:** peptide docking; global optimization; HLA-DRB1 class II protein

## INTRODUCTION

The docking problem has received a lot of attention in the open literature. The presented methods can be classified as shape-based methods that are based on molecular surface representation and energy-based methods that optimize interaction energy in order to determine good dockings. Shape-based methods have the advantage of being less computationally intensive since the number of possible different binding modes can be greatly reduced by using a simplified model for the shapes of the receptor and binder. Based on this idea, are the works of Lee and Richard,<sup>22</sup> Connolly,<sup>8</sup> Bacon and Moul,<sup>5</sup> Jiang and Kim,<sup>19</sup> Kuntz and coworkers.<sup>17</sup>

Energy-based methods on the other hand, represent a more precise way of determining good dockings but they are more computationally demanding. Due to this fact most of the proposed approaches are based on Monte Carlo simulation and Simulated Annealing such as the works of Goodsell and Olson,<sup>12</sup> Hart and Read,<sup>16</sup> and Calfisch and coworkers.<sup>7</sup> Rosenfeld and coworkers,<sup>35</sup> present a peptide binding study based on random selection and minimization among potential peptide structures. More recently, dynamic programming optimization,<sup>21</sup> is used for optimizing the overall free energy based on a fragment assembly algorithm and molecular dynamics simulation is also utilized for studying the binding affinity of the HLA-B\*2705 protein.<sup>34</sup> All the proposed approaches identify the importance of accurate prediction, which leads to the need of establishing efficient and systematic ways of predicting the global energetically most favorable docking mode.

Contract grant sponsor: National Science Foundation; Air Force Office of Scientific Research; American Diabetes Association.

Dr. Androulakis's current address is Corporate Research Science Laboratories, Exxon Research & Engineering Co., Annandale, NJ 08801.

\*Correspondence to: C.A. Floudas, Department of Engineering, Princeton University, Princeton, NJ 08544-5263.

Received 10 December 1996; Accepted 17 April 1997

Histocompatibility molecules or human leukocyte antigens (HLA) are cell surface molecules that form complexes with self- and nonself-peptides. The HLA-peptide complex is recognized by the T-lymphocyte receptor and initiates antigen specific immune responses. HLA molecules are very polymorphic and each of them may interact with large number of peptides. Both characteristics, their polymorphism and their binding promiscuity serve the basic function of presenting a wide range of antigens to the immune system. Appropriately presented antigens induce either tolerance or an active immune response. HLA molecules are classified as class I and class II. This distinction relates to the mode of interaction with peptides as well as to their function and distribution in the tissues. The study of the differences between the two classes as well as a trial to predict the structure of class II MHC based on the structure of class I is the subject of a recent publication.<sup>20</sup>

In this study the presented results involve the class II molecule HLA-DR1 (DRB1\*0101). A deterministic global optimization approach is proposed for determining the conformation of the binding complex with the global minimum of interaction energy. This approach is applied on evaluating the total potential energy of the entire array of amino acids interacting with pocket 1 of HLA-DR1. A detailed description of the HLA-DR1 molecule is presented in the next section. In this paper we present the problem definition and formulation, the proposed global optimization approach, and the results for all naturally occurring amino acids binding in pocket 1 of HLA-DR1 discussed in comparison with provided experimental data.

## PROBLEM DEFINITION AND PROPOSED APPROACH

### HLA-DR1 Molecule

The binding of an influenza virus peptide to the MHC protein HLA-DR1<sup>38</sup> is illustrated in Figure 1. The HLA-DR1 molecule is in white, while the bound peptide is in gray with the different locations of the major protein pockets defined by the residues shown in different shades. Notice that the peptide obtains an extended conformation in the binding groove on this complex. (All pictures were created in the molecular graphics program GRASP.<sup>31</sup>)

Histocompatibility proteins are organized into two major classes. Polymorphic residues in both class I and class II proteins are clustered in the peptide-binding region and are responsible for the different peptide specificities. The major distinctive features of the class I and class II loci are (i) allograft rejection properties, (ii) relative tissue distributions, and (iii) differing chemical compositions.<sup>32</sup> Class I proteins generally bind fragments that range from 8 to 10

residues in length. Additionally, the protein pockets in this class show allele defined tendencies to bind particular amino acid side chains, or to bind unspecifically. In contrast, class II molecules bind much longer fragments and it has proven difficult to define the binding tendencies of the various pockets.<sup>38</sup> The different binding properties of these two classes are conjectured to be a result of the more open structure of class II peptides. This allows longer peptides to be situated in the MHC binding groove.<sup>32</sup>

The HLA-DR1 molecule consists of an  $\alpha$  chain (33–35 kDa), and a  $\beta$  chain (26–28 kDa) consisting of 366 amino acid residues. The  $\beta$ 1 chain of the HLA-DR1 locus is highly variable, while all other regions tend to be relatively invariant.

Crystallographic studies<sup>38</sup> have shown that peptide binding is accommodated by five polymorphic pockets on the surface of the HLA-DR1 molecule. Each of these pockets can accommodate a single amino acid residue when a particular peptide is bound.<sup>38</sup> Accordingly, these pockets play a major role in determining the peptide specificity of class II molecules. Both pocket 1 and pocket 4 have been implicated as playing vital roles in peptide binding and subsequent recognition by T cells.<sup>10,14</sup>

Pocket 1 is the largest and deepest pocket of the HLA-DR1 molecule. The area of contact for potential binders has been estimated at 200 Å<sup>2</sup>.<sup>38</sup> The pocket has been implicated as being an “anchor” peptide. It has been postulated that the residues that bind to the other four pockets are mainly determined by which residue in the binding peptide attaches to this pocket.<sup>14</sup> Pocket 1 consists of hydrophobic residues including several phenylalanine groups. This accounts for the preference of this pocket to accommodate hydrophobic residues, such as tyrosine and phenylalanine.<sup>38</sup> The large size of this pocket makes it the most solvent-accessible of the five pockets.

Pocket 4 is a relatively large pocket that is much shallower than pocket 1. The pocket consists of predominately hydrophobic amino acids, except for a positively charged arginine group. This accounts for this pocket's tendency to bind residues that have large, aliphatic side chains, or negatively charged side chains. Thus, residues such as glutamate and aspartate bind favorably to this pocket while positively charged groups such as lysine or arginine are repelled.<sup>15,38</sup> Pocket 4 has been shown to play an important role in the recognition of the bound peptide by T lymphocytes.<sup>10</sup> Pocket 4, in addition to pockets 6, 7, and 9, has been shown to be 90% inaccessible to solvent.<sup>38</sup>

The other three pockets (6, 7, and 9) are considered to affect to a lesser degree the determination of peptide binding. Pocket 6 is a shallow pocket that prefers smaller residues. The similarly shallow pocket 7 is nondiscerning in its binding activities, and only partially accommodates side chains. Pocket 9 binds

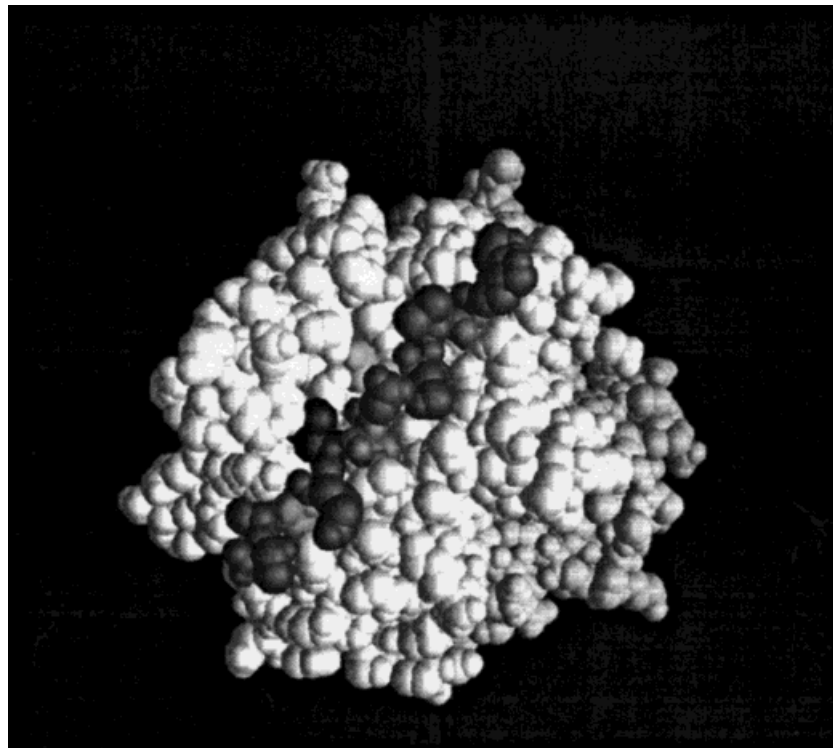


Fig. 1. HLA-DR1 bound to an influenza virus peptide.

aliphatic side chains due to its small, hydrophobic nature.<sup>38</sup>

### Proposed Approach

The modeling and optimization studies of the interactions between the HLA-DR1 protein and a virus peptide are based on a novel decomposition scheme. As it has been described in the previous section, the binding specificity of the HLA-DR1 molecule is mainly determined by the binding characteristics of its five pockets, which enables the investigation of each one separately. This paper concentrates on the study of pocket 1. The key ideas in the proposed decomposition approach are (i) to consider the binding at each pocket separately, (ii) to study the binding of each amino acid to each pocket by considering one at a time, and (iii) to create a rank ordered list of the binding amino acids for each pocket, based on an energetic criterion that reflects the binding specificity.

The proposed decomposition approach consists of the following stages.

#### Stage 1

In stage 1 the pocket of the HLA-DR1 protein is represented by a number of residues as described in detail in the subsection Pocket Definition. The work of Stern and colleagues<sup>38</sup> provides information on the constituent amino acids of each pocket of HLA-DR1 protein. Furthermore, this work provides the carte-

sian coordinates of the atoms that participate in each amino acid of the HLA-DR1 protein.

