

Rochester Institute of Technology

**RIT Digital Institutional Repository**

---

Theses

---

8-15-2006

**Structural analysis of the EGR family of transcription factors:  
Templates for predicting protein - DNA interactions**

Jamie Duke

Follow this and additional works at: <https://repository.rit.edu/theses>

---

**Recommended Citation**

Duke, Jamie, "Structural analysis of the EGR family of transcription factors: Templates for predicting protein - DNA interactions" (2006). Thesis. Rochester Institute of Technology. Accessed from

This Thesis is brought to you for free and open access by the RIT Libraries. For more information, please contact [repository@rit.edu](mailto:repository@rit.edu).

Structural Analysis of the EGR Family of Transcription Factors:  
Templates for Predicting Protein – DNA Interactions

Approved:



Gary R. Skuse, Ph.D.  
Thesis Advisor



Richard L Doolittle, Ph.D  
Head, Department of Biological Sciences

Submitted in partial fulfillment of the requirements for the Master of Science  
degree in Bioinformatics at the Rochester Institute of Technology.

Jamie L. Duke  
May 2006

## Thesis/Dissertation Author Permission Statement

Title of thesis or dissertation: Structural Analysis of the EGR Family of Transcription Factors:  
Templates for Predicting Protein – DNA Interactions.

Name of Author: Jamie L. Duke  
Degree: Master's of Science  
Program: Bioinformatics  
College: Science

I understand that I must submit a print copy of my thesis or dissertation to the RIT Archives, per current RIT guidelines for the completion of my degree. I hereby grant to the Rochester Institute of Technology and its agents the non-exclusive license to archive and make accessible my thesis or dissertation in whole or in part in all forms of media in perpetuity. I retain all other ownership rights to the copyright of the thesis or dissertation. I also retain the right to use in future works (such as articles or books) all parts of this thesis or dissertation.

### *Print Reproduction Permission Granted:*

I, Jamie L. Duke, hereby grant permission to the Rochester Institute of Technology to reproduce my print thesis or dissertation in whole or in part. Any reproduction will not be for commercial use or profit.

Signature of Author: Jamie L. Duke

Date: 20 June 2006

### *Inclusion in the RIT Digital Media Library Electronic Thesis & Dissertation (ETD) Archive*

I, Jamie L. Duke, additionally grant to the Rochester Institute of Technology Digital Media Library (RIT DML) the non-exclusive license to archive and provide electronic access to my thesis or dissertation in whole or in part in all forms of media in perpetuity.

I understand that my work, in addition to its bibliographic record and abstract, will be available to the world-wide community of scholars and researchers through the RIT DML. I retain all other ownership rights to the copyright of the thesis or dissertation. I also retain the right to use in future works (such as articles or books) all or part of this thesis or dissertation. I am aware that the Rochester Institute of Technology does not require registration of copyright for ETDs.

I hereby certify that, if appropriate, I have obtained and attached written permission statements from the owners of each third party copyrighted matter to be included in my thesis or dissertation. I certify that the version I submitted is the same as that approved by my committee.

Signature of Author: Jamie L. Duke

Date: 20 June 2006

**Thesis Advisor**

Dr. Carlos J. Camacho  
Department of Computational Biology  
University of Pittsburgh

**Thesis Committee**

Dr. Gary Skuse  
Department of Biological Sciences  
Rochester Institute of Technology

Dr. James Halavin  
Department of Mathematics & Statistics  
Rochester Institute of Technology

Dr. Paul Craig  
Department of Chemistry  
Rochester Institute of Technology

## **Abstract**

The EGR family of transcription factors is known to be activated in cells exposed to growth factors in a variety of tissues. The overall structure of the family is highly conserved while the amino acid sequence can be quite diverse allowing for a wide array of DNA recognition sequences. Through homology modeling it is possible to reproduce the structure of the DNA binding domain of EGR proteins, which consists of three zinc fingers. It has also been determined through molecular dynamic simulations that most side chains within the domain reach an equilibrium state. However, residues that are essential for DNA binding are seen throughout the simulation as not reaching an equilibrium state, but constantly sampling available conformational space. Furthermore, through cluster analysis the three recognition residues in each zinc finger are found to have side chain conformations that are optimal for DNA recognition. These studies help to show a possible mechanism for zinc finger recognition of DNA and create homology modeled proteins that are able to be used in protein – DNA interaction prediction.

## List of Figures

<b>Figure 1.</b>	Patterns of zf-C2H2 domain.	Page 2
<b>Figure 2.</b>	Representative Structure of the EGR Family of Transcription Factors.	Page 3
<b>Figure 3.</b>	EGR Protein Binding Motif	Page 4
<b>Figure 4.</b>	Flowchart of Methods.	Page 6
<b>Figure 5.</b>	Figures of Isoleucine and Alanine.	Page 7
<b>Figure 6.</b>	Alpha Carbon Aligning.	Page 9
<b>Figure 7.</b>	Sample Alignments from the Multiple Sequence Alignment of EGR Protein Domains.	Page 12
<b>Figure 8.</b>	Full Alignments of EGR Proteins of Interest.	Page 13
<b>Figure 9.</b>	Differences between Homology Models and Crystal Structures.	Page 15
<b>Figure 10.</b>	Hydrogen Bonding in Initial Molecular Dynamic Simulation Data Clustering	Page 17
<b>Figure 11.</b>	Graphs of Domain 1 of GD-AA.	Page 19
<b>Figure 12.</b>	Graphs of Domain 2 of GD-AA.	Page 20
<b>Figure 13.</b>	Graphs of Domain 3 of GD-AA.	Page 21
<b>Figure 14.</b>	Graphs of Domain 1 of AA-GD.	Page 23
<b>Figure 15.</b>	Graphs of Domain 2 of AA-GD.	Page 24
<b>Figure 16.</b>	Graphs of Domain 3 of AA-GD.	Page 25
<b>Figure 17.</b>	Graphs of Domain 1 of ME-AA.	Page 27
<b>Figure 18.</b>	Graphs of Domain 2 of ME-AA.	Page 28
<b>Figure 19.</b>	Graphs of Domain 3 of ME-AA.	Page 29
<b>Figure 20.</b>	GD-AA Clustered Domain Results.	Page 31
<b>Figure 21.</b>	AA-GD Clustered Domain Results.	Page 32
<b>Figure 22.</b>	ME-AA Clustered Domain Results.	Page 33
<b>Figure 23.</b>	Potential Protein – DNA Interactions for GD-AA.	Page 38

<b>Figure 24.</b>	Potential Protein – DNA Interactions for AA-GD.	Page 39
<b>Figure 25.</b>	Favorable Coordination of a Guanine by Residues in Two Domains.	Page 40
<b>Figure 26.</b>	Potential Protein – DNA Interactions for ME-AA.	Page 41
<b>Figure 27.</b>	Favorable and Unfavorable DNA Interactions within ME-AA.	Page 42

### List of Tables

<b>Table 1.</b>	Residues of interest for GD-AA.	Page 34
<b>Table 2.</b>	Residues of interest for AA-GD.	Page 35
<b>Table 3.</b>	Residues of interest for ME-AA.	Page 36
<b>Table 4.</b>	RMS values between completed models and the crystal structure.	Page 37

### List of Equations

<b>Equation 1.</b>	Distance Equation	Page 9
<b>Equation 2.</b>	Root Mean Square Equation	Page 9

### List of Appendices

<b>Appendix A.</b>	Multiple Sequence Alignment of EGR Protein Domains	Page 50
<b>Appendix B.</b>	Consensus Output	Page 63
<b>Appendix C.</b>	Domain Information for Each Homology Modeled EGR Protein	Page 69
<b>Appendix D.</b>	RMS versus Crystal Structure for Clustered Models	Page 70

## **Acknowledgements**

I would like to thank Dr. Carlos Camacho for allowing me to start this project during the summer of 2005 with Bioengineering and Bioinformatics Summer Institute (BBSI) hosted by the Department of Computational Biology at University of Pittsburgh, and for allowing me to continue this project at RIT. Without the NIH and NSF grant through the BBSI at the University of Pittsburgh, I would not have been able to start doing the research that ultimately led to the completion of my thesis. I would also like to thank Dr. Gary Skuse, Dr. Paul Craig and Dr. James Halavin for serving on my thesis committee and helping me through the thesis process. Finally, I would like to thank my parents, Candy and Russ, as well as my sister, Danielle, for all of their encouragement and support while continuing my education.

## Table of Contents

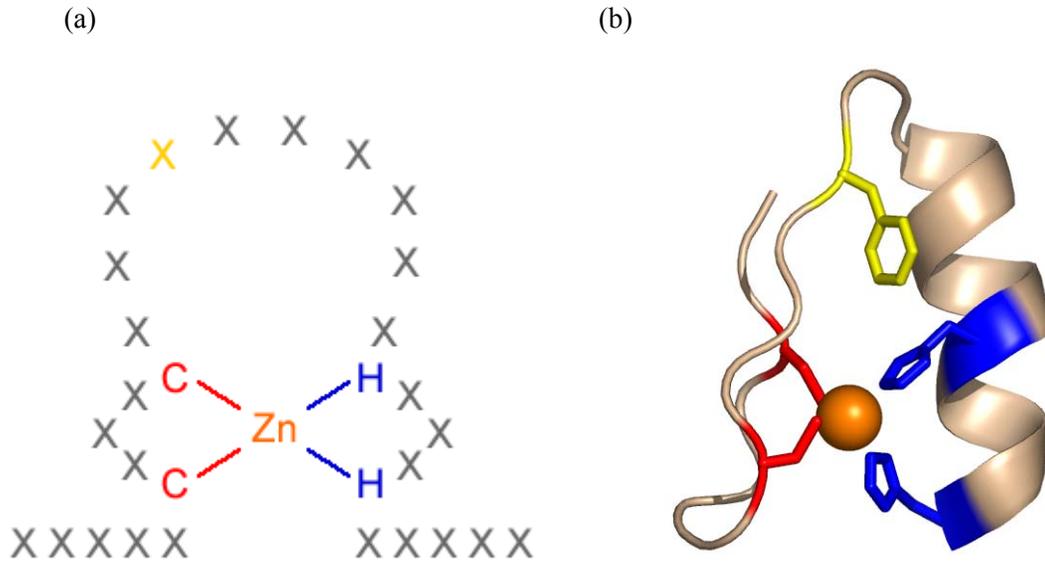
Copyright Release Form	Page ii
Thesis Committee Information	Page iii
Abstract	Page iv
List of Figures	Page v
List of Tables	Page vi
List of Equations	Page vi
Acknowledgements	Page vii
Introduction	Page 1
Methods	Page 5
Results	Page 11
Discussion	Page 37
Conclusion	Page 46
References	Page 47
Appendix	Page 49

## **Introduction**

Understanding interactions between proteins and DNA has become a major area of research in biology within the past few years. Furthermore, the challenges of predicting these interactions have lured both computational biologists and bioinformaticists into joining the expedition. By understanding the general mechanism for protein – DNA binding, it is possible to further understand cell regulation and develop more specific pharmaceuticals for existing diseases.

Focusing on specific families of proteins has allowed for advancements in the field, specifically with the early growth response factor (EGR) family. This family of transcription factors has proved useful due to the fact that it is highly conserved and extensively studied. It has been implicated in a diverse set of cellular processes including cell growth, differentiation, and apoptosis (Gashler and Sukhatme 1995). The EGR family of transcription factors is within the class of immediate early genes (IEGs) which are the group of genes directly activated by extracellular signals. Immediate early genes have been referred to as the “gateway to the genomic response” as they control the regulation of downstream genes, with two subclasses: effector and regulatory genes (Davis et al. 2003). EGR proteins respond to mitogens, differentiation stimuli, tissue injury and signals from neuron excitation (Gashler and Sukhatme 1995). Mitogens are the dominant form of extracellular signals for EGR protein induction and also cause the quickest reaction in EGR production.

All members of the EGR family contain two or three zinc finger domains which are used in conjunction with one another to recognize specific DNA sequences (Benos et al. 2002). Each zinc finger recognizes approximately 5 bases within the DNA strand, and has an overlap in the recognition of bases due to coordination between the other fingers in the protein (Palliard et al. 2004). Christy and Nathans first described EGR proteins in 1989 as containing three repeated zinc finger domains, and they hypothesized that collectively these domains worked as transcription factors to recognize a specific DNA sequence. The gene they describe, *zif268*, is



**Figure 1. Patterns of zf-C2H2 domain.** (a) The two dimensional representation of the zf-C2H2 domain. The four conserved residues that coordinate the zinc ion are shown in the center of the domain. The yellow ‘X’ four residues away from the second cysteine residue is also a conserved aromatic ring essential to the stability of the domain. Adapted from Prosite figure. (b) Three dimensional representation of the zinc finger C2H2 domain. The amino acid and zinc coloration is the same as in (a).

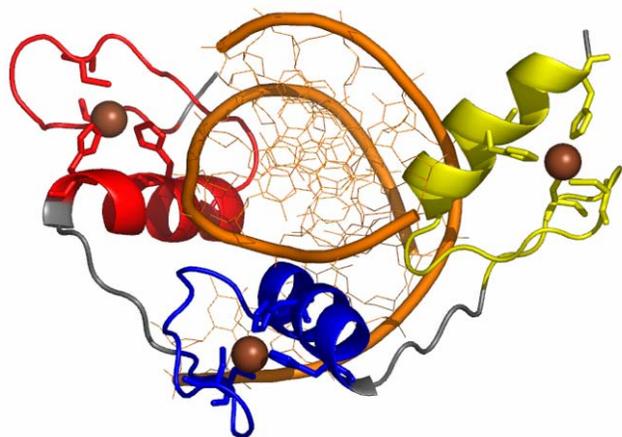
widely expressed in mouse tissues and is also known by the names *NGF-1A*, *egr-1*, and *Krox-24*.

The consensus DNA sequence for zinc finger binding of *zif268* is  $GCG \frac{G}{T} GGGCG$ . A high affinity for guanine in DNA recognition is specific to this family.

Each zinc finger in the EGR family of transcription factors is of the type zf-C2H2, where ‘zf’ stands for zinc finger and C2H2 represents the 2 conserved cysteine and 2 conserved histidine residues. The zf-C2H2 domain has a consensus pattern of: ‘x-C-x(1-5)-C-x(12)-H-x(3-6)-H’ (Hulo et al. 2004). The ‘X’ in the pattern stands for any amino acid with the amount of repeating units in parenthesis. Structural stability in the domain is provided by the four conserved residues working together to coordinate a zinc ion along with a conserved aromatic ring located 4 amino acids away from the second cysteine residue. Figure 1 shows the zf-C2H2 domain in both a two- and three-dimensional manner, as well as the details of how the four conserved residues coordinate the zinc ion. Each zinc finger domain includes one  $\alpha$ -helix and two  $\beta$ -strands. The  $\alpha$ -helix is inserted into the major groove of the DNA double helix and is responsible for recognizing

at most 5 nucleic acids. The two conserved histidine residues responsible for coordinating the zinc ion are located towards the C-terminal of the  $\alpha$ -helix pointing away from the DNA, as can be seen in Figure 2.