#### Stage 2

For the specific pocket interacting with each naturally occurring amino acid a mathematical model is formulated that represents all the energetic atom-to-atom interactions. These interactions are classified as (i) inter-interactions between the atoms of the residues that define the pocket of HLA-DR1 protein and the atoms of the considered naturally occurring amino acid, and (ii) intra-interactions between the atoms of the considered naturally occurring amino acid. These interactions consist of electrostatic, non-bonded, hydrogen-bonding, torsional, and loop-closing components. In addition, solvation energy is also considered based on solvent accessible areas. The detailed mathematical model and potential functions used are described in the subsections Protein Representation and Potential Energy Model.

#### Stage 3

Having a mathematical model which accounts for all the inter- and intra-interactions of the specific pocket and the considered naturally occurred amino acid, in stage 3 we formulate the global optimization problem, which minimizes the total potential energy as it is explained in detail in the subsection Problem Formulation.

#### Stage 4

A deterministic global optimization method,  $\alpha\mathbf{BB}$ ,<sup>3,4,24,25</sup> is adopted at stage 4 for the solution of the resulting nonconvex mathematical model of stage 3. This stage requires the connection of the  $\alpha\mathbf{BB}$  global optimization method with the conformation energy program PACK,<sup>37</sup> which utilizes ECEPP/3,<sup>36</sup> and the program MSEED,<sup>33</sup> that supplies the solvation contribution as described in detail in the section Deterministic Global Optimization.

#### Stage 5

In stage 5, we introduce an energetic-based criterion that allows for the comparison of the binding between a given pocket and each naturally occurring amino acid. This measure, which is denoted as  $\Delta E$ , corresponds to the difference of (i) the global minimum total potential energy that is obtained in stage 4 and which is indicated as  $E_{\text{Total}}$ , and (ii) the global minimum potential energy of the considered naturally occurring amino acid when it is far away from the pocket and which is denoted by  $E_{\text{Res}}^{\circ}$ :

$$\Delta E = E_{\text{Total}} - E_{\text{Res}}^{\circ} \quad (1)$$

Note here that the energies  $E_{\text{Total}}$  and  $E_{\text{Res}}^{\circ}$  include the consideration of the solvation energy as it will be discussed in more detail in the subsection Potential Energy Model. This criterion represents a measure of the binding affinity of each amino acid to the given pocket, in the sense that it quantifies the tendency of an amino acid to bind with the pocket of the HLA-DR1 molecule. The amino acid that exhibits the least  $\Delta E$  corresponds to the one with the best possible binding to that pocket of the HLA-DR1 protein.

#### Stage 6

In stage 6, we repeat the previous stages for each naturally occurring amino acid and hence create a rank ordered list for the binding of each of them to the specific pocket. The detailed results are shown in the section Computational Studies and Discussion.

### MATHEMATICAL MODELING

#### Protein Representation

The geometry of a protein can be fully described by defining the relative cartesian coordinates of each atom. Instead of specifying the coordinate vector for all atoms in a protein, one can specify all bond lengths, covalent bond angles and dihedral angles. Under biological conditions, the bond lengths and bond angles are fairly rigid and thus can be assumed to be fixed at their equilibrium values. Under this assumption, the dihedral angles determine the geo-

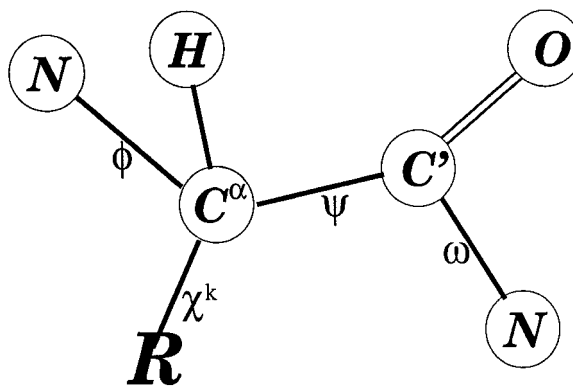


Fig. 2. Dihedral angles of a standard amino acid.

metric shape of the folded protein. The names of the dihedral angles of a folded protein chain follow a standard nomenclature as shown in Figure 2.

If more than one polypeptide is involved then the relative orientations, and locations of these different chains must be defined. This can most easily be accomplished by defining a translation vector and a rotation matrix. The translation is achieved through the cartesian coordinates of the initial nitrogen atom of each independent chain. The Euler angles specify the rotations necessary to orient a particular polypeptide and are defined as the angles between the coordinate axes defined by the initial hydrogen, nitrogen, and alpha carbon of each residue. The detailed determination of the euler angles is given in Appendix A.

#### Pocket Definition

The relative energies of minimization of each of the five protein pockets is mainly determined by the residues that constitute these pockets. A Program for Pocket Definition, denoted as PPD, constructs these pockets through the selection of all residues that are within a radius  $\mathcal{R}$  of the atoms of the crystallographic binder. A range of values for  $\mathcal{R}$  has been evaluated in an attempt to discover a radius that realistically represents the pocket, while limiting the number of residues necessary to define the pocket. The information required by the user is provided in a file containing the coordinates of the HLA-DR1 molecule and a file with the coordinates of the influenza virus binding peptide, as well as the value of radius  $\mathcal{R}$ . Each of the five pockets of the HLA-DR1 molecule were run through PPD for three different radius lengths ( $\mathcal{R} = 4.0, 4.5, 5.0$  Å). The program filters through the  $\alpha$  and  $\beta$  strands of the molecule and returns two output files one having a list of all atoms within a radius  $\mathcal{R}$ , as well as the exact distances and a PDB file with the coordinates of all residues that define a given protein pocket for the specific radius. Table I presents the residues

**TABLE I. PPD Pocket Compositions  
for  $R = 4.0\text{--}5.0 \text{ \AA}$** 

Pocket	$R = 4.0 \text{ \AA}$	$R = 4.5 \text{ \AA}$	$R = 5.0 \text{ \AA}$
1	ile $\alpha$ 31 trp $\alpha$ 43 ser $\alpha$ 53 val $\beta$ 85 phe $\beta$ 89 phe $\alpha$ 32 ala $\alpha$ 52 phe $\alpha$ 54 gly $\beta$ 86	ile $\alpha$ 31 trp $\alpha$ 43 ser $\alpha$ 53 val $\beta$ 85 phe $\beta$ 89 asn $\beta$ 82 phe $\alpha$ 32 ala $\alpha$ 52 phe $\alpha$ 54 gly $\beta$ 86 phe $\alpha$ 24	ile $\alpha$ 31 trp $\alpha$ 43 ser $\alpha$ 53 val $\beta$ 85 phe $\beta$ 89 asn $\beta$ 82 thr $\beta$ 90 phe $\alpha$ 32 ala $\alpha$ 52 phe $\alpha$ 54 gly $\beta$ 86 phe $\alpha$ 24 glu $\alpha$ 55
4	gln $\alpha$ 09 phe $\beta$ 13 arg $\beta$ 71 tyr $\beta$ 78 asn $\alpha$ 62 gln $\beta$ 70 ala $\beta$ 74	gln $\alpha$ 09 phe $\beta$ 13 arg $\beta$ 71 tyr $\beta$ 78 leu $\beta$ 26 asn $\alpha$ 62 gln $\beta$ 70 ala $\beta$ 74 glu $\alpha$ 11	gln $\alpha$ 09 phe $\beta$ 13 arg $\beta$ 71 tyr $\beta$ 78 leu $\beta$ 26 asn $\alpha$ 62 gln $\beta$ 70 ala $\beta$ 74 glu $\alpha$ 11
6	glu $\alpha$ 11 val $\alpha$ 65 leu $\beta$ 11 asn $\alpha$ 62 asp $\alpha$ 66	glu $\alpha$ 11 val $\alpha$ 65 leu $\beta$ 11 asn $\alpha$ 62 asp $\alpha$ 66	glu $\alpha$ 11 val $\alpha$ 65 leu $\beta$ 11 arg $\beta$ 71 asn $\alpha$ 62 asp $\alpha$ 66 phe $\beta$ 13
7	val $\alpha$ 65 glu $\beta$ 28 trp $\beta$ 61 arg $\beta$ 71 asn $\alpha$ 69 tyr $\beta$ 47 leu $\beta$ 67	val $\alpha$ 65 glu $\beta$ 28 trp $\beta$ 61 arg $\beta$ 71 asn $\alpha$ 69 tyr $\beta$ 47 leu $\beta$ 67	val $\alpha$ 65 glu $\beta$ 28 trp $\beta$ 61 arg $\beta$ 71 asn $\alpha$ 69 tyr $\beta$ 47 leu $\beta$ 67
9	ile $\alpha$ 72 met $\alpha$ 73 trp $\beta$ 09 tyr $\beta$ 60 asn $\alpha$ 69 arg $\alpha$ 76 asp $\beta$ 57	ile $\alpha$ 72 met $\alpha$ 73 trp $\beta$ 09 tyr $\beta$ 60 asn $\alpha$ 69 arg $\alpha$ 76 asp $\beta$ 57 trp $\beta$ 61	ile $\alpha$ 72 met $\alpha$ 73 trp $\beta$ 09 tyr $\beta$ 60 leu $\alpha$ 70 asn $\alpha$ 69 arg $\alpha$ 76 asp $\beta$ 57 trp $\beta$ 61

defining each of the protein pockets for the various radii considered. There are several important observations that should be made. First, there is the intuitive trend of the pockets becoming more complex with increased radius. Second, note that pocket 1 is composed of a significantly larger number of residues than any of the other pockets. Finally, note that in some cases an increase in radius does not

cause the inclusion of additional amino acids. An example of this is pocket 7, where the pocket is identically defined across the entire range of radii.

The resulting PDB file is then translated to the internal coordinate system by a program denoted as ARAS (Amino acid Residue Angle Solver). The output file obtained by this program is the one required by the conformational energy program PACK to evaluate the potential energy as is described in detail in the next section.