The EGR family of proteins has a set method of DNA recognition that is coordinated between the three zinc fingers, with each finger specifically recognizing three nucleotides. Each domain has three key residues for interaction with DNA, the negative one position, the third position and the sixth position in the  $\alpha$ -helix, and can be seen in Figure 3. DNA recognition is specific to the residues that are located within the three key residues, and the given recognition sequence can be changed by mutating the protein binding region (Palliard et al. 2004). Secondary



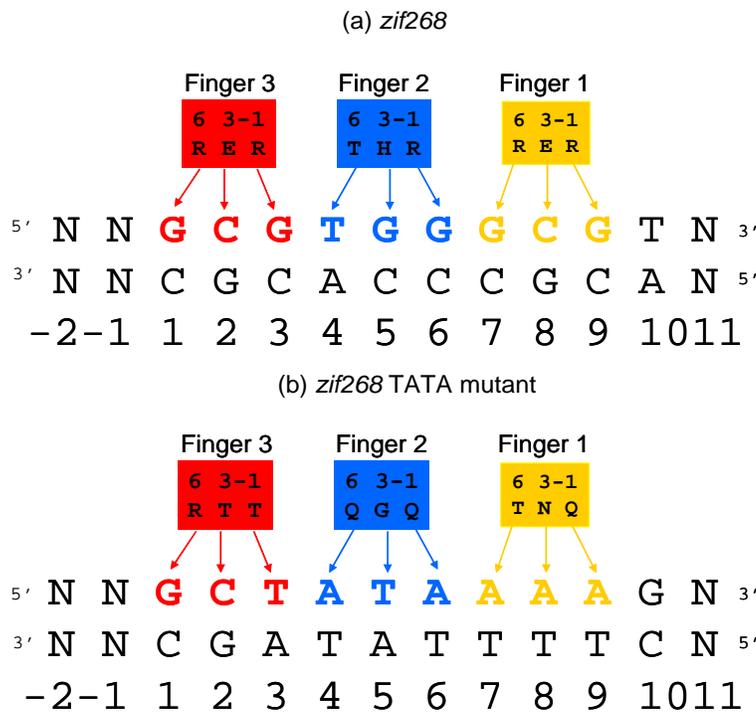
**Figure 2. Representative Structure of the EGR Family of Transcription Factors.** Seen here is the EGR protein from *Mus musculus* (PDB identifier: 1A1F). The DNA structure is shown in orange, with the backbone outlined. Each of the zinc finger domains is highlighted, starting with the N-terminus domain in yellow, followed by the second domain in blue, and the C-terminus domain in red. The four conserved amino acids, two cysteine and two histidine, can be seen coordinating the space-filled zinc atom. The conserved histidine residues are located within the  $\alpha$ -helix, pointing away from the DNA chain. As can be seen from the figure, each of the zinc finger domains is inserted into the major groove of the DNA. Each domain overlaps the recognized bases with the subsequent domain.

interactions also occur within the binding region through water molecules. These secondary interactions help to overlap the recognition of DNA between the zinc finger domains and involve additional residues located within the  $\alpha$ -helix. When secondary interactions are included, each domain can recognize between four and seven nucleotides.

The proteins that include the classic zinc finger domain are quite diverse. According to Pfam version 17.0 (Bateman et al. 2004), there are a total of 32,784 different protein domains that have been identified with

a zf-C2H2 domain; 1,390 of those sequences are from humans and 1,085 are from *Mus musculus*. The various zf-C2H2 protein sequences are broken up into 235 different architectures, including the ‘zf-C2H2, zf-C2H2, zf-C2H2’ architecture that includes the EGR family of transcription factors. This architecture has 723 members of which 128 are from humans. There are 5 recognized human EGR proteins; however there may potentially be an additional 26 human proteins that could belong to the EGR family, determined by sequence length and proximity of zf-C2H2 domains to one another. The remaining human sequences do not exhibit similarity to the EGR family of proteins.

Although proteins with zinc finger domains are highly studied, there are few structures available relative to the number of sequences that contain the domain. There are a total of 43



**Figure 3. EGR Protein binding motif.** (a) DNA binding motif for *zif268*, the representative for the EGR family. (b) DNA binding motif for a *zif268* variant recognizing the TATA box. The key interactions are shown for each domain with the solid arrows representing hydrogen bonds. Adapted from a figure by Palliard et al. 2004.

structures in the Protein Data Bank (PDB) with resolved zinc finger domains as of October 2005. A preliminary BLASTp search of the PDB yielded a total of 9 to 11 structures that would be eligible for use as a template for structure prediction for the EGR family. Of these possible templates, one is of human origin and the remaining are from *Mus*

*musculus*. Four of the mouse structures are variants of the wild-type EGR family proteins. With the lack of structure data for this large family of proteins, structure prediction is necessary for being able to predict the potential interactions with DNA.

The main goals of this project consist of investigating the diversity of the EGR family in humans including basic bioinformatics analysis, carrying out homology modeling between resolved structures and known human EGR proteins and analyzing the structures to determine if protein – DNA interactions are possible.

## **Methods**

This project has been working towards creating a pseudo-pipeline of programs that is able to take the DNA binding region of an EGR protein and be able to accurately determine a template for homology modeling in hopes of finding the most likely amino acid conformations to determine a structure for use in protein – DNA interaction prediction. While this process is not completely automated, significant progress has been made towards reducing the amount of human intervention in the process. There are numerous points in the following methods which need to be examined for quality purposes. Figure 4 shows the general flow of the methods from an amino acid sequence to a clustered protein structure.

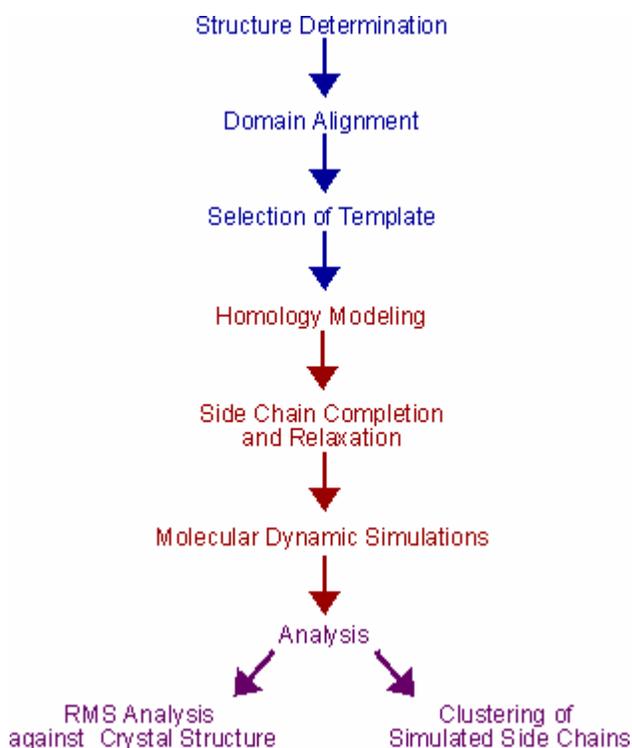
### *EGR Family Analysis*

Before beginning with homology modeling of the EGR family of transcription factors, it was first necessary to gain an understanding of the diversity its members. Members of the family were gathered through the Pfam database, starting with all proteins that have a zinc finger domain of the C2H2 type (Pfam ID: PF00096). It was determined through conversing with Dr. Takis Benos of the University of Pittsburgh that the proteins of interest for the purpose of the project were only the human zinc finger proteins with either two or three zf-C2H2 domains. To gain an understanding of the diversity of the individual domains, the domains were split apart from the

protein and a multiple sequence alignment was performed using ClustalW (Thompson et al. 1994). The alignment was performed using an opening gap penalty of 4.0, a gap extension penalty of 10.0 and the PAM matrix series. The goal of the alignment was to align the conserved cysteine and histidine residues in the zinc finger.

### *Template Determination*

The first step involved in this project was to determine the proper structure(s) to use as a



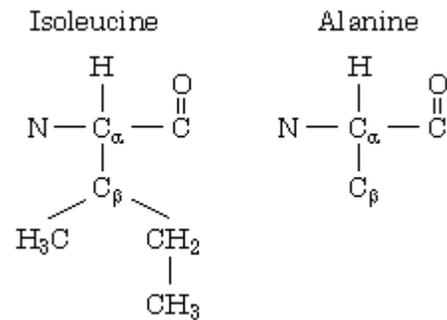
**Figure 4. Flowchart of Methods.** There are three groups of methods used: template determination, structure modeling, and structural analysis. Template determination, shown in blue, is the first three steps and is vital to having an accurate prediction. Structure modeling, shown in red, is the next three steps and is where all of the primary data are generated. Finally, the analysis is shown in purple and is broken into two different parts: analysis against the crystal structure and clustering of the simulated structure.

homology modeling template(s) for a given query sequence. With so few predetermined structures for the EGR family, a combination of multiple template structures may be needed to determine an accurate predicted structure. Basic analysis of EGR sequences were done using the results from BLASTp against the PDB database, and performing multiple sequence alignments to determine the best possible alignment and template (Altschul et al. 1997).

### *Structure Modeling*

After the proper template for the query sequence was determined, the program Consensus was used to create

the homology model (Prasad et al. 2003). Although multiple templates may be beneficial during homology modeling, Consensus can only accommodate one structural template. Instead, Consensus uses five different methods to create an alignment between the template structure and the sequence to be modeled with the structure being determined through threading algorithms. Not only does Consensus perform its own sequence alignment between the template and query, it also predicts the accuracy of the secondary structure assignment for the alignment allowing the user to easily assess the prediction. Although Consensus selectively removes regions of dissimilarity and splits loosely connected domains to minimize potential misalignments, the zinc finger domains of the EGR proteins were not broken apart during the homology modeling process due to the high similarity between the models and the templates.



**Figure 5. Figures of Isoleucine and Alanine.** The structures of these amino acids are only similar through the C- $\beta$  atom, and thus Isoleucine can be modeled only through its C- $\beta$  leaving C- $\gamma$ 1, C- $\gamma$ 2, and C- $\delta$  not in the model.

The structure that was returned from Consensus was incomplete in respect to the side chains of the amino acids. Since threading algorithms are used, side chains for the query sequence can only be predicted in the model as much as they are in agreement with the template. For example, if an isoleucine from the query sequence is aligned with an alanine in the template, the predicted structure for isoleucine would contain the side chain through the C<sub>β</sub> atom, as that is where the structural similarity between the two concludes. Figure 5 further demonstrates this point. Thus, completion of the side chains must be performed to be able to use the structure in further analysis. For this application, side chains were completed and minimized through the use of CHARMM (Brooks et al. 1983). Most importantly, CHARMM added the necessary hydrogen atoms to the amino acids, which are required to be able to do molecular dynamic simulations. The three domains of each EGR protein were manually broken apart for the side chain

completion and molecular dynamic simulations. The independent domains overlap one another by three residues.

The final modeling of the EGR protein involved molecular dynamic (MD) simulations to allow for fluctuation of the individual amino acids. According to Camacho, molecular dynamic simulations can help to predict native-like rotamer conformations of the side chains (2005). Proper prediction of side chain conformation is essential to predict interactions, whether they are between proteins or proteins and DNA. If one side chain is modeled improperly the resulting prediction of interactions could be incorrect.

Simulations are to be performed using GROMACS, in which the protein is simulated using the default settings (Lindahl et al., 2001). The protein is solvated in a dodecahedral water layer where there is a minimum distance from the protein to the edge of the water layer. The size of the water layer was varied, and the optimal water layer was found to be 1.75 angstroms when the simulation commences. Throughout the simulation, the carbonyl carbon, carbonyl oxygen, and nitrogen atoms of the backbone are constrained to keep the protein in the general state that was predicted through homology modeling. The simulation for each zinc finger domain of the protein lasted a total of 4.0 ns with each time step lasting 2 femtoseconds. The first 0.2 ns are removed from the simulation data for equilibrium purposes. Following the simulation, snapshots were taken every 2.5 ps recording structural information to be used for analysis of the simulation.

### *Structure Analysis*

There were two different classes of structural analysis that have taken place: analysis against the crystal structure and analysis of the most highly sampled rotamer conformation. These two different forms of analysis are vital as this method is being verified using known structures for the query protein to have a measurement of the predictive abilities of the homology modeling and the molecular dynamic simulations.

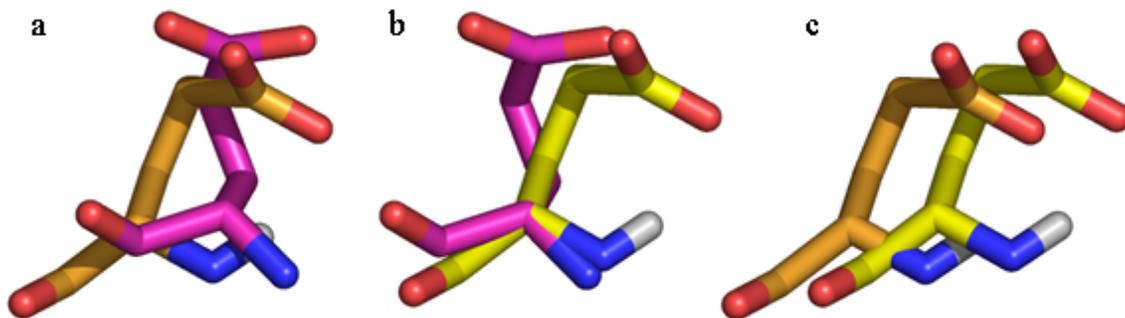
The first form of analysis compared each amino acid for each snapshot of the MD simulation against the crystal structure. Each amino acid was individually analyzed by computing the distance between each heavy atom in the predicted structure and the crystal

$$d = (|x_1 - x_1| + |y_1 - y_2| + |z_1 - z_2|) \quad (\text{Eq. 1})$$

structure, defined as the sum of the absolute value of the differences of each of the three dimensions. Before the distance was calculated, the amino acids were aligned to minimize the fluctuations the protein may have encountered during the MD simulation. To align the amino acids, the distance was calculated for each axis in space between the alpha carbons ( $C_\alpha$ ) and was subtracted from each atom in the predicted structure to align the two amino acids. Figure 6

$$RMS = \sqrt{\frac{1}{N} \sum_{i=1}^N d^2} \quad (\text{Eq. 2})$$

further demonstrates this point visually. The root mean square (RMS) was then calculated, with  $d$  being the distance between  $N$  pairs of structurally equivalent atoms. Additionally, the RMS was computed a second time for arginine residues to account for the symmetry in the terminal atoms. For example, atoms that are labeled the same were compared during the RMS calculation,



**Figure 6. Alpha Carbon Centering.** (a) Pictured are two unaligned glutamate residues, one from the crystal structure shown in magenta, and the other from the MD simulation shown in orange. Computing the RMS between these two residues is going to return a value that is misleading. (b) The two residues are aligned, with the crystal structure still shown in magenta, and the modeled structure is shown in yellow. Here the RMS is going to give much more insight into the difference between the two structures. (c) Shown is the difference between the unaligned and the aligned structure.

but in the case of arginine the two terminal atoms are both nitrogen, and are randomly assigned the designation of 'NH1' and 'NH2'. Thus when the terminal atoms were compared to the crystal structure it was also necessary to compare 'NH1' to 'NH2' and vice versa.

The second form of analysis aimed to complete the structural prediction process by clustering the individual side chains of the modeled protein. A neighbor clustering method was employed to gain general insight into the structure of a given amino acid. Kozakov *et al.* have described a clustering method for amino acid side chains that is applicable to this project (2005). Through clustering, we were looking to determine the most frequently sampled conformation throughout the MD simulation. This procedure relied on the calculation of the RMS for the heavy atoms in a side chain between pairs of time points. The procedure for calculating the RMS in this step of the analysis was the same as comparing the simulated structures to the crystal structure. The first 0.2 ns were dropped from the simulated data before clustering commenced to allow for equilibrium purposes. Approximately  $1.16 * 10^6$  RMS calculations were required to compare each snapshot to every other snapshot for one residue. Each modeled protein has approximately 90 residues, equating to  $1.04 * 10^8$  RMS calculations. Due to the magnitude of the RMS calculations required, three different clustering techniques were attempted throughout the course of this research, one technique used PERL, another used Java, and a third also used PERL but utilizing a different method. Between all three different techniques, conformations were clustered together if the RMS was found to be within a clustering radius of 0.75 angstroms.