### Potential Energy Model

Molecular dynamics calculations employ an empirical derived set of potential energy contributions for approximating the force field of the protein system. These energy functions are based on specific types of interactions instead of being associated with a particular molecule. The parameters for these correlations have been determined to provide the best possible agreement with experimental data. Many different parameterizations have been proposed for approximating the *force field* in protein folding calculations. Some of the most popular ones are ECEPP,<sup>26</sup> MM2,<sup>1</sup> ECEPP/2,<sup>30</sup> CHARMM,<sup>6</sup> DISCOVER,<sup>9</sup> AMBER,<sup>41,42</sup> GROMOS,<sup>39</sup> ENCAD,<sup>23</sup> MM3,<sup>2</sup> AND ECEPP/3.<sup>29</sup> In this work the ECEPP/3<sup>29</sup> detailed potential model is utilized. In this potential model, it is assumed that the covalent bond lengths and angles are fixed at their equilibrium values and the conformational energy is treated as the sum of electrostatic, nonbonded, hydrogen bonding, torsional, and cystine contributions.

The potential function of ECEPP/3 includes the following terms:

$$\begin{aligned}
 E = & \sum_{(i,j) \in ES} 332.0 \frac{q_i q_j}{D r_{ij}} \quad (\text{Electrostatic}) \\
 & + \sum_{(i,j) \in NB} F \frac{A}{r_{ij}^{12}} - \frac{C}{r_{ij}^6} \quad (\text{Nonbonded}) \\
 & + \sum_{(hx) \in HX} F \frac{A'}{r_{hx}^{12}} - \frac{B}{r_{hx}^{10}} \quad (\text{Hydrogen bonding}) \\
 & + \sum_{k \in TOR} \left( \frac{E_0}{2} \right) (1 \pm \cos n_k \theta_k) \quad (\text{Torsional}) \\
 & + \sum_{i \in COOP} B_L \sum_{i=1}^{i=3} (r_{i_1} - r_{i_0})^2 \quad (\text{Cystine loop-closing}) \\
 & + \sum_{i \in COOP} A_L (r_{4_1} - r_{4_0})^2 \quad (\text{Cystine torsional})
 \end{aligned}$$

In addition, the solvation energy is also considered through the utilization of the program MSEED,<sup>33</sup>

which supplies solvent accessible areas. Once these areas have been calculated, the following formula can be utilized to define the solvation potential:

$$E_{\text{SOL}} = \sum_{i=1}^n \sigma_{k(i)} A_i \quad (2)$$

where  $n$  equals the total number of atoms in the molecule,  $\sigma_{k(i)}$  is a coefficient dependent upon the atom type, and  $A_i$  is the solvent accessible area of the  $i$ th atom. The  $\sigma$  coefficients were determined by the research performed in Ref. 40.

The solvent accessible area is determined by rolling a spherical test probe over the surface of the molecule (see Fig. 3). The areas of direct contact between the molecule and the probe define the accessible surface. Additionally, the area of the bottom most part of the probe traces the surface in inaccessible cavities of the protein. The probe radius is equivalent to the van der Waals radius of a water molecule, which is equivalent to 1.4 Å. These empirical solvent-accessible surface areas are calculated by the program MSEED. This program utilizes Connolly's analytical algorithm, which is described in Ref. 33. Note that  $E_{\text{SOL}}$  is only added to this overall potential at local minima, and hence is not explicitly stated in the above equation. This is done because the parameters of the JRF set used in Ref. 40 were derived based on a set of tetrapeptide conformations that correspond to local minima of the ECEPP potential energy.<sup>38</sup> The total energy  $E_{\text{Total}}$  is then defined as:

$$E_{\text{Total}} = E + E_{\text{SOL}}$$

### Problem Formulation

As it was described in the subsection Protein Representation a particular amino acid chain could be defined by a translation vector, a rotation matrix, and the corresponding set of dihedral angles. The translation vector will be defined as the coordinates of the nitrogen atom on the first residue of a chain, while the rotation matrix will be defined by the Euler angles. Since the pocket is considered to be rigid, the only variables will be the amino coordinates, Euler angles, and dihedral angles of the amino acid binder.

Let  $k = 1, \dots, K$ , where  $K$  is the total number of side chain angles of the amino acid residue that attempts to bind the pocket. Then, the set of variable dihedral angles would include the backbone angles ( $\phi$ ,  $\psi$ , and  $\omega$ ), and the side-chain angles ( $\chi^k$ ). The cartesian coordinates of the amino translation vector will be defined by the variables  $N_x$ ,  $N_y$ , and  $N_z$ . Similarly, the cartesian coordinates of the backbone carboxyl carbon are represented by  $C'_x$ ,  $C'_y$ , and  $C'_z$ . Finally, the Euler angles will be represented by  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3$ . Utilizing the above definitions the potential

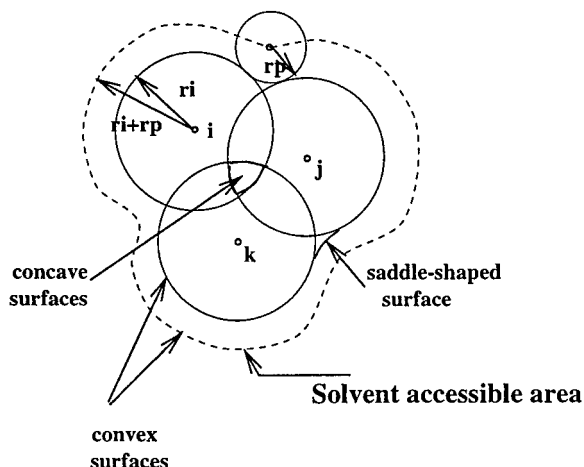


Fig. 3. Determination of solvent-accessible area.

energy minimization problem can be formulated as follows:

$$\min E(\phi, \psi, \omega, \chi^k, N_x, N_y, N_z, \epsilon_1, \epsilon_2, \epsilon_3) \quad (3)$$

$$\text{s.t. } -\pi \leq \phi \leq \pi \quad (4)$$

$$-\pi \leq \psi \leq \pi \quad (5)$$

$$-\pi \leq \omega \leq \pi \quad (6)$$

$$-\pi \leq \chi^k \leq \pi, k = 1, \dots, K \quad (7)$$

$$-\pi \leq \epsilon_1 \leq \pi \quad (8)$$

$$-\pi \leq \epsilon_2 \leq \pi \quad (9)$$

$$-\pi \leq \epsilon_3 \leq \pi \quad (10)$$

$$N_x^l \leq N_x \leq N_x^u \quad (11)$$

$$N_y^l \leq N_y \leq N_y^u \quad (12)$$

$$N_z^l \leq N_z \leq N_z^u \quad (13)$$

$$C'_x{}^l \leq C'_x(\phi, \psi, \omega, \chi^k, N_x, N_y, N_z, \epsilon_1, \epsilon_2, \epsilon_3) \leq C'_x{}^u \quad (14)$$

$$C'_y{}^l \leq C'_y(\phi, \psi, \omega, \chi^k, N_x, N_y, N_z, \epsilon_1, \epsilon_2, \epsilon_3) \leq C'_y{}^u \quad (15)$$

$$C'_z{}^l \leq C'_z(\phi, \psi, \omega, \chi^k, N_x, N_y, N_z, \epsilon_1, \epsilon_2, \epsilon_3) \leq C'_z{}^u \quad (16)$$

Note that the superscripts  $u$  and  $l$  denote upper and lower bounds, respectively, for the cartesian coordinates of both the amino nitrogen and the carboxyl carbon. In addition to the constraints on the amino nitrogen, note that in the above formulation there are additional constraints on the carboxyl carbon. It has been assumed due to the decomposition employed that enables the consideration of each pocket separately, that the conformational move-

ments of the binding peptide are only constrained by the locations of these two atoms. In the original problem though, the binding residue is part of a longer antigen peptide. The rest of the binding peptide is assumed to bind normally so even if the binding residue is changed these backbone atoms will be relatively confined to their initial positions due to their peptidic linkages. Although the constraints of amino nitrogen can be directly considered in the above formulation since they correspond to problem variables, the  $C'$  coordinates are not explicit variables and consequently they must be defined as a function of the other variables (see Appendix B). Note that  $E$  is a nonconvex function involving numerous local minima that correspond to metastable states of the specific amino acid binding to the pocket 1. A single global minimum defines the energetically most favorable peptide conformation. In establishing a ranked-order list of binding peptides, one needs to identify rigorously the best conformation of (i) the binding residue far from the pocket and (ii) the complex of pocket 1 with the binding residue. Consequently, there is a need for a method that can guarantee convergence to the global minimum potential energy conformation and which is capable of solving large scale constrained optimization problems. In this paper, the global optimization approach  $\alpha\mathbf{BB}$ ,<sup>3,4,25</sup> described in more detail in the next section, has been extended to peptide systems interacting with realistic atomistic potential energy models (e.g., ECEPP/3<sup>36</sup>), which include solvation contributions via surface accessible area as the solvent method (e.g., MSEED<sup>33</sup>).

### DETERMINISTIC GLOBAL OPTIMIZATION Global Optimization Approach

The global optimization scheme  $\alpha\mathbf{BB}$ <sup>3,4,25</sup> is a deterministic branch and bound algorithm for locating the global optimum based on the construction of converging lower and upper bounds. Upper bounds can be simply obtained by minimizing  $E$  using local methods. Lower bounds can be evaluated by constructing the convex underestimator,  $L$ , of the original function  $E$  and evaluating the single global minimum of the resulting convex problem.