The first two methods of clustering are very similar, but were instantiated in different programming languages. The clustering method was first developed using PERL, and is suitable for smaller data sets of approximately 400 models. This program holds every coordinate of every atom for each model throughout the simulation in memory, and then stores the RMS between each residue. Although it was memory intensive, the program was successful in clustering simulations of 1.0 ns in length, however when applied to a 4.0 ns simulation, the program quickly maximized memory usage, and took upwards of 24 to 36 hours to cluster one residue, which

would have equated to approximately three months to cluster one simulated protein. This result was highly unsatisfactory.

Next, a similar clustering algorithm was constructed in Java using the class structure and taking advantage of the hierarchical nature of object oriented programming. The system developed had classes for models within a simulation, residues within a model, and atoms within a residue. There was also a clustering class that applied the previously developed classes. This system was slightly different in the fact that only a subset of residues were accessed and held in memory at one point in time to keep memory usage low. The class method was also not ideal due to the fact that clustering still took approximately a week for one protein.

The third method of analysis developed utilized a program that computes the RMS between two PDB files. In the two previous clustering methods, a majority of the time was spent calculating the RMS, which was not done in an efficient method. By computing the RMS externally from the clustering algorithm and storing the results in a file, the time it takes for computing the RMS and clustering is significantly reduced. Computing the RMS was reduced to approximately one day per domain, and multiple domains were run simultaneously. For each comparison, a RMS file was created, which was accessed in the clustering algorithm. Following calculation of the RMS, the clustering program read in the RMS calculations, and clustered residues that were within  $0.75\text{\AA}$  radius of one another.

## **Results**

### *EGR Family Analysis*

A multiple sequence alignment was performed between all of the human EGR proteins listed in the Pfam database. A total of 719 domains were aligned encompassing a total of 287 EGR proteins. The alignment is a total length of 32 residues, with each domain being 23 residues in length. By having a low gap opening penalty of 4 and a larger gap extension penalty of 10,

CLUSTAL W is able to align three of the four conserved residues: the two cysteine residues and the first histidine. The fourth residue is the final histidine residue in the  $\alpha$ -helix which is noted in the profile of the family as sometimes being replaced by a cysteine, which occurs in 2.2% of the zf-C2H2 domains. These two residues are not seen as easily interchangeable by CLUSTAL W in the alignment procedure as cysteine is polar neutral and histidine is polar basic (Figure 7a). Additionally, the distance between the histidine residues can vary by a difference of four locations at the tail end of the domain. Since there is variability in the location of the histidine, CLUSTAL W prefers to keep the end of the domain intact unless inserting a gap is more favorable than misaligning a residue (Figure 7b). The full alignment can be seen in Appendix A.

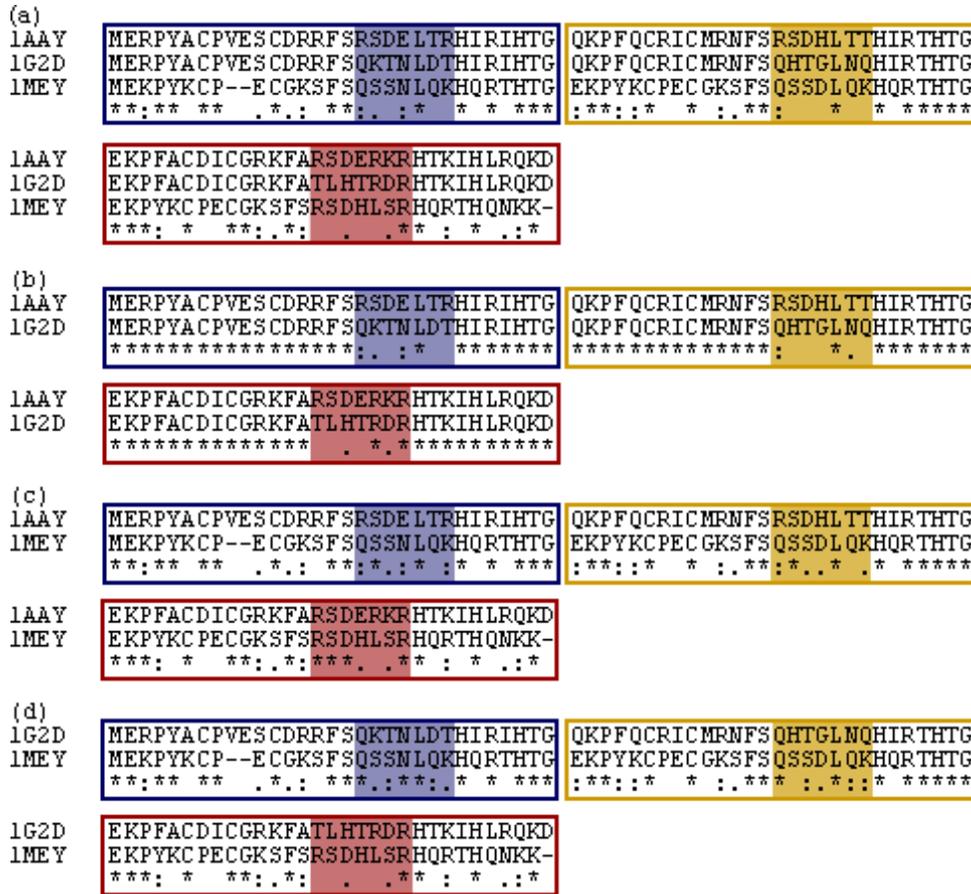
### *Homology Modeling*

The three structures under investigation are a C2H2 zinc finger protein, a C2H2 zinc finger variant, and a C2H2 designed zinc finger. The C2H2 zinc finger protein (PDB 1AAY), was crystallized in 1997 with a 1.60 Å resolution and was chosen because of its sequence complementarities to the EGR family consensus (Elrod-Erickson et al., 1996). The second

```
(a)
Z297B_HUMAN_428-450      YECN--ICAKRFMWRDS--FHRHVTS-C----
Q6ZN18_HUMAN_45-70      YNCCWDQCQACFNSSPD--LADHIRS-IH---
Q6ZN29_HUMAN_3-26       YQCKK--CNVVFPRIFD--LITHQKKQ-C---
                        *      *              :
```

```
(b)
HELI_HUMAN_140-162      FHCNQ--CGASFTQKGN--LLRH--IKLH---
Q6ZMZ8_HUMAN_69-92     FHCDQ--CSYSCKRKDN--LNLH--KCLKH-
FOG2_HUMAN_363-385     FRCN--HCHFQGFQTQRE--LLQHQ--ELH---
Q9C0D4_HUMAN_982-1004  FKCW--FCGRLYEDQEE--WMSHGQR--H---
```

**Figure 7. Sample Alignments from the Multiple Sequence Alignment of EGR Protein Domains.** (a) The misalignment between the final histidine and cysteine residues. Although gaps are inserted into the sequence, if the gaps were removed, the domain would be aligned correctly in this instance. (b) CLUSTAL W prefers to keep the carboxyl terminal of the domain together and places a priority in aligning the residues occurring before the final histidine.



**Figure 8. Full Alignments of EGR proteins of interest.** (a) Multiple sequence alignment of all EGR Proteins in use. Domain 1 is outlined in blue, domain 2 in yellow, and domain 3 in red. The DNA binding region is also highlighted for each domain. (b) Pairwise alignment between 1AA Y and 1G2D. The only difference between these two proteins lies in the DNA binding region. (c) Pairwise alignment between 1AA Y and 1ME Y. While there is more similarity in the DNA binding regions than between 1AA Y and 1G2D at 86% there is only 38% identity. The overall identity and similarity is significantly lower at 49% and 77% respectively. (d) Pairwise alignment between 1G2D and 1ME Y. Overall similarity between 1G2D and 1ME Y is 76%, while the identity is 42%, and the DNA binding domain is 76% similar and 29% identical.

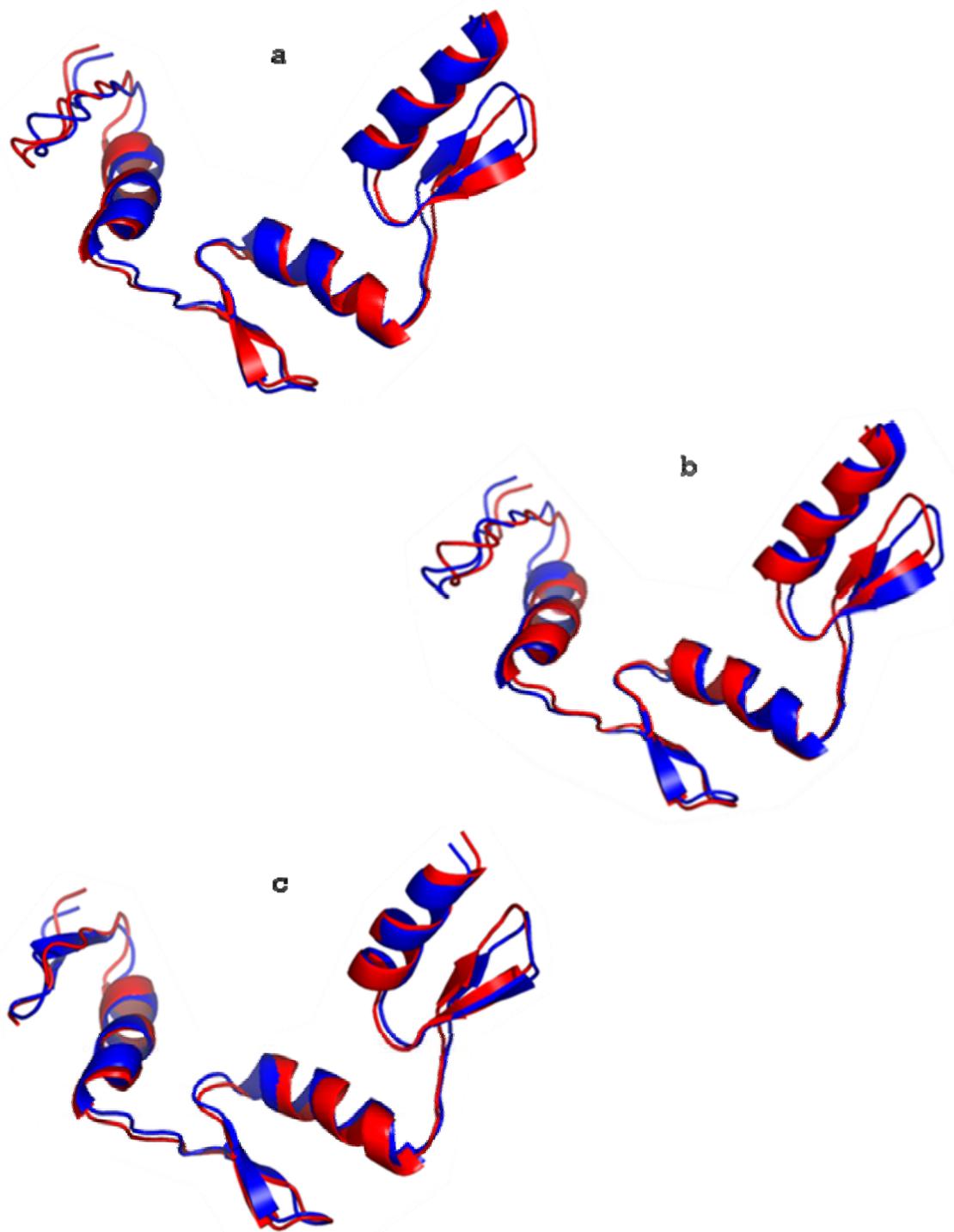
structure, the C2H2 zinc finger variant (PDB 1G2D), was crystallized in 2000 with a 2.20 Å resolution and was chosen for its relatively high sequence identity with 1AA Y and its divergence of the binding domain (Wolf et al., 2001). These proteins have 79% sequence identity, but are only 14% identical in the DNA binding domains. The final structure that was chosen for analysis, the C2H2 designed zinc finger (PDB 1ME Y), was crystallized in 1996 at 2.20 Å resolution (Kim and Berg, 1996). This protein has 49% identity with 1AA Y, but unlike 1G2D,

1MEY has 38% identity to 1AAY in the binding domains. 1MEY was chosen to be modeled off of 1AAY due to the fact that similarity was slightly higher than that of 1MEY with 1G2D. The pairwise and multiple sequence alignments of these domains can be found in Figure 8. The three homology models that were constructed are the prediction of 1G2D with 1AAY as the template (GD-AA), 1AAY with 1G2D as the template (AA-GD), and 1MEY with 1AAY as the template (ME-AA).

Each model utilizes the sequence from an EGR protein with a known structure and goes through the process of homology modeling, side chain completion, molecular dynamic simulations and clustering. Upon completion of homology modeling, Consensus reports the strength of the alignment in terms of the compatibility of the five alignment processes, the secondary structure prediction, as well as two different structure predictions: the full structure prediction and the structure of the reliable regions. For the purposes of this project the full structure prediction was used. Output from Consensus for all three EGR modeled proteins can be found in Appendix B.

For GD-AA, Consensus had a strong alignment between the target and template, but was only confident about the structure prediction for the first domain and the beginning half of the second domain. This result is unexpected as the only regions of nonidentity are in the first seven amino acids of the  $\alpha$ -helix for each domain. It is possible that there is a shift in the secondary structure prediction, which could have caused the lower confidence. However, the RMS between the backbones of GD-AA and 1G2D is 1.29 Å, which confers confidence to the structure prediction.

The structure prediction of AA-GD is more favorable than the prediction of GD-AA. Although the alignments between 1G2D and 1AAY are the same in both homology models, Consensus was more confident in predicting AA-GD and this was seen in the selection status of the output. The only area of difficulty for the prediction lies in the C terminus with the final



**Figure 9. Differences between Homology Models and Crystal Structures.** (a) GD-AA (red) modeled against 1G2D. (b) AA-GD (red) modeled against 1AAY (blue). (c) ME-AA (red) modeled against 1MEY (blue). Major differences in the homology models against the crystal structure is seen more prominently in the first and third domains, leftmost and rightmost domains respectively, particularly in the  $\beta$ -strands. The second domains appear to have the best fit between all three represented figures.

eleven amino acids of the protein. The RMS is also between the backbone of AA-GD and 1AAY is also very favorable at 1.29 Å.

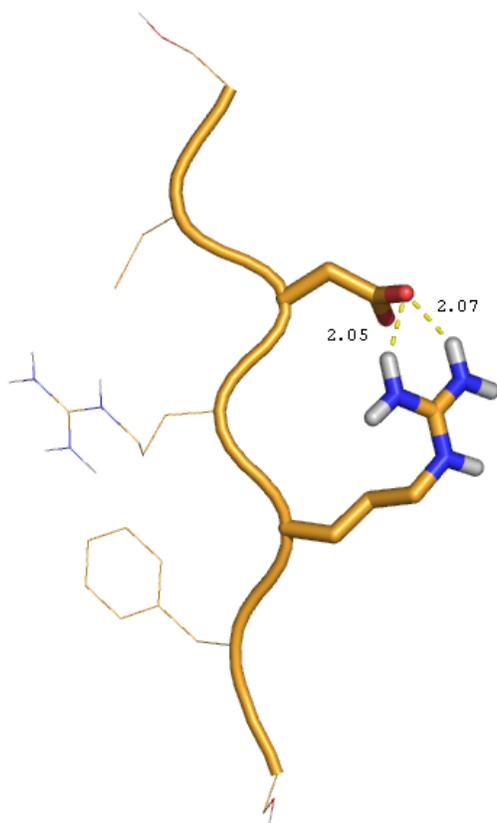
Although there is less similarity between 1AAY and 1MEY than between 1G2D and 1AAY, Consensus has a very strong structural prediction for ME-AA. Much like AA-GD, ME-AA has a strong alignment throughout the whole protein, and Consensus is confident about the structural prediction through most of the model. Unlike AA-GD, ME-AA does not have a confident secondary structure prediction at either the N or C terminals of the protein, lacking confidence for 9 and 10 amino acids respectively. As the confidence is slightly lower for this model than the previous two, the RMS deviation is therefore slightly higher between ME-AA and 1MEY at 1.35 Å. The predicted models and crystal structures can be seen in Figure 9.