A convex lower bounding function  $L$  of potential energy function  $E$  can be defined by augmenting  $E$  using the ideas of the approach introduced in Ref. 25:

$$\begin{aligned}
 L = & E + \alpha[(\phi^l - \phi)(\phi^u - \phi) + (\psi^l - \psi)(\psi^u - \psi) \\
 & + (\omega^l - \omega)(\omega^u - \omega) + \sum_{k=1}^K (\chi^{k,l} - \chi^k)(\chi^{k,u} - \chi^k) \\
 & + (N_x^l - N_x)(N_x^u - N_x) + (N_y^l - N_y)(N_y^u - N_y) \\
 & + (N_z^l - N_z)(N_z^u - N_z) + (\epsilon_1^l - \epsilon_1)(\epsilon_1^u - \epsilon_1) \\
 & + (\epsilon_2^l - \epsilon_2)(\epsilon_2^u - \epsilon_2) + (\epsilon_3^l - \epsilon_3)(\epsilon_3^u - \epsilon_3)]
 \end{aligned}$$

where  $\alpha$  is a nonnegative parameter which must be greater or equal to the negative one half of the minimum eigenvalue of the hessian of  $E$  over the rectangular under consideration described by the lower and upper bounds of the involved variables defined by the superscripts  $l$  and  $u$ , respectively. The following properties of function  $L$  will enable the construction of a global optimization algorithm. These properties whose proof is given in Ref. 25 demonstrate that

1.  $L$  is always a valid underestimator of  $E$ .
2.  $L$  matches  $E$  at all corner points of the box constraints.
3.  $L$  is convex.
4. The maximum separation between  $L$  and  $E$  is bounded and proportional to  $\alpha$  and to square of the diagonal of the current box constraints. This property ensures that an  $\epsilon_f$  feasibility and  $\epsilon_c$  convergence tolerances can be reached for a finite size partition element.
5. The underestimators  $L$  constructed over supersets of the current set are always less tight than the underestimator constructed over the current box constraints for every point within the current box constraints.

These bounds are successively refined by iteratively partitioning the initial feasible region into smaller ones. The feasible region partition is achieved by subdivision of a rectangle into two subrectangles by halving along the longest side of the initial rectangle (bisection). At each iteration the lower bound would be the minimum over all the minima in every subrectangle composing the original domain. Therefore, a simple way to produce a nondecreasing sequence of lower bounds is to halve only the subrectangle responsible for the infimum of the minima. A nonincreasing sequence of upper bounds can also be produced by solving locally the nonconvex problem and selecting the minimum over the previously recorded upper bounds. Based on this procedure a fathoming step of the algorithm leads to no further consideration of a subrectangle where the minimum is greater than the current upper bound. Convergence proof to an  $\epsilon$ -global solution in finite steps is given in Ref. 25.

### Algorithmic Description

The proposed approach for the determination of the global minimum of  $E$  that corresponds to the peptide conformation binding to the pocket 1 of HLA-DR1 as posed in the section Problem Definition and Proposed Approach, necessitates the development of an optimization interface that combines the global optimization program  $\alpha\mathbf{BB}$ , the conformational energy program PACK which utilizes ECEPP/3, the solvation program MSEED, and the local optimization solver NPSOL. Additional pro-

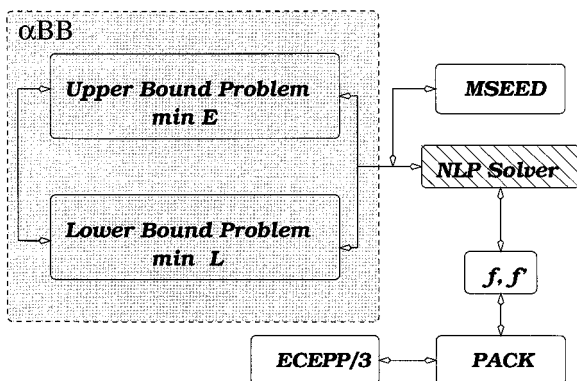


Fig. 4. The interface for global optimization.

gram files serve to link these programs. A schematic diagram of the interface between the used programs is shown in Figure 4.

The following steps are required in the calculation of the global minimum:

1. The local solver (NPSOL<sup>11</sup>) obtains a local minimum of the potential function supplied by PACK in a domain (rectangle) defined by the original lower and upper bounds of the variables (bounds are supplied by  $\alpha\mathbf{BB}$ ). PACK determines the energies of individual chains through repeated calls to ECEPP/3.
2. The solvation energy at this local minimum is calculated by MSEED. This hydration energy is added to the potential function to yield  $E_{\text{Total}}$ , which will serve as an upper bound on the global minimum solution in the current rectangle.
3. The current best upper bound is updated to be the minimum of those thus far stored.
4. The current rectangle is partitioned by bisection along the longest side.
5. The convex function  $L$  is minimized in each rectangle and the solvation energy is added at the minimum. If a solution is greater than the best upper bound it will be eliminated, otherwise it will be kept on the stack.
6. The rectangle with the current minimum solution for  $L$  is selected for further partitioning.
7. If the best upper and lower bounds are within  $\epsilon$  the program will terminate, otherwise it will proceed to step 1.

It should be noted that ECEPP/3 calculates the potential energy function for a polypeptide chain,<sup>36</sup> while PACK is a program that inputs multiple chain data and makes the appropriate calls to ECEPP/3 for calculation of the interaction energies.<sup>37</sup> Special type penalty functions had to be added to the upper bound function,  $E$ , in order to implement the previously discussed constraints on  $C'$  (Equations 14–16). The

TABLE II. Bounds on  $N$  and  $C'$  for Pockets 1 of HLA-DRB1

Bounds	Lower	Upper
$N$		
$x$	-9.2	-8.4
$y$	23.3	24.1
$z$	17.3	18.1
$C'$		
$x$	-10.0	-9.0
$y$	25.0	27.0
$z$	16.0	18.0

modified objective function takes then the following form:

$$\begin{aligned}
 E' = E + \beta \{ & (C'_x{}^l - C'_x) + \langle C'_x - C'_x{}^u \rangle \\
 & + \langle C'_y{}^l - C'_y \rangle + \langle C'_y - C'_y{}^u \rangle \\
 & + \langle C'_z{}^l - C'_z \rangle + \langle C'_z - C'_z{}^u \rangle \}
 \end{aligned}$$

The  $\langle \rangle$  function is defined as follows:  $\langle A \rangle$  equals  $A$  if  $A$  is greater than zero, otherwise  $\langle A \rangle$  equals zero. Thus, as long as the coordinates are within the defined bounds the objective function will not be modified. Yet, if a particular coordinate falls outside of the bounds, the function will be increased by the value of the transgression multiplied by the arbitrarily large constant  $\beta$ .

Since the pocket is assumed rigid, the optimization variables are the dihedral angles, the translation vector and the euler angles of the amino acid under consideration. These variables are partitioned into three sets. The first one (i.e., global variables) consists of the variables where branching occurs; the second set (i.e., local variables) consists of the variables where branching is not performed, and the third set (i.e., fixed variables) includes the variables for which there exists sufficient experimental evidence for keeping them fixed.

The information required by the user is provided in four different files. The first one is required by PACK and contains information about the different protein chains to be considered. The second one is needed by ECEPP/3 and contains information about the sequence and number of the amino acid residues and the type of end groups. It also initializes the dihedral angles, translation vector and Euler angles. The third file provides the bounds on the amino nitrogen and carboxyl carbon of the binding residue. These bounds are defined around the cartesian coordinates of these two atoms of the influenza virus peptide.<sup>38</sup> Thus, for each pocket bounds were set in the  $x$ ,  $y$ ,  $z$  directions around the coordinates of the corresponding atoms of the influenza virus peptide presented by Stern and colleagues.<sup>38</sup> These bounds are given in Table II.



## COMPUTATIONAL STUDIES AND DISCUSSION

In this section, each one of the naturally occurring amino acids will be examined and accessed regarding its binding affinity with pocket 1 of HLA-DR1 molecule. Before presenting the results obtained by the proposed approach, the following point regarding the amino acid polarity should be made. The aliphatic side chains the amino acids Ala, Val, Ile, and Leu can be clearly considered as nonpolar ones, whereas, at the opposite end of the polarity scale, are the charged residues Glu, Asp, Arg, Lys. Asn and Gly, which have amide side chains, as well as the hydroxylic amino acids Ser and Thr are polar and expected to interact strongly with water and have high solubility. The polarity of the rest of the amino acids is more ambiguous. Cys and His have  $pK_a$  values close to 7 and may actually be charged in many proteins under physiological conditions. In our computational studies we consider as positively charged residues the Arg+, His+, Lys+, and as negatively charged residues the Asp- and Glu-, using the parameters of ECEPP/3 for evaluating their energy contributions. In Tyr the aromatic ring compensates for the hydroxyl and makes Tyr a nonpolar residue. Gly and Pro are special but from the way they behave in proteins can be considered as nonpolar and polar, respectively.

### Individual Solvated Residues Away From Pocket 1

As it has been mentioned in the section Problem Definition and Proposed Approach, in order to evaluate the energy of interaction of the 20 naturally occurring amino acids within a pocket, the intramolecular energy due to atomic interactions between the atoms of the single residue far away from the pocket has to be calculated. Thus, the global minimum energy for each residue in isolation is found with the consideration of the solvation contribution as described in the subsection Potential Energy Model. The results obtained by applying the global optimization algorithm,  $\alpha\mathbf{BB}^{3,4,25}$  are shown in Table III, where, based on the previous remark regarding amino acid polarity, some of them are considered charged.

### Complex of Pocket 1 and Binding Amino Acids

As mentioned earlier,  $\Delta E$ , defined as the difference  $E_{\text{Total}} - E_{\text{Res}}^0$ , has been considered to represent the binding potential of a specific naturally occurring amino acid to the pocket considered. As shown in Table I, the number of amino acids included in pocket 1 increase as R increases. Specifically from 4.5 to 5.0 Å the amino acids that were added involve threonine and glutamic acid. Note that Glu is negatively charged and is an important factor for evaluat-

**TABLE III. Standard Energies for Individual Solvated Residues**

Residue	Code	$E_{\text{Res}}^{0,S}$ (kcal/mol)
Ala	A	-42.143
Asn	N	-94.376
Cys	C	-74.667
Gln	Q	-86.964
Gly	G	-56.260
Ile	I	-17.074
Leu	L	-23.166
Met	M	-46.269
Phe	F	-160.850
Ser	S	-82.476
Thr	T	-69.423
Trp	W	-184.230
Tyr	Y	-178.950
Val	V	-25.055
Glu-	E-	-56.607
Asp-	D-	-67.416
His+	H+	-125.460
Arg+	R+	-105.800
Lys+	K+	-27.706

ing the interactions with positive charged residues as illustrated in Table V.

The results for pocket 1 with  $R = 5.0$  Å are presented in Tables IV and V and in Figure 5a. Trp, Tyr and Phe are found to have the strongest binding affinities with interaction energies in a range of -20.0 to -16.950. At lower positions in the middle of the list there are the Leu, Ile, and Val having interaction energies between -12.481 and -11.209. At the bottom of the list are the negative charged residues Glu- and Asp- with 40% smaller interaction energy than that of Val. An interesting result of the theoretical studies is the one obtained for the positive charged residues that appear to be the most unfavorable binders for this pocket.

A series of competitive binding assays was performed that involved analogs of the HA peptide (306-318) and the DR1 molecule.<sup>28</sup> Since the HA (306-318) peptide residue that interacts with Pocket 1 is Y(308) a number of analog peptides were synthesized that substituted the Y(308) with 11 different amino acids. The relative binding affinity was derived from the reciprocal of 50% inhibitory concentration (IC50) of each analog peptide in a logarithmic scale. Figures 5 and 6 show the results of both the theoretical and the experimental results, respectively.

Based on the competitive binding assays shown in Figure 6, three groups of binding affinities have been identified. The first group includes the amino acids Trp, Tyr, and Phe, that are the residues with the highest affinity to DR1. The second group includes the amino acids Ile, Leu, and Val and are characterized by an intermediate level of affinity to DR1. The third group finally consists of low level affinity amino

**TABLE IV. Relative Energies for Solvated Residues in Pocket 1 ( $R = 5.0 \text{ \AA}$ )**

Residue	$E_{\text{Total}}^S$ (kcal/mol)	$E_{\text{Res}}^{0,S}$ (kcal/mol)	$\Delta E$ (kcal/mol)
Tyr	-198.950	-178.950	-20.00
Phe	-180.475	-160.850	-19.625
Trp	-201.180	-184.230	-16.950
Gln	-102.360	-86.964	-15.396
Met	-60.212	-46.269	-13.943
Asn	-108.160	-94.376	-13.784
Thr	-82.718	-69.423	-13.297
Leu	-35.647	-23.166	-12.481
Ile	-29.539	-17.074	-12.465
Ser	-94.033	-82.476	-11.557
Cys	-85.947	-74.667	-11.280
Val	-36.264	-25.055	-11.209
Ala	-52.498	-42.143	-10.355
Gly	-66.351	-56.260	-10.091
Glu-	-64.531	-56.607	-7.744
Asp-	-69.847	-67.416	-2.431

**TABLE V. Relative Energies for Solvated Positive Charged Residues in Pocket 1 ( $R = 5.0 \text{ \AA}$ )**

Residue	$E_{\text{Total}}^S$ (kcal/mol)	$E_{\text{Res}}^{0,S}$ (kcal/mol)	$\Delta E$ (kcal/mol)
His+	-58.374	-124.870	+66.496
Lys+	+196.15	-27.706	+223.856
Arg+	+182.78	-105.640	+288.420

acids. This group involves the charged residues Asp-, Glu-, Arg+, His+ as well as the amino acids serine and threonine.

Based on the theoretical predictions, shown in Figure 5, Trp, Tyr, and Phe are at the top positions of the rank-ordered list of the examined naturally occurring amino acids a result that is further supported from the strong preference of this pocket for large hydrophobic side chains. Furthermore, the amino acids Leu, Ile, and Val were found by the optimization studies to be characterized by potential energies that correspond to 7th, 8th, and 11th position on the ordered list, respectively. The binding assays resulted in intermediate affinities for these amino acids. The  $\Delta E$  value of  $-11.809$  kcal/mol for the binding of Val, reflects an approximate increase of 43% as compared to Tyr. Provided that an increase of 15% in potential energy defines the group of strong binders (trp, tyr, phe), an increase of up to 43% could very well reflect a group of intermediate level binders. At the bottom of Table IV, the global optimization studies put the charged residues which is also in agreement with experimental data. An increase of approximately 31% between  $\Delta E$  values of Val and Glu- reflects the low-affinity group of amino acids.

The laboratory studies present serine and threonine as relatively weak binders. The hydroxyl groups on both of these residues would favor interaction with polar molecules. Thus, weak interactions with the hydrophobic pocket 1 would be a predictable

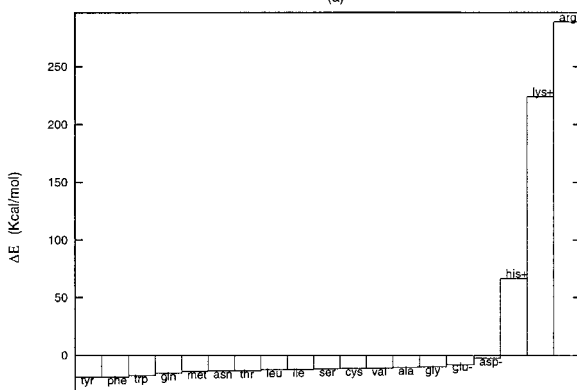
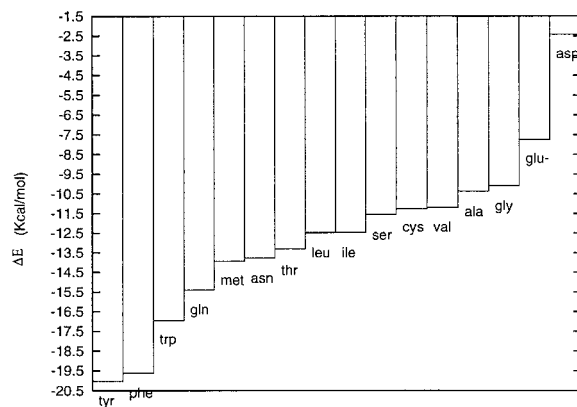
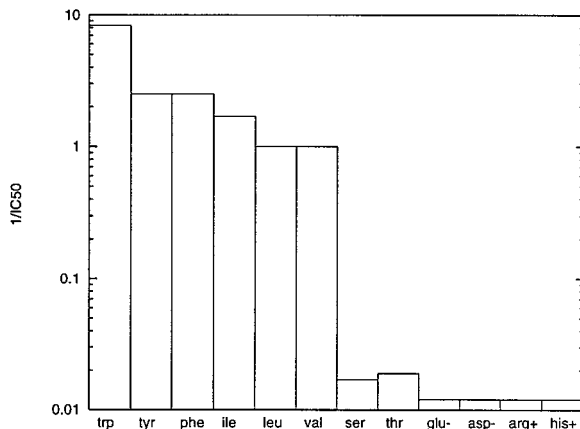
Fig. 5.  $\Delta E$  (kcal/mol) of the naturally occurring amino acids.

Fig. 6. Experimental data for the naturally occurring amino acids.

consequence. Although, valine appears to have comparable binding energy with these two residues again the optimization results support the observation of these residues being weaker binders than the small aliphatic residues.

Finally, there is a number of amino acids including Gln, Lys+, Met, Asn, Cys, Gly and Ala for which no analogs were synthesized. However, it has been reported in Refs. 18 and 13 that the peptides with

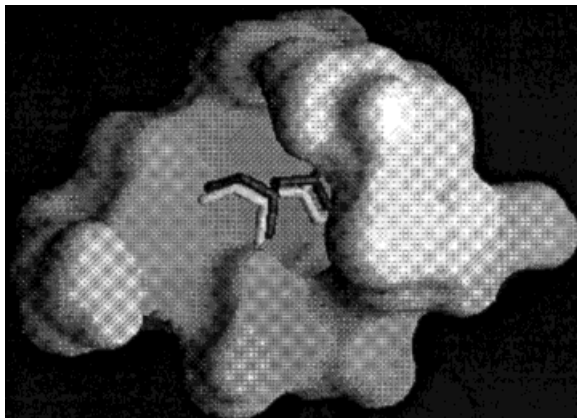


Fig. 7. Comparison of tyrosine binding to pocket 1.

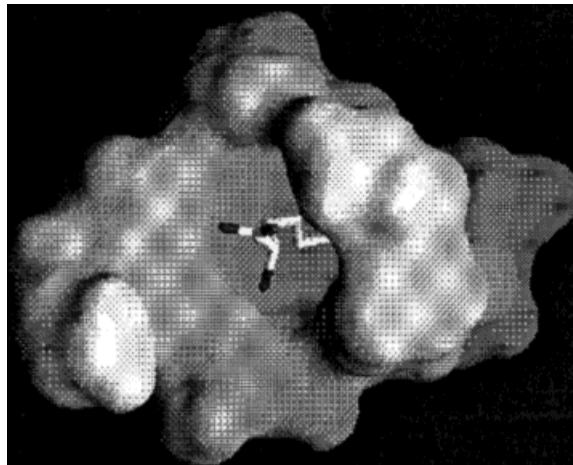


Fig. 9. Phenylalanine binding to pocket 1.

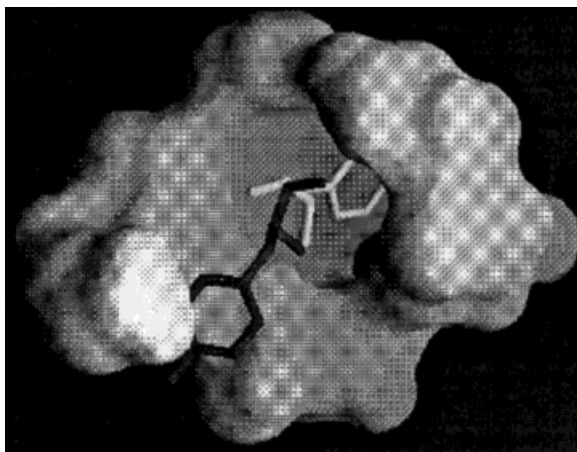


Fig. 8. Local vs global minimum configuration of tyrosine.

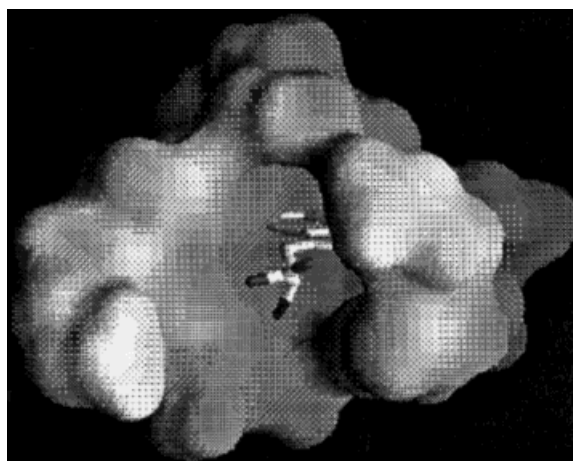


Fig. 10. Tryptophan binding to pocket 1.