The molecular dynamic simulations for each EGR protein zinc finger domain yields structure and energy information for each time step throughout the simulation. To make the data more manageable, snapshots were taken every 2.5 ps creating 1600 total structures for the 4.0 ns simulation and were used for both types of analysis. Initially, the simulations were lacking a layer of water large enough to allow for proper side chain fluctuations, which was made evident through examining the clustering results. The initial water layer was 0.75 Å, which is insufficient for amino acids with larger side chains, for example arginine and lysine. Figure 10 shows potential hydrogen bonding that was occurring in the second zinc finger domain in the GD-AA simulated structure. These results showed neighboring side chains hydrogen bonding to one another, which was unsubstantiated in the literature for the crystal structure. By expanding the water layer the side chains are given more flexibility and less chance to hydrogen bond with neighboring residues.

## Structure Analysis

The structures resulting from the MD simulations were first analyzed against the corresponding crystal structure. For each model there are nine different residues of interest, negative one position, third position and sixth position of the  $\alpha$ -helix.

The first model analyzed was GD-AA against the crystal structure for 1G2D. The three positions in first domain exhibit three different characteristic states as can be seen in Figure 11. The first binding position is the negative one position, or the residue before the  $\alpha$ -helix begins.



**Figure 10. Hydrogen Bonding in Initial Molecular Dynamic Simulation Data Clustering.** The aspartate and arginine residues from the second  $\beta$ -strand, shown in stick representation, have potential hydrogen bonds as they are located within  $2.1\text{\AA}$ . The oxygen atom from the aspartate could hydrogen bond with either  $\text{NH}_2$  group in the arginine.

The glutamine residue has very high fluctuation in the RMS throughout the course of the simulation. There is no pattern within the fluctuations, but the residue does seem to be going between two different states, one state that is close to the crystal structure, approximately  $1.0\text{\AA}$  difference, and another that is much further removed at approximately  $2.5\text{\AA}$ . The third position is occupied with an asparagine which resonates around and equilibrium point of  $1.2\text{\AA}$  difference from the crystal structure. Most residues in the domain are expected to reach some sort of equilibrium state, which residue 21 exemplifies. The sixth position of the helix is inhabited by a threonine. It is interesting that this residue clearly

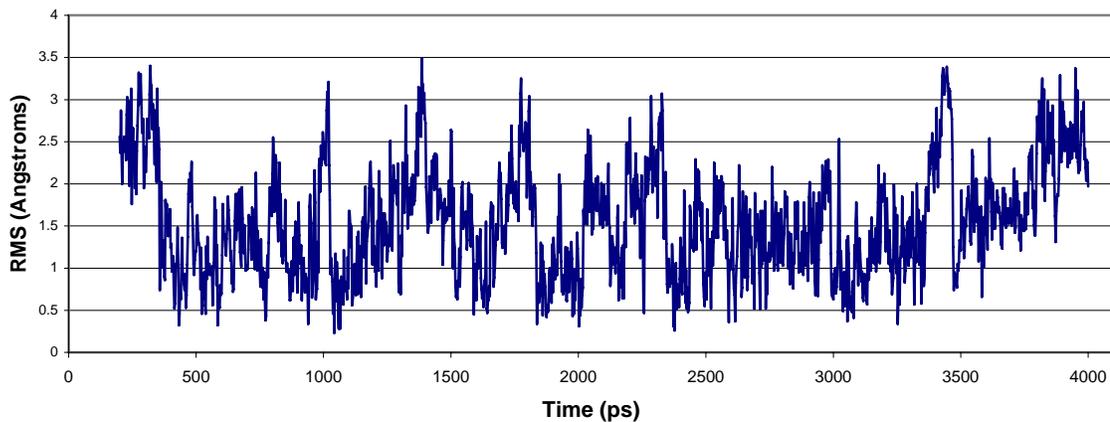
fluctuates between two apparent states, one state existing around 0.3Å difference from the crystal structure and the other state having a 1.1Å difference.

The residues in the second domain are not as interesting as the first domain, seen in Figure 12. Residue 46, a glutamine, in the negative one position in the  $\alpha$ -helix and is very similar to the glutamine in the -1 position of the first domain. There seems to be fluctuation between two states, but it is not as distinct as other residues previously seen. The first state is approximately 1.0Å away from the crystal and the second state is approximately 2.25Å away, with the residue spending a significant amount of time in between the two states. The third position in the helix does not yield much information as the residue is a glycine; the fluctuation seen in Figure 12b is between atoms located in the backbone of the structure. The sixth position of the  $\alpha$ -helix is also a glutamine residue, which has a similar RMS profile as residue 46. The lower state fluctuates just below 1.0Å difference from the crystal structure and the higher state is situated around 2.25Å. Unlike residue 46, residue 52 does not spend much time in an intermediary state.

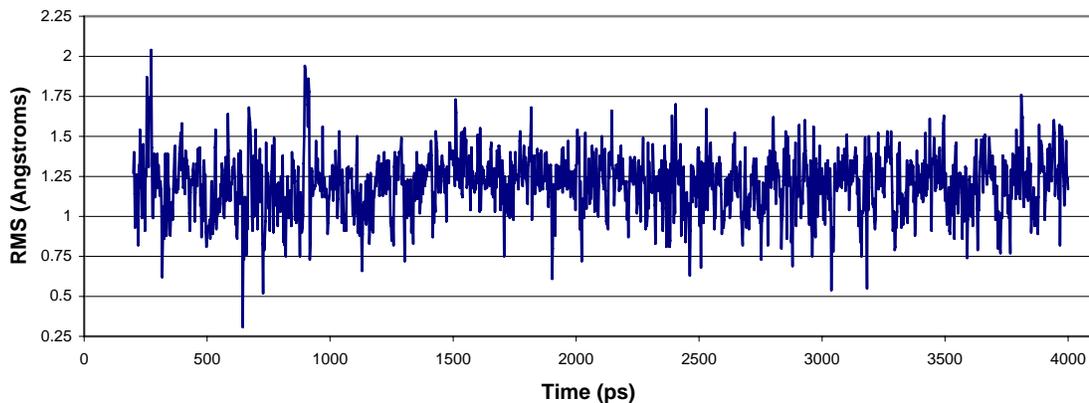
The third domain is comprised of two threonine residues and an arginine residue, seen in Figure 13. Residue 74 in negative one position of the  $\alpha$ -helix has two very distinct states, the first fluctuating around 0.45Å and the second state existing at approximately 1.1Å difference from the crystal structure. Residue 77, the threonine in the third position of the helix is in an equilibrium state for most of the simulation around 0.4Å difference from the crystal structure, but in the beginning and the end of the simulation it jumps to the higher state at 1.3Å difference. The final residue of interest for GD-AA is residue 80, an arginine. For the majority of the simulation, the residue is located between 3.0Å and 4.5Å away from the crystal structure. While this residue has a more bulky side chain, it is not expected that it would be so far removed from the crystal structure.

The residues of interest for the first domain of AA-GD are fairly far removed from the crystal structure as can be seen in Figure 14. The negative one position fluctuates around 1.5Å away from the crystal structure with the higher state above 2.5Å difference. Although the lowest

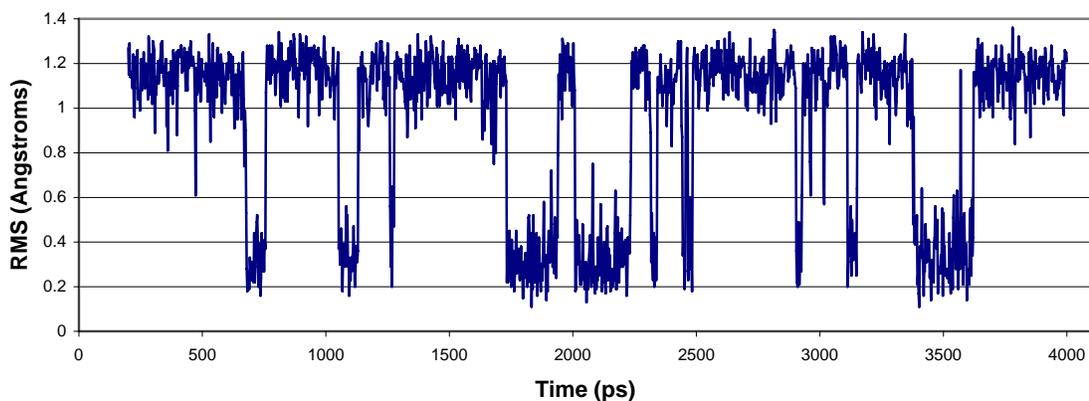
(a) Residue 18 Glutamine (Domain 1 Position -1)



(b) Residue 21 Asparagine (Domain 1 Position 3)

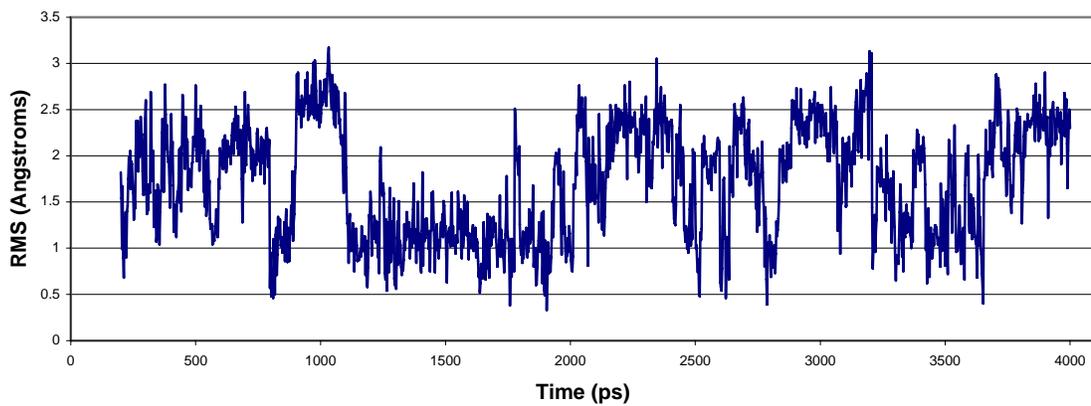


(c) Residue 24 Threonine (Domain 1 Position 6)

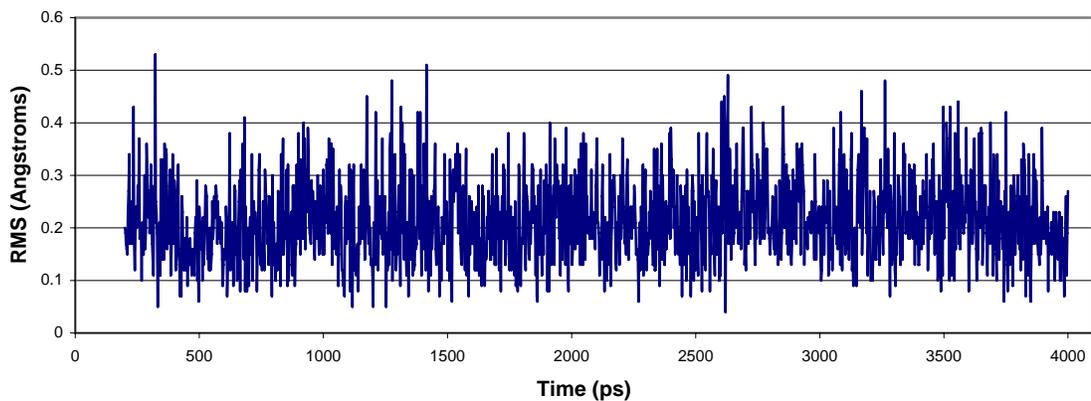


**Figure 11. Graphs of Domain 1 of GD-AA.** All graphs are the RMS between the crystal structure of 1G2D against the MD simulation data of GD-AA against time in picoseconds.

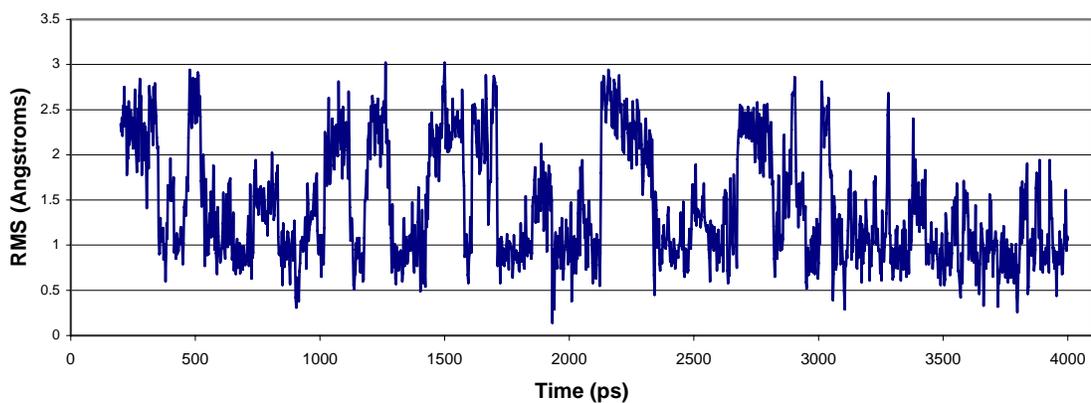
(a) Residue 46 Glutamine (Domain 2 Position -1)



(b) Residue 49 Glycine (Domain 2 Position 3)

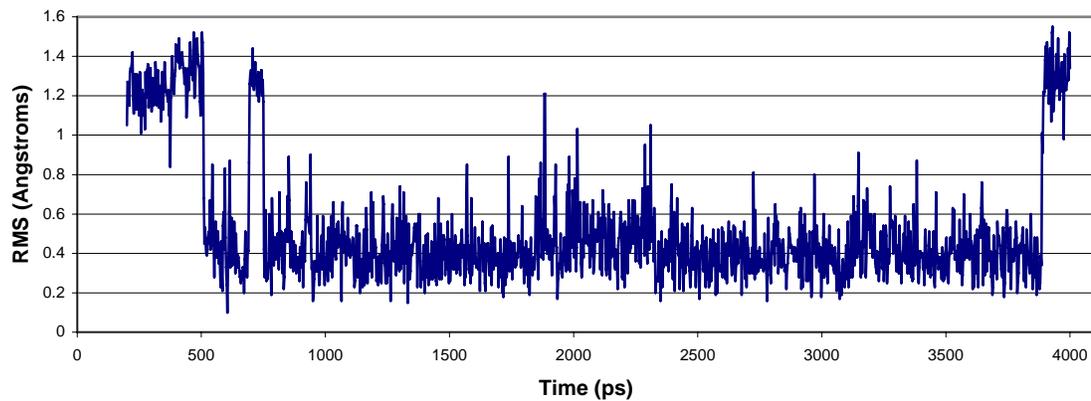


(c) Residue 52 Glutamine (Domain 2 Position 6)

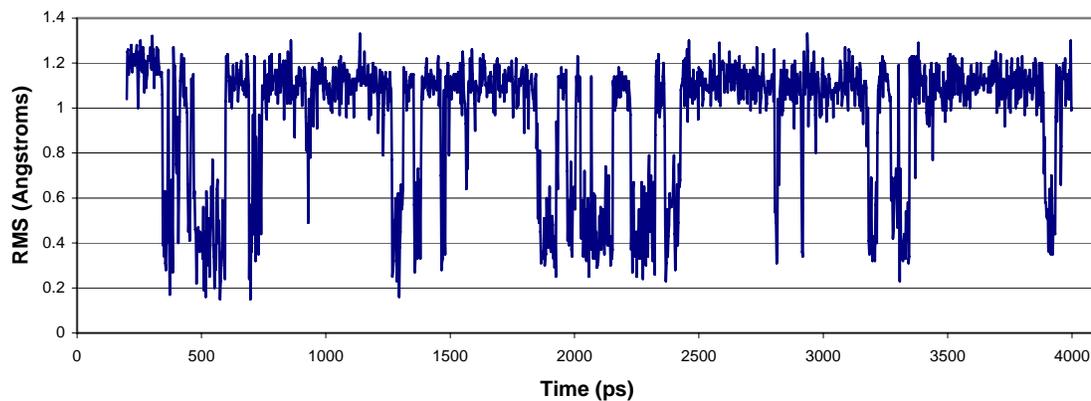


**Figure 12. Graphs of Domain 2 of GD-AA.** All graphs are the RMS between the crystal structure of 1G2D against the MD simulation data of GD-AA against time in picoseconds.

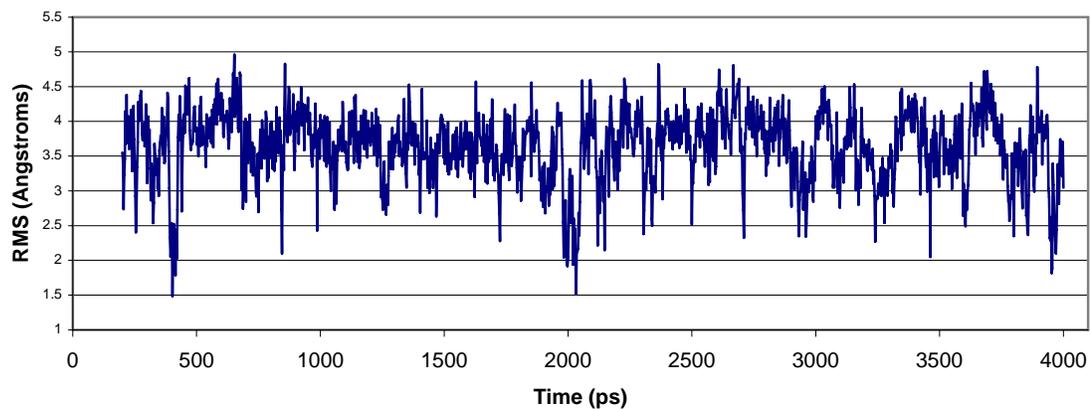
(a) Residue 74 Threonine (Domain 3 Position -1)



(b) Residue 77 Threonine (Domain 3 Position 3)



(c) Residue 80 Arginine (Domain 3 Position 6)



**Figure 13. Graphs of Domain 3 of GD-AA.** All graphs are the RMS between the crystal structure of 1G2D against the MD simulation data of GD-AA against time in picoseconds.

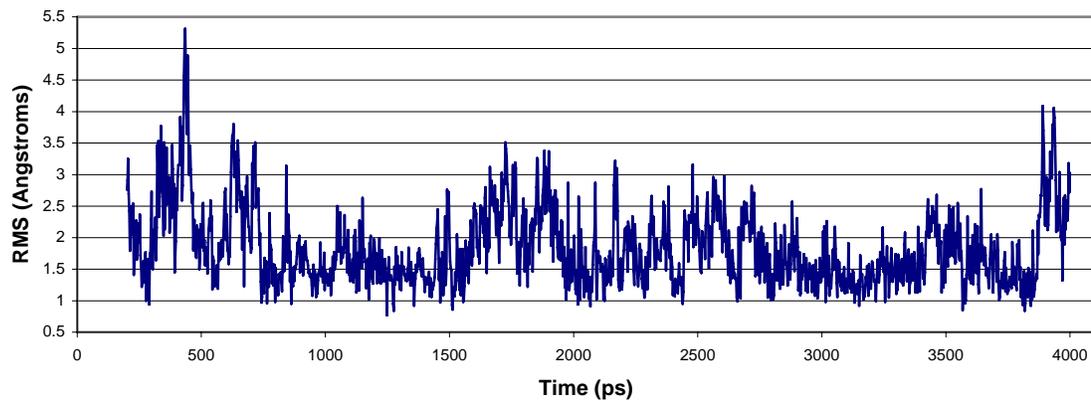
state is removed from the crystal structure, it is more desirable than the results for the previous arginine residue, GD-AA domain 3 position 6. Residue 21 vacillates between two states, and also spends time in an intermediary state. The low state is approximately  $0.5\text{\AA}$  difference from the crystal structure, where the higher state is  $1.1\text{\AA}$  away. The sixth position of the helix is occupied by another arginine residue, but unlike the arginine residue in the negative one position, this residue prefers to exist in a higher state around  $2.5\text{\AA}$  away from the crystal structure.

The residues of the second domain seem to exist in states of equilibrium, except for the first residue as can be seen in Figure 15. The arginine residue at the negative one position does not prefer to exist in any given state as it continually changes states. Towards the end of the simulation, the residue goes through a series of gradual increases followed by a sharp decline in RMS. This trend is interesting as the residue can quickly jump  $3.0\text{\AA}$  from a high state to a low state. Residue 49, a histidine, highly prefers a state that is  $1.25\text{\AA}$  away from the crystal structure. This residue, like the threonine residues seen previously also jumps sharply into a higher state of  $2.5\text{\AA}$  difference for brief periods of time. The sixth position of the  $\alpha$ -helix is a threonine and prefers a state that is  $0.4\text{\AA}$  away from the crystal structure. The higher state resonates around  $1.2\text{\AA}$  away from the crystal structure.

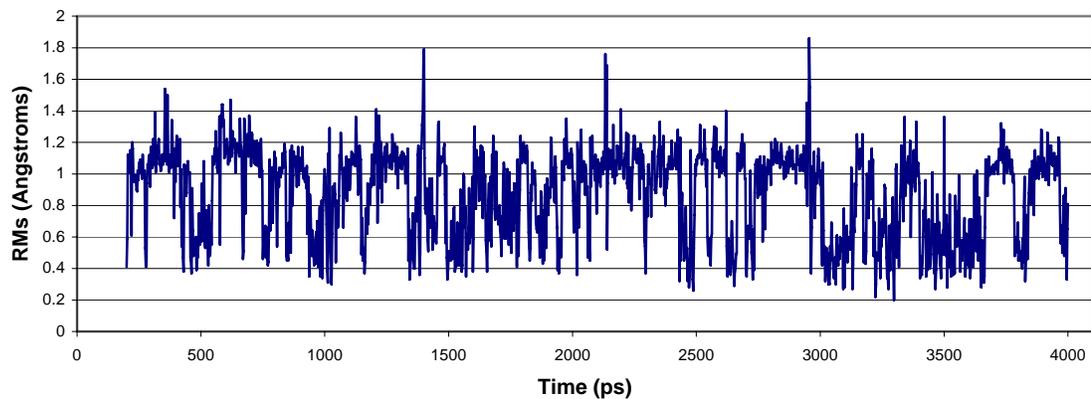
Much like the first domain, the third domain has the same sequence of residues, and has very similar RMS profiles as seen in Figure 16. The negative one position is an arginine residue and exists around  $1.5\text{\AA}$  away from the crystal structure. The third position of the helix fluctuates between  $0.5\text{\AA}$  and  $1.25\text{\AA}$  away from the crystal structure. The sixth position of the helix, an arginine residue, starts at a lower state around  $1.5\text{\AA}$ , and increases to a higher state around  $2.75\text{\AA}$  away from the crystal structure.

The model of ME-AA introduces new residues of aspartate and lysine into the DNA binding positions. The negative one position of first domain is a glutamine and is in a state of flux throughout the simulation not preferring one particular state, but the average RMS value

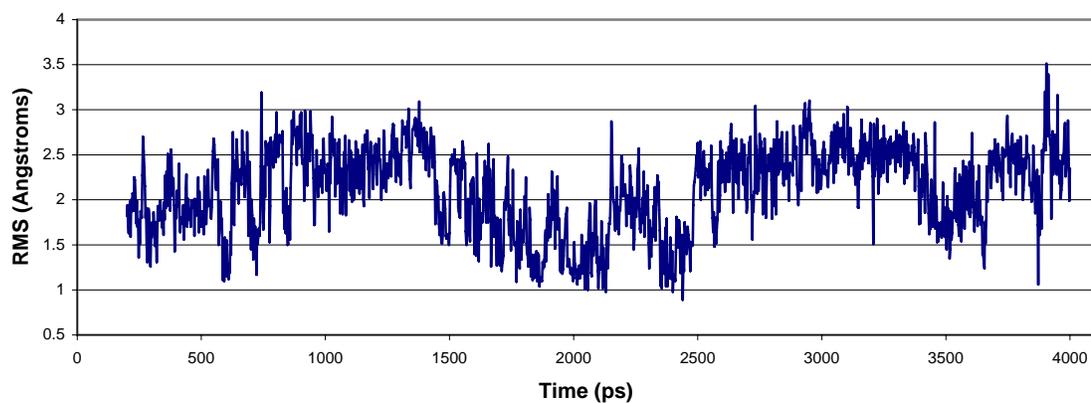
(a) Residue 18 Arginine (Domain 1 Position -1)



(b) Residue 21 Glutamate (Domain 1 Position 3)

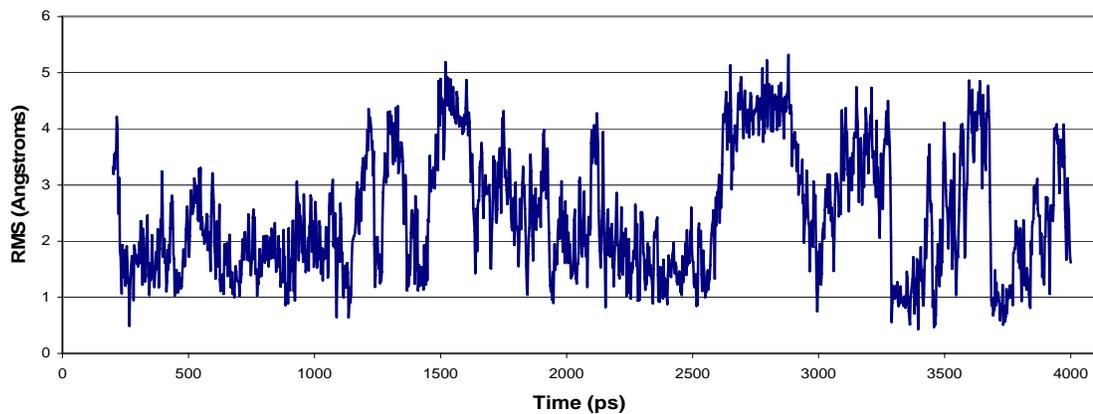


(c) Residue 24 Arginine (Domain 1 Position 6)

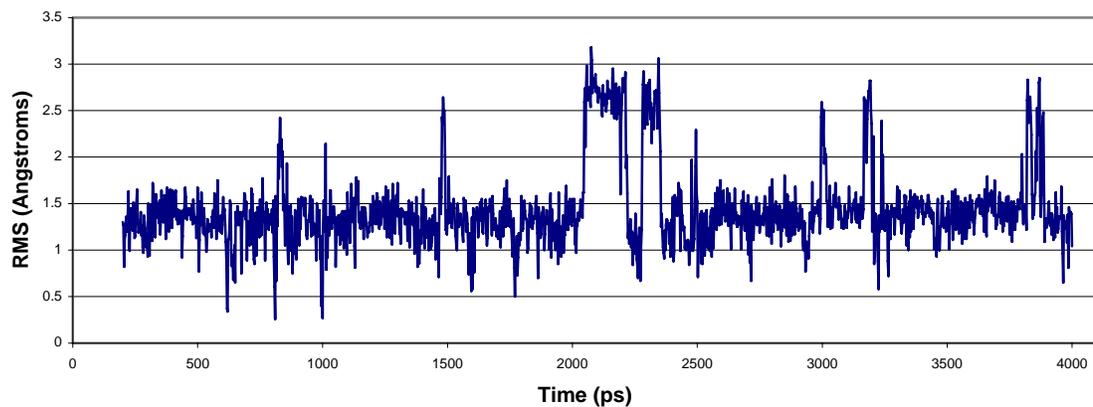


**Figure 14. Graphs of Domain 1 of AA-GD.** All graphs are the RMS between the crystal structure of 1AAY against the MD simulation data of AA-GD against time in picoseconds.

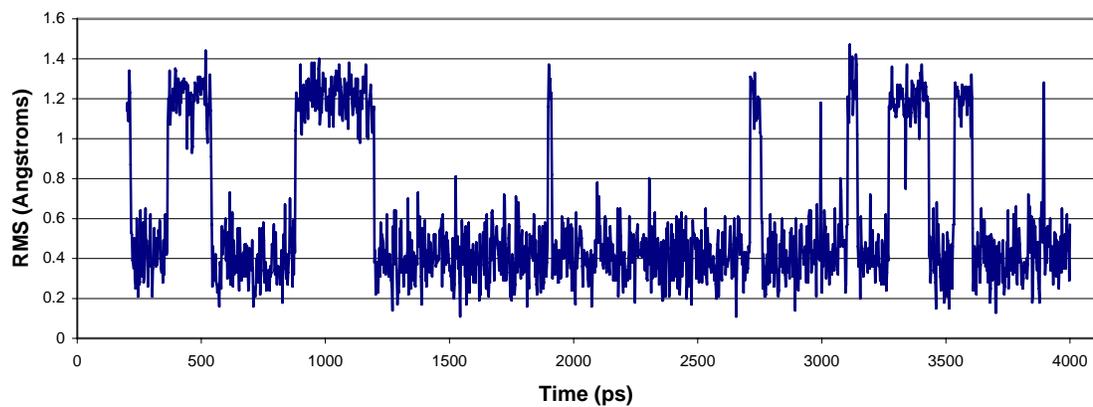
(a) Residue 46 Arginine (Domain 2 Position -1)



(b) Residue 49 Histidine (Domain 2 Position 3)

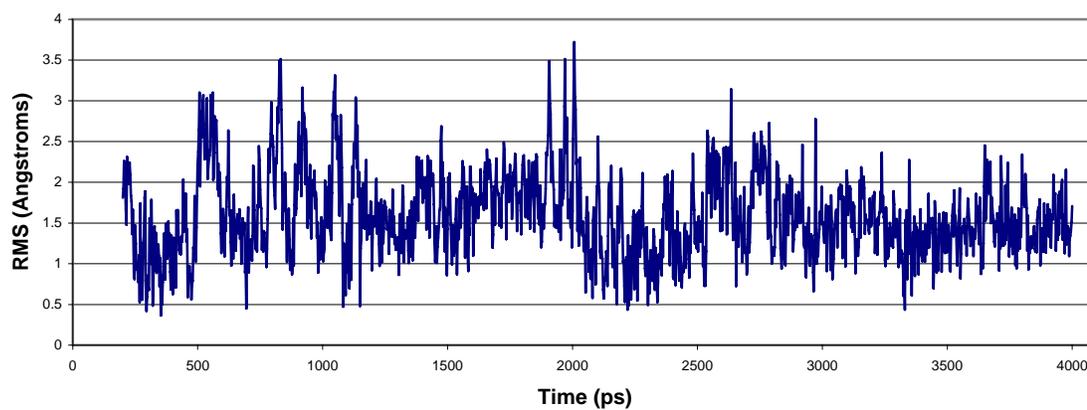


(c) Residue 52 Threonine (Domain 2 Position 6)

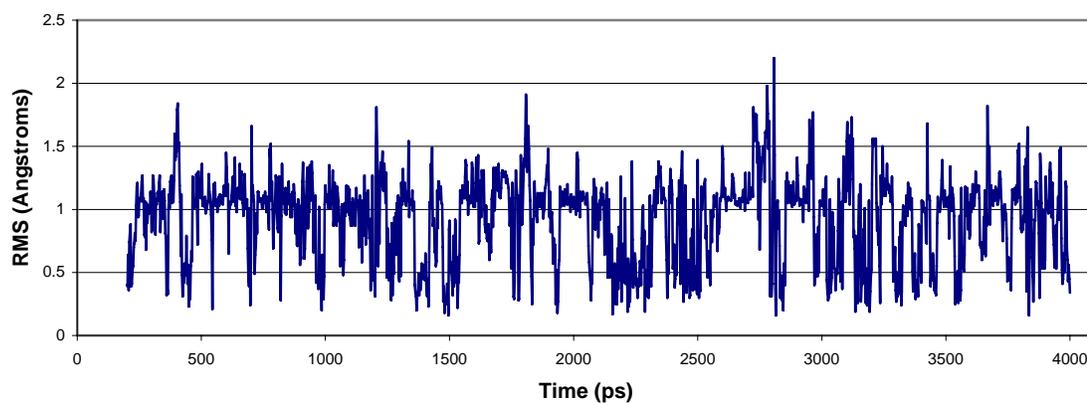


**Figure 15. Graphs of Domain 2 of AA-GD.** All graphs are the RMS between the crystal structure of 1AAY against the MD simulation data of AA-GD against time in picoseconds.