Met are intermediate binders while peptides with Ala result in loss of peptide binding. The values of  $\Delta E$  equal to  $-13.943$ ,  $-10.355$  for Met and Ala, respectively, found from our global optimization studies are consistent with the reported binding studies.

Therefore, the theoretical results are in excellent agreement with those obtained by the experimental approach of competitive binding assays.<sup>28</sup>

Moreover, since the optimization interface produces a PDB file of the coordinates of the minimum energy conformation of the binder a direct comparison with the crystallographic data can be made for the tyrosine residue that binds to pocket 1. Figure 7 shows the HA peptide binder (tyrosine 308) in white and the minimum conformation of tyrosine for pocket 1 in gray. An almost identical orientation with  $1.28 \text{ \AA}$  is observed. Figures 9 and 10 illustrate the orientation of phenylalanine and tryptophan in comparison to the virus peptide binding shown in Figure 11, suggest that these residues are in fact very strong binders.

The need for determining the global minimum conformation is illustrated in Figure 8 where a local

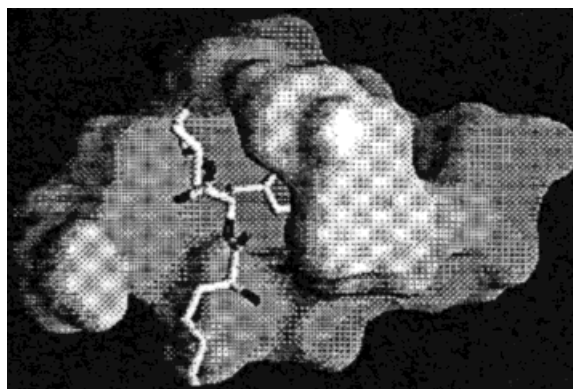


Fig. 11. Influenza virus peptide binding to pocket 1.

minimum conformation of tyrosine corresponding to  $-196.637 \text{ kcal/mol}$ , that is having only 1.16% difference from the global minimum of  $-198.95 \text{ kcal/mol}$  is illustrated with the grey whereas the global minimum conformation is shown with white. Note that

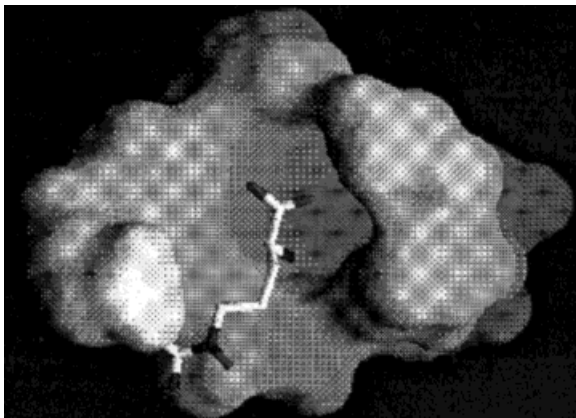


Fig. 12. Arginine+ binding to pocket 1.

the local minimum energy conformation of tyrosine is very different than the global minimum despite its proximity to the global minimum energy value.

For the charged residues, experimental data suggest their weak binding affinity. The fact that pocket 1 is extremely hydrophobic region intuitively verifies this result. Charged residues are not stabilized by the weak van der Waals interactions that stabilize the conglomeration of hydrophobic residues. The hypothesis that results from this knowledge is that the infusion of charge on the five aforementioned residues should greatly destabilize their interactions with pocket 1, which corresponds to an increase in their overall conformational energies ( $\Delta E$ ). The theoretical results obtained for the negative charged residues support this idea. Positive charged residues have large binding energy due to large electrostatic contribution from the interaction with the negative charged glutamic within the pocket, which, however, does not suggest favorable binders as their orientation shown in Figures 12–14 for arginine+, histidine+, and lysine+, respectively. The enforcement of these residues inside the pocket gives rise to large (positive) energies (Table V, Fig. 5b) that indicate highly unfavorable residues (see Figs. 15–17).

#### SUMMARY AND CONCLUSIONS

In this paper, a novel predictive method is proposed for modeling and studying the binding affinity of different naturally occurring amino acids with pocket 1 of the HLA-DR1 protein. First, the composition of the pocket is identified together with the cartesian coordinates of the atoms that participate in each amino acid of the pocket 1 of the HLA-DR1 protein. Second, explicit relations for all the energetic inter- and intra-interactions between the atoms of the residues that define the pocket of the HLA-DR1 protein and the atoms of the considered naturally occurring amino acid were derived. Moreover, solvation energy was also taken into account based on solvent-accessible area method. Then, the docking problem is formulated as a nonconvex opti-

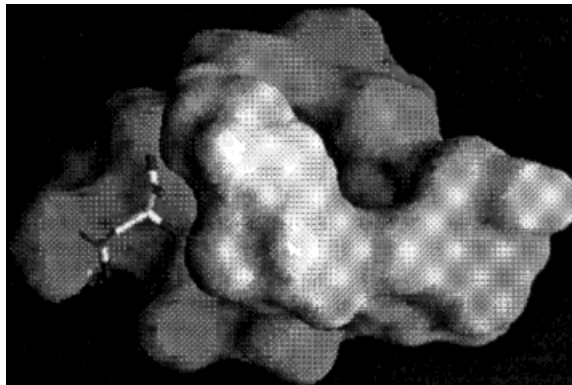


Fig. 13. Histidine+ binding to pocket 1.

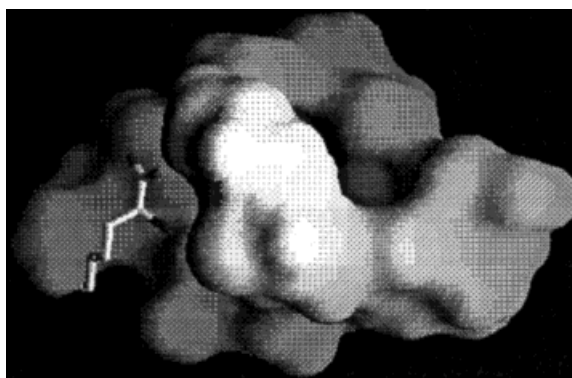


Fig. 14. Lysine+ binding to pocket 1.

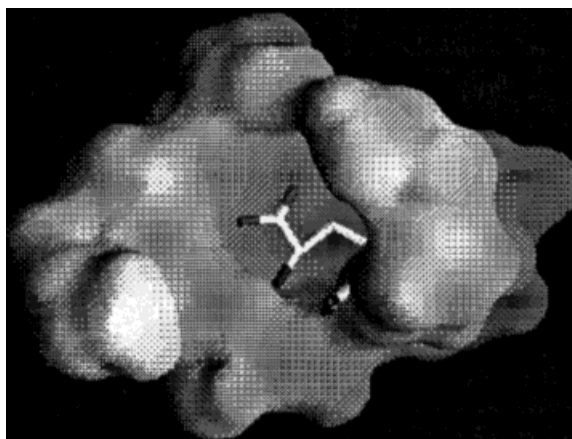


Fig. 15. Arginine+ forced within pocket 1.

mization problem on a set of independent dihedral angles, Euler angles and translation variables. The deterministic global optimum method,  $\alpha\mathbf{BB}$ , is then adopted for the solution of the resulting problem which is based on the generation of a sequence of converging upper and lower bounds found from the local solution of the nonconvex problem and the convex lower bounding problem which is constructed based on eigenvalue analysis of the nonconvex poten-

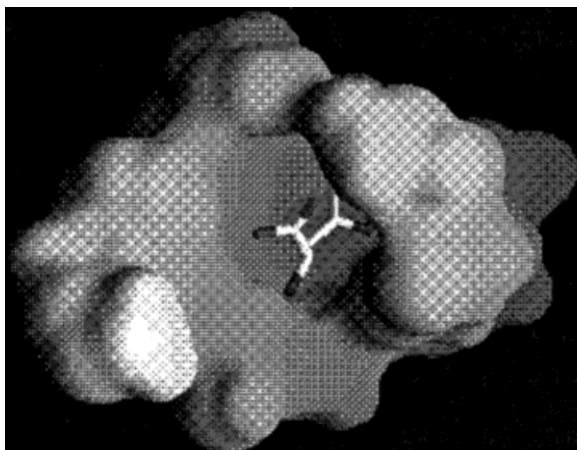


Fig. 16. Histidine+ forced within pocket 1.

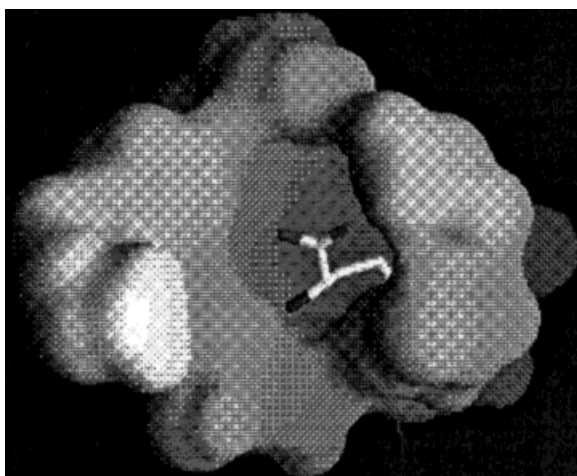


Fig. 17. Lysine+ forced within pocket 1.

tial energy function. The final step of the proposed approach consists of evaluating the interaction energy for all naturally occurring amino acids and generate a ranked-order list.

The results of the proposed approach were found to agree very well with the experimental competitive binding assays. It should be emphasized that, although in this paper only one of the binding sites of the HLA-DR1 protein is examined, the approach is applicable to predict the binding affinity of the amino acid residues in the different pockets. Our current work focuses on studying the binding information for the different pockets in order to be able to predict the binding of the whole peptide to the HLA-DR1 protein.

#### ACKNOWLEDGEMENTS

Financial support from the National Science Foundation and the Air Force Office of Scientific Research is gratefully acknowledged.

#### REFERENCES

- Allinger, N.L. Conformational analysis. MM2: a hydrocarbon force field utilizing  $V_1$  and  $V_2$  torsional terms. *J. Am. Chem. Soc.* 99:8127, 1977.
- Allinger, N.L., Yuh, Y.H., Liu, J.-H. Molecular mechanics: the MM3 force field for hydrocarbons. *J. Am. Chem. Soc.* 111:8551, 1989.
- Androulakis, I.P., Maranas, C.D., Floudas, C.A.  $\alpha$ BB: A global optimization method for general constrained nonconvex problems. *J. Global Optimiz.* 7:337-363, 1995.
- Androulakis, I.P., Maranas, C.D., Floudas, C.A. Prediction of oligopeptide conformations via deterministic global optimization. *J. Global Optimiz.* 11:1-34, 1997.
- Bacon, D.J., Moulton, J. Docking by least-square fitting of molecular surface patterns. *J. Mol. Biol.* 225:849-858, 1992.
- Brooks, B., Bruccoleri, R., Olafson, B., States, D., Swaminathan, S., Karplus, M. CHARM: a program for macromolecular energy minimization and dynamics calculation. *J. Comp. Chem.* 8:132, 1983.
- Calfisch, A., Niederer, P., Anliker, M. Monte Carlo docking of oligopeptides to proteins. *Proteins* 13:223-230, 1992.
- Connolly, M.L. Solvent-accessible surfaces of proteins and nucleic acids. *Science* 221:709, 1983.
- Dauber-Osguthorpe, P., Roberts, V.A., Osguthorpe, D.J., Wolff, J., Genest, M., Hagler, A.T. Structure and energetics of ligand binding to peptides: *Escherichia coli* dihydrofolate reductase-trimethoprim, a drug receptor system. *Proteins* 4:31, 1988.
- Fu, X., Bono, C., Woulfe, S., Swearingen, C., Summers, N., Sinigaglia, F., Sette, A., Schwartz, B., Carr, R.W. Pocket 4 of the HLA-DR molecule is a major determinant of T cell recognition of peptide. *J. Exp. Med.* 181:915-926, 1995.
- Gill, P., Murray, W., Saunders, M., Wright, M. User's Guide for NPSOL (Version 4.0): A Fortran Package for Nonlinear Programming. Stanford University Department of Operations Research, January 1986.
- Goodsell, D.S., Olson, A.J. Automated docking of substrates to proteins by simulated annealing. *Proteins* 8:195-202, 1990.
- Hammer, J., Bolin, C., Papadopoulos, D., Walsky, J., Higelin, J., Danho, W., Sinigaglia, F., Nagy, Z.A. High-affinity binding of short peptides to major histocompatibility complex class II molecules by anchor combinations. *Proc. Natl. Acad. Sci. U.S.A.* 91:4456, 1994.
- Hammer, J., Bono, E., Gallazzi, F., Belunis, C., Nagy, Z., Sinigaglia, F. Precise prediction of major histocompatibility complex class II-peptide interaction based on peptide side chain scanning. *J. Exp. Med.* 180:2353-2358, 1994.
- Hammer, J., Gallazzi, F., Bono, E., Karr, R., Guenot, J., Valsasini, Nagy, Z. Peptide binding specificity of HLA-DR4 molecules: correlation with Rheumatoid Arthritis Association. *J. Exp. Med.* 181:1847-1855, 1995.
- Hart, T.N., Read, R.J. A multiple-start Monte-Carlo docking method. *Proteins* 13:206-222, 1992.
- Kuntz, I.D., Blaney, J.M., Oatley, S.J., Langridge, R., Ferrin, T.E. A geometric approach to macromolecule-ligand interactions. *J. Mol. Biol.* 161:269-288, 1982.
- Jardesky, T.S., Gorga, J.C., Bush, R., Rothbard, J., Strominger, J.L. Wiley, D.C. Peptide binding to HLA-DR1: a peptide with most residues substituted to alanine retains MHC binding. *EMBO J.* 9:1797, 1990.
- Jiang, F., Kim, S.H. Soft docking: Matching of molecular surface cubes. *J. Mol. Biol.* 219:79-102, 1991.
- Nauss, J.L., Reid, R.H., Sadegh-Nasseri, S. Accuracy of a structural homology model for a class II histocompatibility protein, HLA-DR1: comparison to the crystal structure. *J. Biomol. Struct. Dyn.* 12:1213-1233, 1995.
- Gulukota, K., Vajda, S., Delisi, C. Peptide docking using dynamic programming. *J. Comp. Chem.* 17:418-428, 1996.
- Lee, B., Richards, F.M. The interpretation of protein structures: estimation of static accessibility. *J. Mol. Biol.* 55:379-400, 1971.
- Levitt, M. Protein folding by restrained energy minimization and molecular dynamics. *J. Mol. Biol.* 170:723, 1983.
- Maranas, C.D., Androulakis, I.P., Floudas, C.A. A determin-

- istic global optimization approach for the protein folding problem. DIMACS Ser. Discrete Math. Theor. Comput. Sci. 23:133–150, 1995.
25. Maranas, C.D., Floudas, C.A. Global minimum potential energy conformations of small molecules. *J. Global Optimiz.* 4:135–170, 1994.
  26. Momany, F.A., Burgess, A.W., McGuire, R.F., Scheraga, H.A. Energy Parameters in polypeptides. VII. Geometric parameters, partial atomic charges, nonbonded interactions, hydrogen bond interactions, and intrinsic torsional potentials for the naturally occurring amino acids. *J. Phys. Chem.* 79:2361, 1975.
  27. Momany, F.A., Carruthers, L.M., McGuire, R.F., Scheraga, H.A. Energy parameters in polypeptides. *J. Phys. Chem.* 78:1595, 1974.
  28. Monos, D., Soulika, A., Argyris, E., Gorga, J., Stern, L., Magafa, V., Cordopatis, P., Androulakis, I.P., Floudas, C.A. HLA diversity functional and medical implications. *Proc. Int. Histocomp. World Conf.* 12:xx–xx, 1997.
  29. Némethy, G., Gibson, K.D., Palmer, K.A., Yoon, C.N., Paterlini, G., Zagari, A., Rumsey, S., Scheraga, H.A. Energy parameters in polypeptides. 10. Improved geometrical parameters and nonbonded interactions to use in the ECEPP/3 algorithms, with applications to proline-containing peptides. *J. Phys. Chem.* 96:6472, 1992.
  30. Némethy, G., Pottle, M.S., Scheraga, H.A. Energy parameters in polypeptides. 9. Updating of geometrical parameters, nonbonded interactions and hydrogen bond interactions for the naturally occurring amino acids. *J. Phys. Chem.* 89:1883, 1983.
  31. Nicholls, A. GRASP: Graphical Representation and Analysis of Surface Properties, October 1992.
  32. Paul, W. E. "Fundamental Immunology." Philadelphia: Lippincott-Raven, 1993.
  33. Perrot, G., Cheng, B., Gibson, K.D., Palmer, K.A., Vila, J., Nayeem, A., Maigret, B., Scheraga, H.A. MSEED: a program for the rapid analytical determination of accessible surface areas and their derivatives. *J. Comput. Chem.* 13:1–11, 1992.
  34. Rogman, D., Scapozza, L., Folkers, G., Daser, A. Molecular Dynamics simulation of MHC-peptide complexes as a tool for predicting potential T cell epitopes. *Biochemistry* 33: 11476–11485, 1994.
  35. Rosenfeld, R., Q. Zheng, Vajda, S., DeLisi, C. Computing the structure of bound peptides: applications to antigen recognition by class I major histocompatibility complex receptors. *J. Mol. Biol.* 234:515–521, 1993.
  36. Scheraga, H.A. ECEPP/3 USER GUIDE. Cornell University Department of Chemistry, January 1993.
  37. Scheraga, H.A. PACK: Programs for Packing Polypeptide Chains, (online documentation), 1996.
  38. Stern, L., Brown, J., Jardetzky, T., Gorga, J., Urban, R., Strominger, L., Wiley, D. Crystal structure of the human class II MHC protein HLA-DR1 complexed with an influenza virus peptide. *Nature* 368:215–221, 1994.
  39. van Gunsteren, W.F., Berendsen, H.J.C. "GROMOS: Groningen Molecular Simulation." Groningen: The Netherlands, 1987.
  40. Vila, J., Williams, R.L., Vasquez, M., Scheraga, H.A. Empirical solvation models can be used to differentiate native from non-native conformations of bovine pancreatic trypsin inhibitor. *Proteins* 10:199–218, 1991.
  41. Weiner, S., Kollmann, P., Case, D.A., Singh, U.C., Ghio, C., Alagona, G., Profeta, S., Weiner, P. A new force field for molecular mechanical simulation of nucleic acids and proteins. *J. Am. Chem. Soc.* 106:765, 1984.
  42. Weiner, S., Kollmann, P., Nguyen, D., Case, D. An all-atom force field for simulations of proteins and nucleic acids. *J. Comp. Chem.* 7:230, 1986.

## APPENDIX A

### Determination of Euler Angles

The calculation of the Euler angles is complicated by the fact that the location of the hydrogen atom in the  $(x, y, z)$  space is not known. For this reason, a

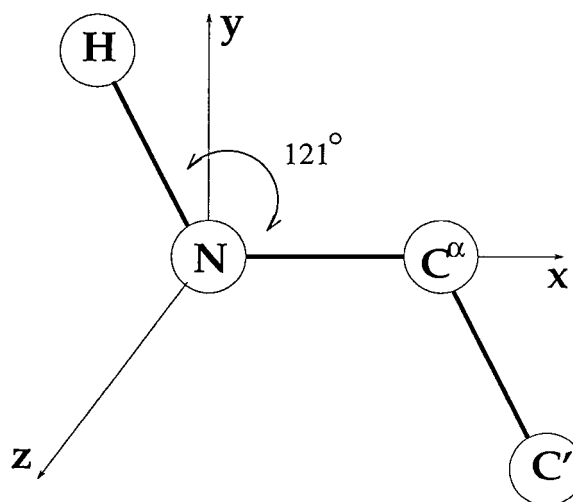


Fig. 18. Axes defined by the relative positions of C', N, and C<sup>α</sup>.

discussion of the method of hydrogen position determination will precede the Euler angle theory.