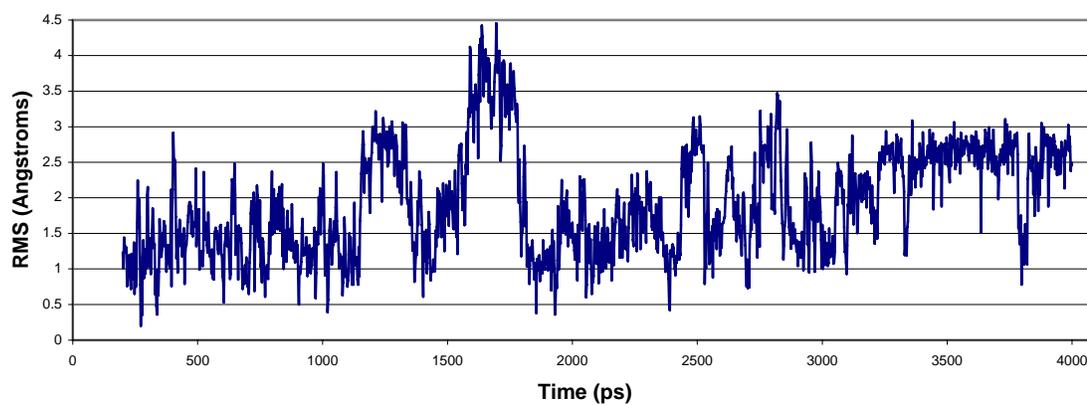
(a) Residue 74 Arginine (Domain 3 Position -1)



(b) Residue 77 Glutamate (Domain 3 Position 3)



(c) Residue 80 Arginine (Domain 3 Position 6)



**Figure 16. Graphs of Domain 3 of AA-GD.** All graphs are the RMS between the crystal structure of 1AAY against the MD simulation data of AA-GD against time in picoseconds.

throughout the simulation is 1.83Å difference from the crystal structure as can be seen in Figure 17. Residue 17, an asparagine, fluctuates around an RMS of 1.2Å. The sixth position of the  $\alpha$ -helix is occupied by a lysine residue, which is functionally similar to an arginine and a strong residue for binding to molecules with negative charges. Unlike other arginine residues that have been analyzed, the lysine here fluctuates between 1.0Å and 3.0Å difference from the crystal structure.

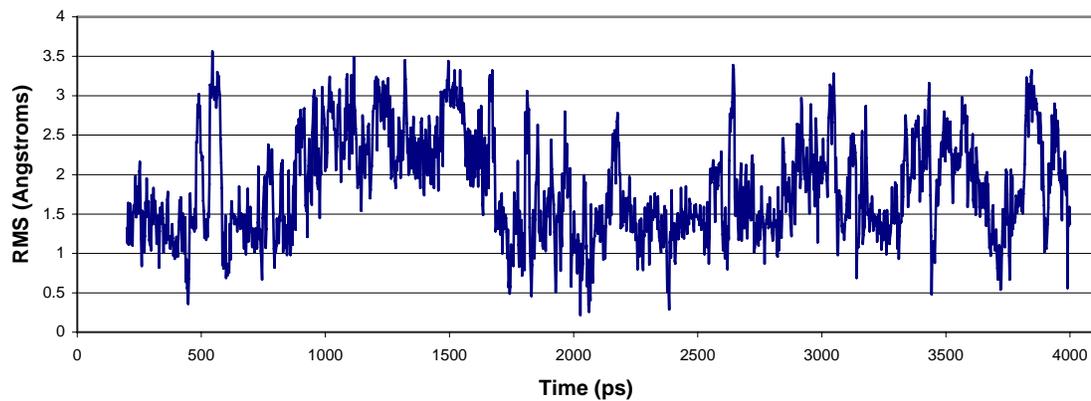
The second domain is interesting in the way it is constructed and thus how it behaves within the simulation. The negative one position is a neutrally charged glutamine which has a good deal of fluctuation, but is not uncommon when compared to other glutamine residues modeled and is seen in Figure 18. The third position of the  $\alpha$ -helix is an aspartate, which is negatively charged. It fluctuates around the 1.0Å difference from the crystal structure. What is interesting is that the negative residue is generally positioned between two positively charged residues, as domains 1 and 3 of AA-GD exhibit, not between a neutral and positively charged residue. The lysine in position 6 of the  $\alpha$ -helix resonates around 2.1Å difference from the crystal structure.

The third domain of ME-AA has a unique RMS profile and can be seen in Figure 19. The negative one position arginine exists primarily around 1.5Å away from the crystal structure, but also fluctuates to a much higher state above 3.0Å. The histidine residue in position 3 of the  $\alpha$ -helix fluctuates randomly between two states, one approximately at 1.0Å difference from the crystal structure, and the second state around 2.25Å difference. The final residue of interest, an arginine in position 6 of the  $\alpha$ -helix prefers a state that is over 4.0Å away from the crystal structure, but the RMS will periodically fall towards the observed state in the crystal.

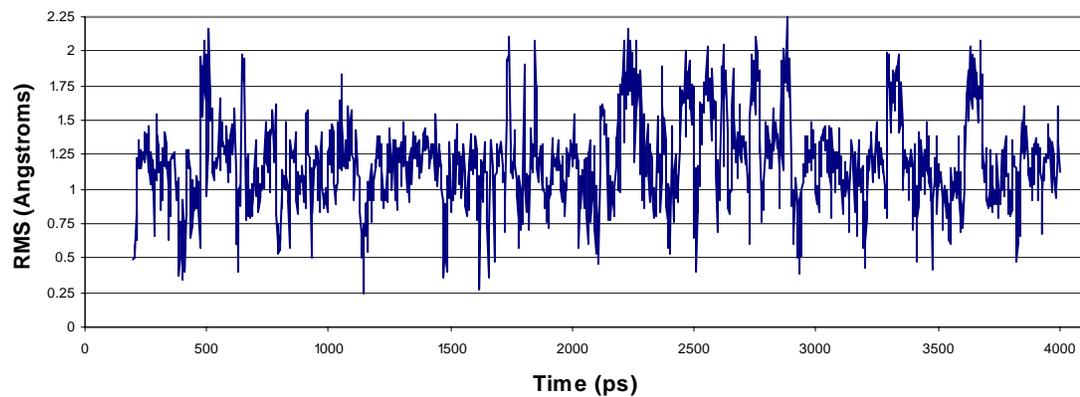
### *Clustering Results*

Clustering the side chains is the final step in determining the overall structure of the homology modeled protein. Since we are trying to determine the structure that is closest to the

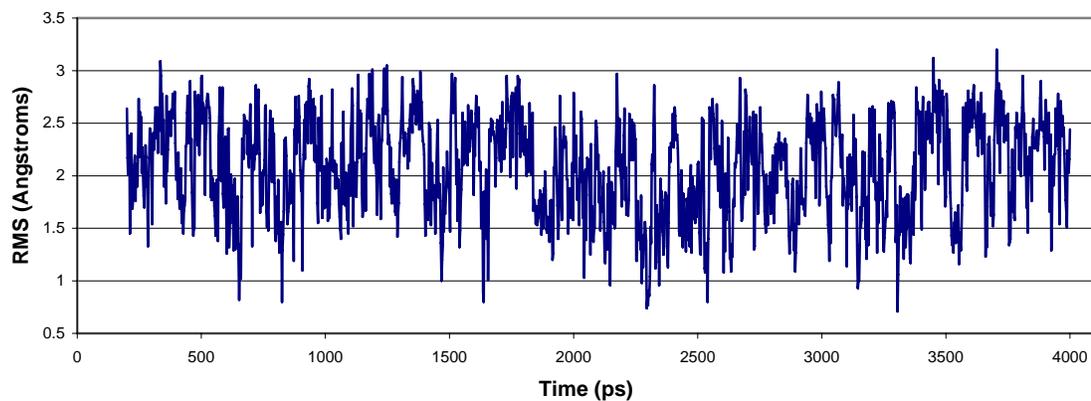
(a) Residue 16 Glutamine (Domain 1 Position -1)



(b) Residue 19 Asparagine (Domain 1 Position 3)

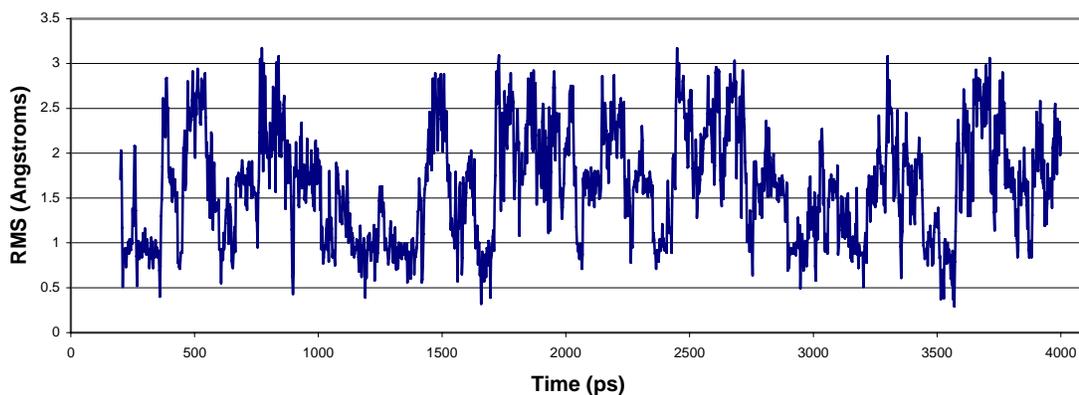


(c) Residue 22 Lysine (Domain 1 Position 6)

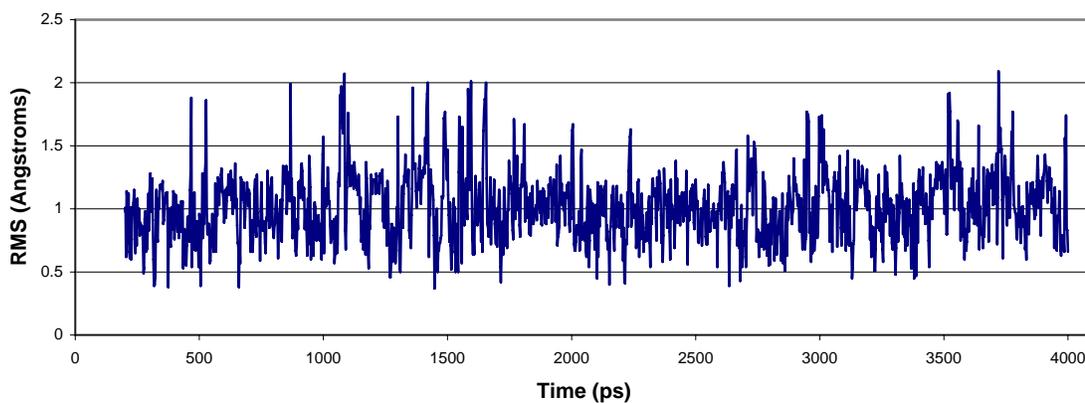


**Figure 17. Graphs of Domain 1 of ME-AA.** All graphs are the RMS between the crystal structure of 1MEY against the MD simulation data of ME-AA against time in picoseconds.

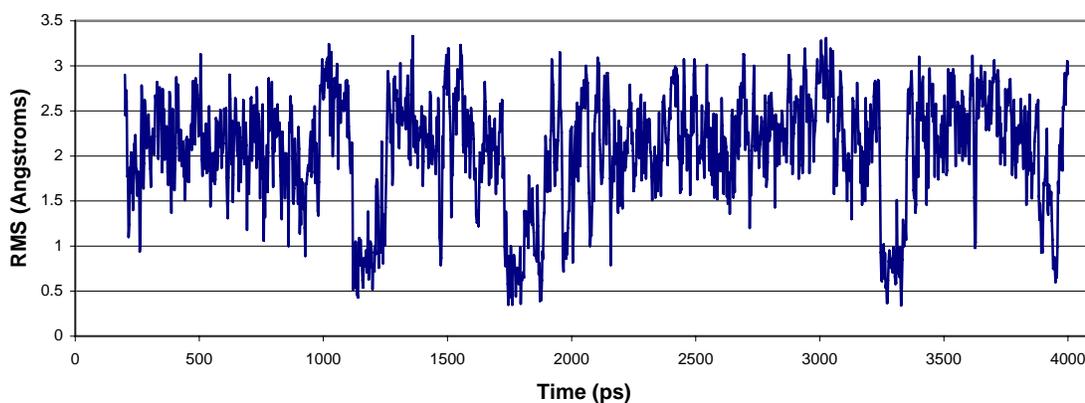
(a) Residue 44 Glutamine (Domain 2 Position -1)



(b) Residue 47 Aspartate (Domain 2 Position 3)

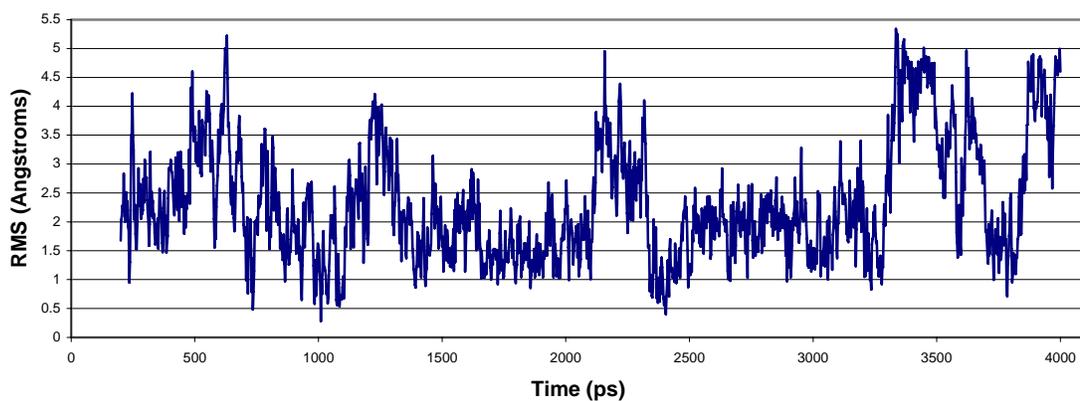


(c) Residue 50 Lysine (Domain 2 Position 6)

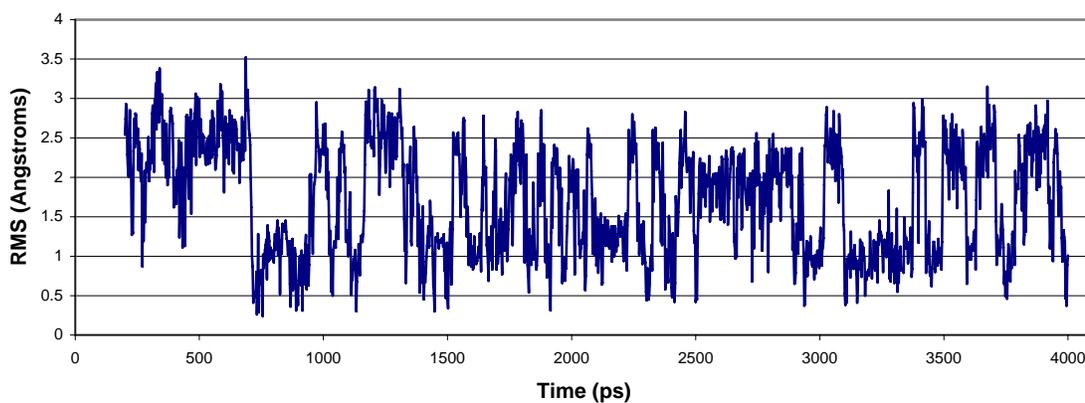


**Figure 18. Graphs of Domain 2 of ME-AA.** All graphs are the RMS between the crystal structure of 1MEY against the MD simulation data of ME-AA against time in picoseconds.

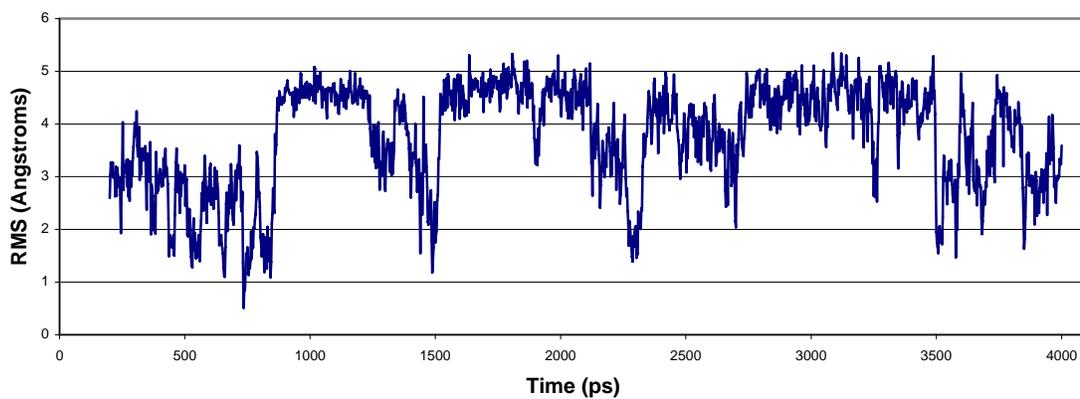
(a) Residue 72 Arginine (Domain 3 Position -1)



(b) Residue 75 Histidine (Domain 3 Position 3)



(c) Residue 78 Arginine (Domain 3 Position 6)



**Figure 19. Graphs of Domain 3 for ME-AA.** All graphs are the RMS between the crystal structure of IMEY against the MD simulation data of ME-AA against time in picoseconds.

crystal structure, the model for a given amino acid with the highest amount of neighbors through RMS analysis within a 0.75Å radius is selected. The residues from the PDB files of various models are concatenated together. The reason this is possible is due to the fact that the backbone is held constant throughout the simulation so the residues, when concatenated, form one connected protein.

In creating one protein of three domains, the loop regions between domains are more difficult to reconstruct as there is duplication of residues between the different domains. Information on the individual domains for each homology model and the overlapping residues for each domain can be found in Appendix C. For purpose of simplicity, the first residue of overlap between the two domains is captured from the N-terminal domain and the second residue of overlap is captured from the C-terminal domain. For example, if the sequence for domain one is ...**GQKP** and the sequence for domain two is **QKPF**... the glycine and glutamine from domain one and the lysine and phenylalanine from domain two are kept. By selecting residues in this manner any terminal residue from each domain is discarded. Additionally, the terminal residues generally have a high difference in the RMS from the crystal structure conferring additional support towards selecting non-terminal residues. Examining the completed structures uses seven residues of interest per domain: the three residues used for DNA binding, the two conserved cysteine residues and two conserved histidine residues.

The completed structure for GD-AA can be seen in Figure 20 as having fairly good conservation of the side chains from the crystal structure. As can be seen in Figure 20a, the residues of interest are almost identical in structure to the crystal. There is a rotation of the threonine in position 6 of the  $\alpha$ -helix and the foreground cysteine (C1) is slightly rotated from the crystal, but still conforming to an architecture that is appropriate for binding the zinc ion. Most of the residues that are seen in a conformation differing from the crystal structure lie within the  $\beta$ -strands, which is expected as they are exposed to solvent. The second domain, Figure 20b has a

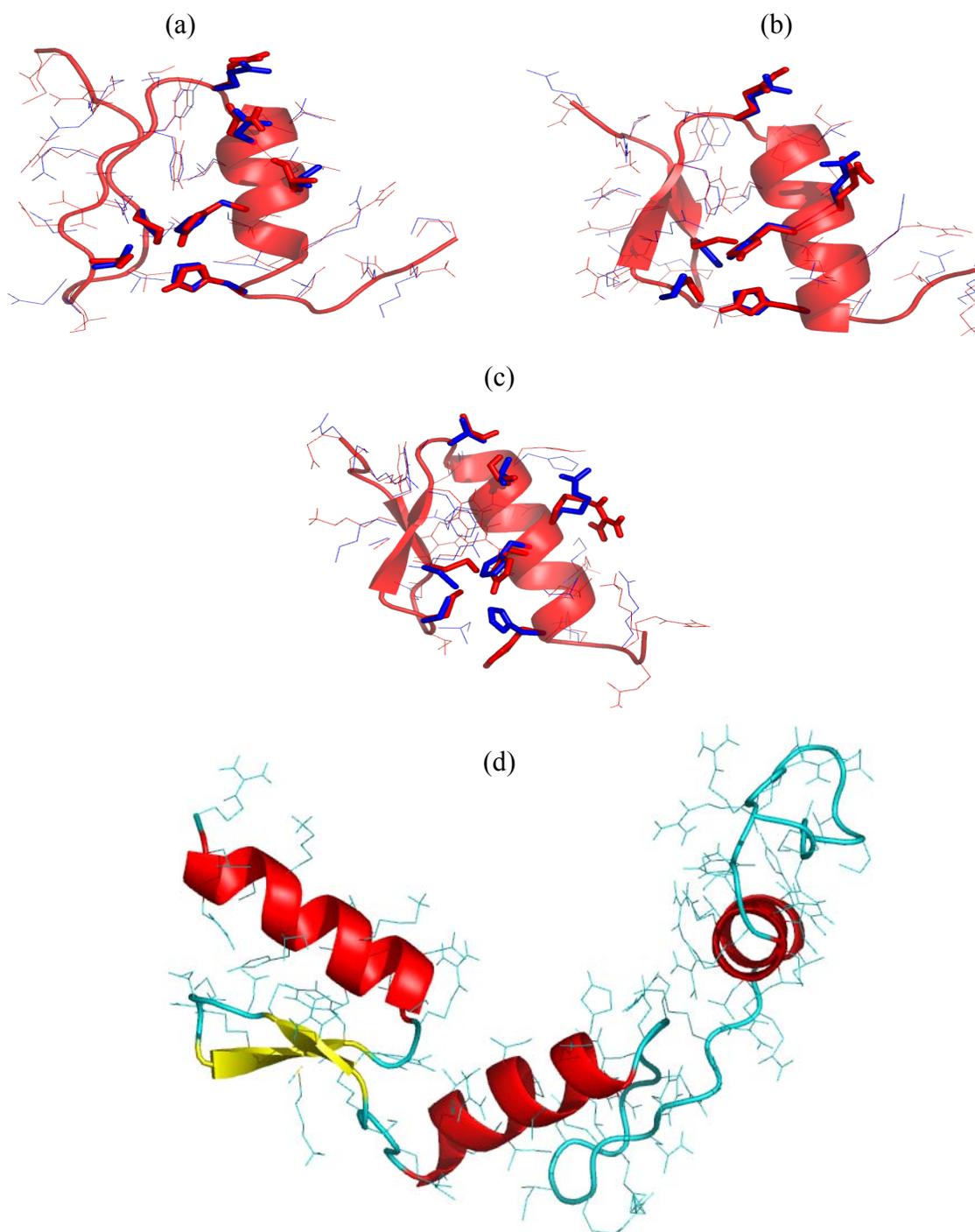


Figure 20. GD-AA Clustered Domain Results. The homology model is shown in red with the crystal structure superimposed in blue with the conserved residues and DNA binding residues in stick form for (a), (b), and (c). (a) Domain 1. (b) Domain 2. (c) Domain 3. (d) All domains for GD-AA concatenated together and colored by secondary structure: yellow is for a  $\beta$ -sheet, red for an  $\alpha$ -helix, and cyan for the remaining residues. Domain 1 is on the left.

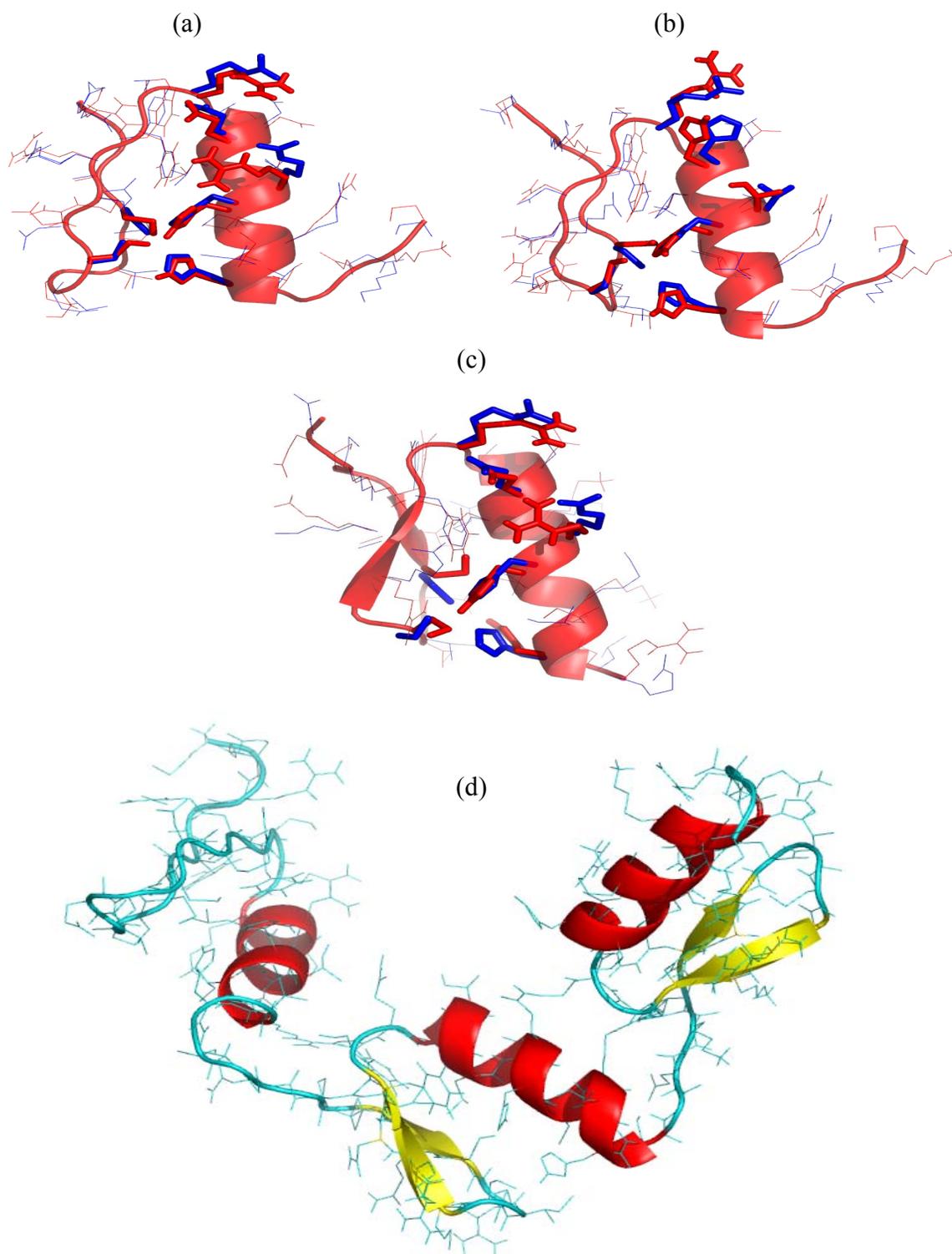


Figure 21. AA-GD Clustered Domain Results. The homology model is shown in red with the crystal structure superimposed in blue with the conserved residues and DNA binding residues in stick form for (a), (b), and (c). (a) Domain 1. (b) Domain 2. (c) Domain 3. (d) All domains for AA-GD concatenated together and colored by secondary structure: yellow is for a  $\beta$ -sheet, red for an  $\alpha$ -helix, and cyan for the remaining residues. Domain 1 is on the left.

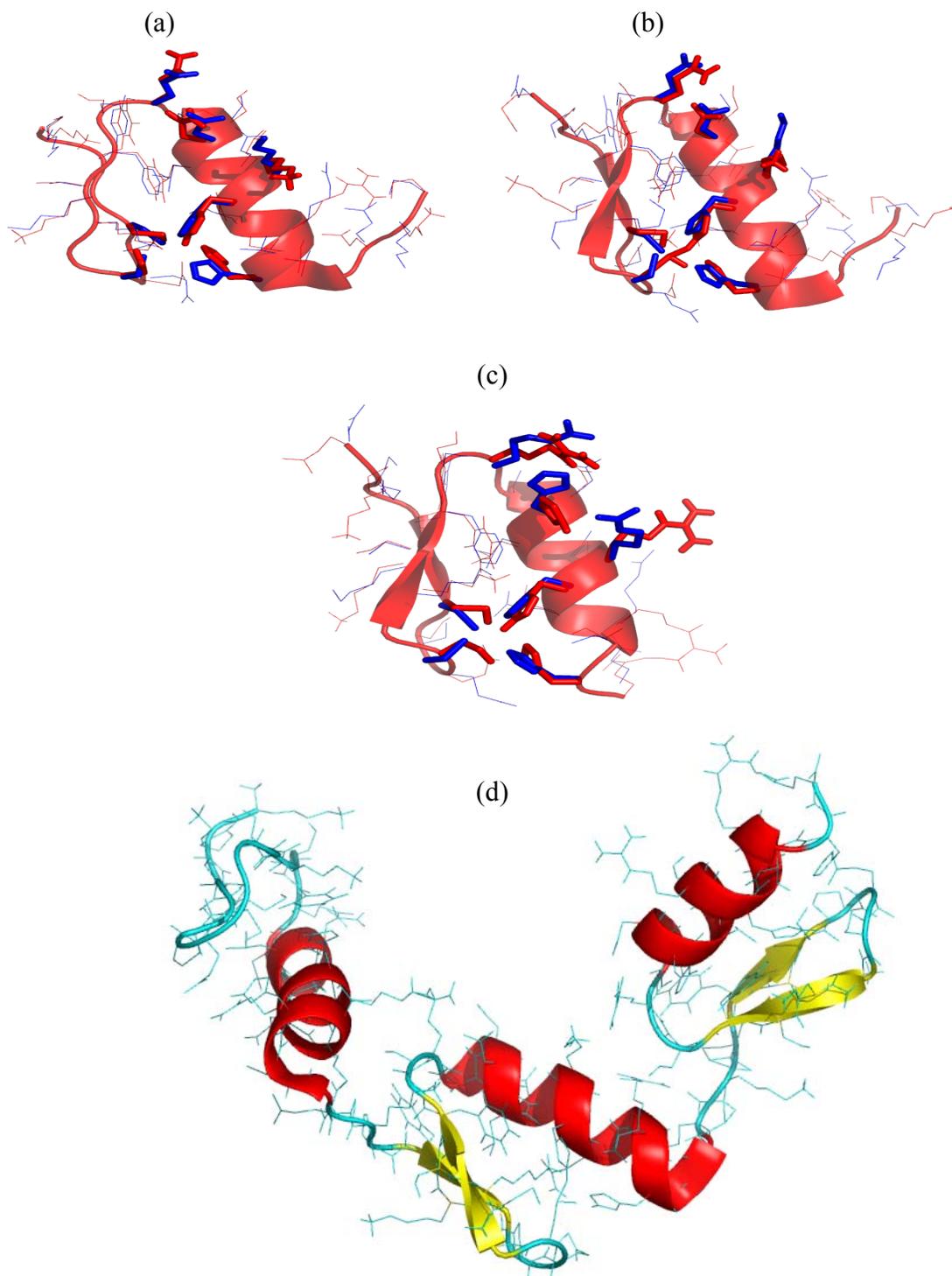


Figure 22. ME-AA Clustered Domain Results. The homology model is shown in red with the crystal structure superimposed in blue with the conserved residues and DNA binding residues in stick form for (a), (b), and (c). (a) Domain 1. (b) Domain 2. (c) Domain 3. (d) All domains for GD-AA concatenated together and colored by secondary structure: yellow is for a  $\beta$ -sheet, red for an  $\alpha$ -helix, and cyan for the remaining residues. Domain 1 is on the left.