### Determination of Hydrogen Location

The basic steps behind finding the location of a hydrogen is the following: (1) define a basis system, (2) find the relative position of the hydrogen in this system, and (3) translate and rotate this position to a new basis defined by a particular molecule. This will become clear as the explanation proceeds.

The diagram in Figure 18 shows how the axes are defined in relation to the prime carbon, nitrogen and alpha carbon. The phi angle is always approximately equal to 180°. Hence the N—C<sup>α</sup> bond defines the x-axis, while the C<sup>α</sup>—C' and N—H bonds lie in the XY plane.

The position of the hydrogen is initially unknown. The position of the nitrogen is taken to be at the origin. The position of the alpha carbon is known because the N—C bond length is approximately 1.435 Å. By this definition, the positions of the atoms are N(0, 0, 0), and C<sup>α</sup>(1.435, 0, 0). The position of the hydrogen can be found by the knowledge that the H—N—C<sup>α</sup> bond angle is 121°. This angle is illustrated by the arrow in Figure 18. Since the hydrogen lies in the x, y plane then the hydrogen will lie in the quadrant defined by negative x and positive y. Hence the position is easily found remembering that the N—H bond length is ~1.0 Å. The x and y coordinates are just the cosine and sine of the angle of 121°, respectively. Solving this system yields: H(−0.5150, 0.8570, 0), or since nitrogen is defined as the origin in this basis, the following vector is defined:

$$\mathbf{NH} = -0.4226\mathbf{i} + 0.3791\mathbf{j} + -0.8232\mathbf{k} \quad (17)$$

where  $\mathbf{i}$ ,  $\mathbf{j}$ , and  $\mathbf{k}$  are the defined unit vectors for the basis.

Now if N, C', and C $\alpha$  are at a random orientation in space, the magnitudes (or distances from the defined origin, nitrogen) do not change, but the directions of the unit normals do change. Hence if  $\mathbf{i}$ ,  $\mathbf{j}$ , and  $\mathbf{k}$  are defined in the new orientation the vector will still be defined as above, but will have a different value due to the change in the unit vectors.

The final step is to define these unit vectors as the basis of the orientation of a particular molecule, where the coordinates of N, C', and C $\alpha$  are known. The vectors NC $\alpha$ , and C $\alpha$ C' are easily defined by subtraction of coordinates. Then the unit vectors in the axis directions can be defined as follows.

In the  $x$  direction,

$$\mathbf{i} = \frac{\mathbf{NC}^\alpha}{|\mathbf{NC}^\alpha|} \quad (18)$$

In the  $z$  direction,

$$\mathbf{z} = (\mathbf{C}^\alpha\mathbf{C}') \times (\mathbf{NC}^\alpha) \quad (19)$$

$$\mathbf{k} = \frac{\mathbf{z}}{|\mathbf{z}|} \quad (20)$$

In the  $y$  direction,

$$\mathbf{j} = \mathbf{k} \times \mathbf{i} \quad (21)$$

Then for the particular orientation:

$$\mathbf{NH} = -0.5150\mathbf{i} + 0.8570\mathbf{j} + 0.00\mathbf{k} \quad (22)$$

Finally, the exact coordinates of hydrogen can be determined by adding  $\mathbf{NH}$  to the given coordinates of nitrogen for the system.

### The Euler Angle Theory

The Euler angles are found by comparing the angles between the unit normals defining the coordinate axes. The basis coordinate system is defined as follows and is subscripted with a 1 when mentioned later: N(0, 0, 0), C $\alpha$ (1.435, 0, 0), H(-0.515, 0.857, 0).

So initially for given coordinates of the H, N, and C $\alpha$  atoms, the unit vectors defining the coordinate axes must be found. First, the vectors  $\mathbf{NC}^\alpha$  and  $\mathbf{NH}$  are found by subtracting the respective coordinates. These vectors define the axes as shown in Figure 18. Once again the N—C $\alpha$  bond lies on the  $x$ -axis, and the N—H bond is defined as lying in the  $xy$  plane. This orientation is shown on the coordinate axes in Figure 18. The unit vectors on the axes are easily described by the following equations.

In the  $x$  direction,

$$\mathbf{i} = \frac{\mathbf{NC}^\alpha}{|\mathbf{NC}^\alpha|} \quad (23)$$

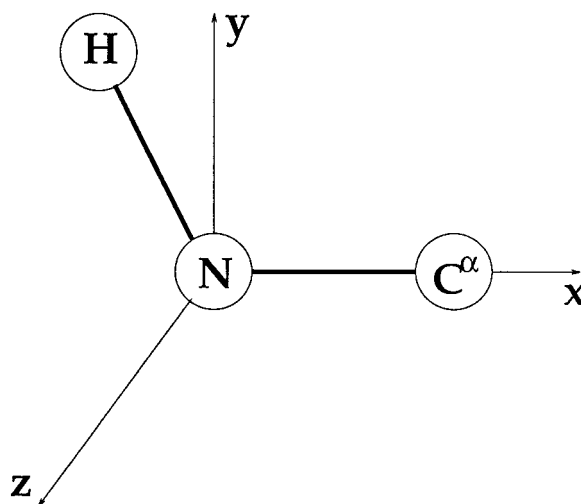


Fig. 19. Axes defined by the relative positions of H, N, and C $\alpha$ .

In the  $z$  direction,

$$\mathbf{k} = \frac{\mathbf{n}}{|\mathbf{n}|} \quad (24)$$

where

$$\mathbf{n} = \mathbf{i} \times (\mathbf{NH}) \quad (25)$$

In the  $y$  direction,

$$\mathbf{j} = \mathbf{k} \times \mathbf{i} \quad (26)$$

After performing the above operations on the initial basis and a residue there are two coordinate systems to compare. By definition the following relations hold for the Euler angles in the basis system defined above:

$$\sin \theta_1 = -\mathbf{i}_1 \cdot \mathbf{k}_2 \quad (27)$$

$$\cos \theta_1 = -\mathbf{j}_1 \cdot \mathbf{k}_2 \quad (28)$$

$$\cos \theta_2 = \mathbf{k}_1 \cdot \mathbf{k}_2 \quad (29)$$

$$\sin \theta_3 = -\mathbf{k}_1 \cdot \mathbf{i}_2 \quad (30)$$

$$\cos \theta_3 = \mathbf{k}_1 \cdot \mathbf{j}_2 \quad (31)$$

$$\sin^2 \theta_2 = 1 - \cos^2 \theta_2 \quad (32)$$

The appropriate euler angles ( $\theta_1$ ,  $\theta_2$ ,  $\theta_3$ ) can easily be found by taking the arctangent of the ratio of sine to cosine. The signs on the Euler angles are determined by using the signs of the sine and cosine of the angle to determine the quadrant where the angle is defined.

## APPENDIX B

### Determination of C' Location

As mentioned in the subsection Problem Formulation the coordinates of the backbone carboxyl carbon have to be expressed as a function of other variables since they do not correspond to explicit optimization variables. This can be made using analytical geometry and linear algebra. In particular, for each amino acid the coordinates of C' atom can be expressed as a function of the coordinates of the N atom (i.e., the translation vector) denoted as  $n_x, n_y, n_z$ , the Euler angles  $\theta_1, \theta_2, \theta_3$  and the  $\phi$  dihedral angle, and requires the knowledge of the following parameters the bond lengths  $NC^\alpha$  and  $C^\alpha C'$  as well as the value of the angle  $NC^\alpha C'$  taken from the literature.<sup>27</sup>

Given the above information and based on the graphic representation of protein shown in Figure 7, the following expressions are derived for the coordinates of C' atom:

$$\begin{aligned} x = & n_x + sc \times temp \times \sin(\theta_1) \times \sin(\theta_2) \\ & + sa \times sc \times \cos(\phi) \times temp \\ & \times \frac{[-\cos(\theta_2) \times \cos(\theta_3) \times \sin(\theta_1) - \cos(\theta_1) \times \sin(\theta_3)]}{\sqrt{[1 - \cos^2(\phi)]}} \\ & + [(NC^\alpha) + (C^\alpha C') \times \cos(\theta)] \times (\cos(\theta_1) \times \cos(\theta_3) \\ & - \cos(\theta_2) \times \sin(\theta_1) \times \sin(\theta_3)) \end{aligned}$$

$$\begin{aligned} y = & n_y - sc \times temp \times \cos(\theta_1) \times \sin(\theta_2) \\ & + sa \times sc \times \cos(\phi) \times temp \\ & \times \frac{[\cos(\theta_1) \times \cos(\theta_2) \times \cos(\theta_3) - \sin(\theta_1) \times \sin(\theta_3)]}{\sqrt{[1 - \cos^2(\phi)]}} \\ & + [(NC^\alpha) + (C^\alpha C') \times \cos(\theta)] \\ & \times (\sin(\theta_1) \times \cos(\theta_3) + \cos(\theta_1) \times \cos(\theta_2) \times \sin(\theta_3)) \\ z = & n_z + sc \times temp \times \cos(\theta_2) + sa \times sc \times \cos(\phi) \\ & \times temp \times \cos(\theta_3) \times \frac{\sin(\theta_3)}{\sqrt{[1 - \cos^2(\phi)]}} \\ & + [(NC^\alpha) + (C^\alpha C') \times \cos(\theta)] \times (\sin(\theta_2) \times \sin(\theta_3)) \end{aligned}$$

where

$$temp = \sqrt{\frac{(C^\alpha C')^2 - (C^\alpha C')^2 \times \cos^2(\theta)}{1 + \frac{\cos^2(\phi)}{1 - \cos^2(\phi)}}$$

$$sa = 1.0, sc = -1.0 \quad \text{if } \phi > 0$$

$$sa = -1.0, sc = 1.0 \quad \text{if } \phi \leq 0$$

and  $\theta$  is the angle  $NC^\alpha C'$ .