**Table 1. Residues of Interest for GD-AA.** Listed for each domain is the time in which the residue was found and the RMS difference from the crystal structure.

Domain 1			Domain 2			Domain 3		
Residue	Time (ps)	RMS	Residue	Time (ps)	RMS	Residue	Time (ps)	RMS
C1	1820	1.08	C1	2432.5	0.45	C1	485	0.95
C2	815	0.91	C2	325	1.26	C2	2010	0.49
H1	1320	1.23	H1	2105	0.64	H1	1072.5	0.99
H2	3302.5	1.18	H2	2885	1.24	H2	1665	2.80
Helix -1	1120	0.97	Helix -1	377.5	0.86	Helix -1	2927.5	0.39
Helix 3	2622.5	0.25	Helix 3	3382.5	1.14	Helix 3	2820	0.71
Helix 6	1352.5	0.84	Helix 6	1605	1.25	Helix 6	515	3.62

less conserved structure than the first domain. The DNA binding residues are highly similar to those in the crystal. The third position of the  $\alpha$ -helix cannot be seen due to the fact that the residue is glycine and the side chain is only a hydrogen atom. The glutamine in the sixth position of the  $\alpha$ -helix is rotated about the  $C_{\delta}$  which puts the residue in a less desirable position to make contact with the DNA. The third domain, Figure 20c, is the least desirable clustering formation for the protein. Positions negative one and three of the  $\alpha$ -helix are identical to the crystal, but the sixth position, an arginine, is rotated in the exact opposite direction of the crystal structure. The RMS between these two residues is 3.62Å, and can be seen at time point 1287.5ps in Figure 13c. This is also the only of the domains in this structure to such a conflict between the conserved residues and the crystal structure. The RMS and the time point for the residues of interest can be found in Table 1. Both the cysteine (C1) in the background and the histidine (H2) in the foreground are rotated out of position and thus would not bind to the zinc atom if present. The residues in the  $\beta$ -strands have much higher similarity to the crystal structure than the residues in either domain one or two. A full table of RMS versus the crystal structure and time point for each residue within each domain can be found in Appendix D. The complete model for GD-AA can be seen in Figure 20d having the  $\beta$ -strands only recognized in the first domain, which is the only domain that did not have the  $\beta$ -strands recognized after homology modeling (Figure 9a).

**Table 2. Residues of Interest for AA-GD.** Listed for each domain is the time in which the residue was found and the RMS difference from the crystal structure.

Domain 1			Domain 2			Domain 3		
Residue	Time (ps)	RMS	Residue	Time (ps)	RMS	Residue	Time (ps)	RMS
C1	1095	0.94	C1	3792.5	1.07	C1	2592.5	0.87
C2	280	2.05	C2	2437.5	0.49	C2	1410	0.84
H1	882.5	1.77	H1	3847.5	1.52	H1	1117.5	1.25
H2	2850	1.27	H2	2490	1.44	H2	655	2.84
Helix -1	3102.5	1.13	Helix -1	2970	2.3	Helix -1	3402.5	1.39
Helix 3	237.5	0.97	Helix 3	1087.5	1.27	Helix 3	2592.5	1.07
Helix 6	1585	2.18	Helix 6	837.5	0.67	Helix 6	1415	1.55

The completed structure for AA-GD is slightly more disorganized than the completed structure for GD-AA. The first domain, Figure 21a, does not show much similarity between the residues of interest and the crystal structure. Although the arginine (position negative one) and the glutamate (position three) in the  $\alpha$ -helix are quite similar to the crystal structure, the arginine in the sixth position is removed from crystal and bent inward towards the middle of the domain. The RMS between this residue and the crystal is 2.18Å. Additionally, the conserved residues are not in favorable positions to coordinate a zinc ion. The cysteine residues are angled away from one another, and the second histidine is turned down away from the crystal. The RMS and time positions for the AA-GD can be found in Table 2. The second domain, Figure 21b, is more ordered than the first domain. The DNA binding residues in this domain oriented in a similar manner to the crystal structure with a slight difference in the positioning of the alpha carbon. The only residues with positions that could be of concern are the C1 and H2 residues, which are both angled in a different manner than the crystal structure but in relative position to continue to coordinate a zinc ion. The quality of the third domain lied between the first and second domains. The arginine (position negative one) and glutamate (position three) are in line with the crystal structure. Once again, the arginine in the sixth position is turned inward towards the middle of the domain, leaving it in a difficult, but not impossible position for binding to DNA. The RMS between the crystal and the clustered arginine residue is 1.55Å, which is low for this residue. In

**Table 3. Residues of Interest for ME-AA.** Listed for each domain is the time in which the residue was found and the RMS difference from the crystal structure.

Domain 1			Domain 2			Domain 3		
Residue	Time (ps)	RMS	Residue	Time (ps)	RMS	Residue	Time (ps)	RMS
C1	2385	0.5	C1	3605	0.64	C1	1707.5	0.6
C2	3412.5	1.25	C2	1627.5	0.62	C2	2287.5	0.98
H1	2847.5	0.49	H1	1412.5	0.67	H1	1265	0.93
H2	1170	2.3	H2	1665	2.22	H2	660	1.88
Helix -1	735	1.18	Helix -1	1320	0.71	Helix -1	2680	1.98
Helix 3	2420	1.11	Helix 3	3735	1.07	Helix 3	2262.5	2.08
Helix 6	3075	2.21	Helix 6	957.5	2.4	Helix 6	2005	4.51

addition, both cysteine residues are angled away from one another, and the H2 residue is once again angled downwards.

The completed structure for ME-AA has good structural preservation of the conserved residues as can be seen in Figure 22. The first domain, Figure 22a, is conserved in the residues of interest except for position six of the  $\alpha$ -helix, a lysine residue. The lysine has a kink in the side chain near the beta carbon and is directed away from the supposed location of the DNA and has an RMS of 2.21Å against the crystal structure. The second domain is similar to the first domain, except that the lysine in the sixth position of the  $\alpha$ -helix is rotated about the beta-carbon to produce a side chain that is directed away from the location of the DNA with an RMS difference of 1.34Å difference from the crystal. The third domain is slightly more chaotic in terms of the side chains of interest on the  $\alpha$ -helix. All three residues on the  $\alpha$ -helix are not aligned with the crystal structure. Detailed RMS descriptions can be found in Table 3. All three of the residues are directed away from where the DNA would be located within the structure. The conserved residues also have more fluctuation from the crystal structure than the other two domains, but are still oriented properly for zinc coordination.

## Discussion

### *Protein Structure Quality*

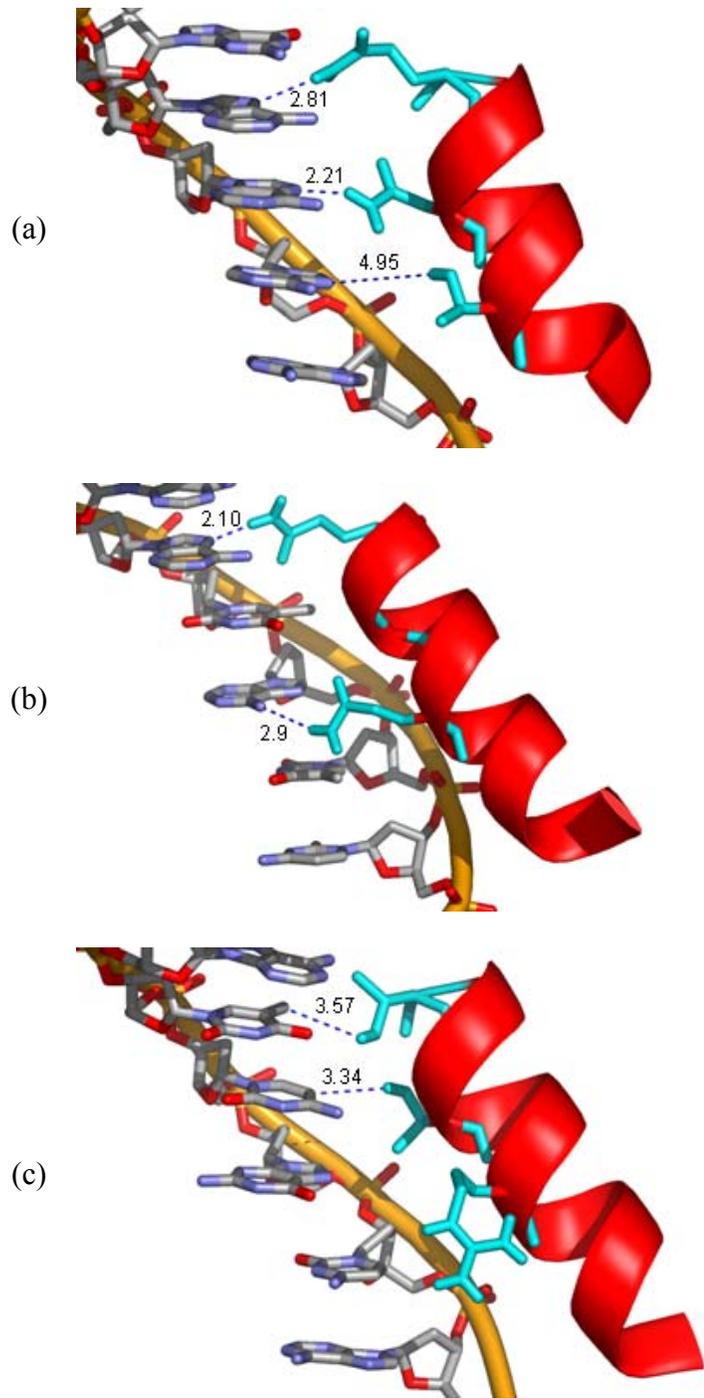
The structures of the final state of the clustered proteins are favorably aligned with the respective crystal structure. When looking at the whole DNA binding region of the protein, the final structures exhibit RMS values that are desirable for the purpose of protein – DNA interaction prediction (Table 4). The structures have much better RMS scores in between the backbone than the initial structures from homology modeling through Consensus. Given that the RMS scores between all atoms within the models is around 1.7Å for each completed structure, there is confidence that these structures are of sound quality, but that there is some difference between the structures. There may possibly be regions within the protein that are more dissimilar than other regions that can contribute to the increased RMS score. Suspect areas are the loop regions between the domains, the loop region between the  $\beta$ -strands as well as the N and C terminals, as the homology modeling process ranked those regions as not having confidence in the prediction. Differences can be visualized by referring back to Figure 9.

### *Protein – DNA Interaction*

Although the completed proteins were not used in protein – DNA interaction prediction programs, it is still possible to observe potential interactions between the protein domains and the DNA present in the crystal structures. By superimposing the DNA from the crystal structure onto

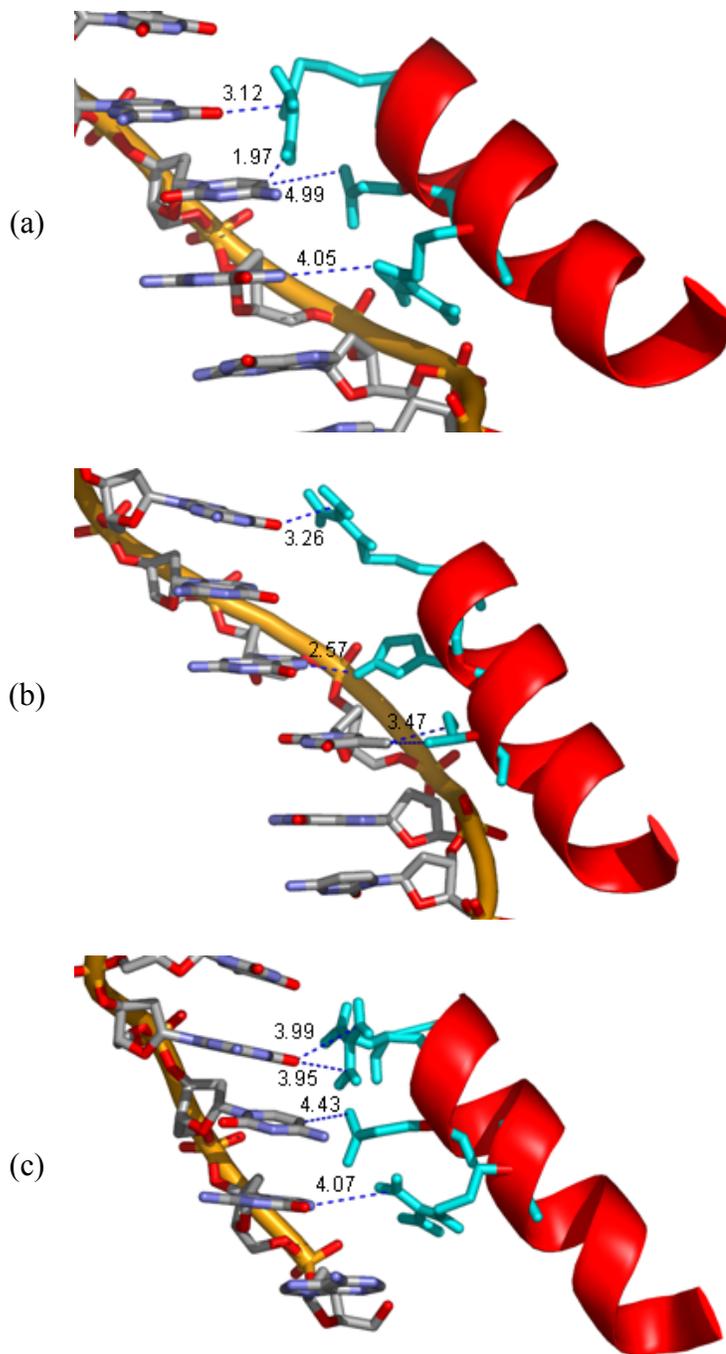
**Table 4. RMS values between completed models and the crystal structure.** Three different RMS calculations were made between a given crystal structure and completed model: between all atoms (amino acids), between C- $\alpha$  atoms, and between atoms in the backbone.

RMS	1G2D vs. GD-AA	1AAY vs. AA-GD	1MEY vs. ME-AA
Amino Acids	1.68 Å	1.75 Å	1.70 Å
C- $\alpha$	0.52 Å	0.59 Å	0.60 Å
Backbone	0.57Å	0.63Å	0.66Å



**Figure 23. Potential Protein - DNA Interactions for GD-AA.** The DNA helix from the crystal structure for 1G2D was superimposed upon the completed structure for GD-AA. For each figure the three residues in the  $\alpha$ -helix are the -1, 3 and 6 positions from top to bottom. (a) Domain 1: Residues Q, N, and T. (b) Domain 2: Residues Q, G (not pictured), Q. (c) Domain 3: Residues T, T, R.

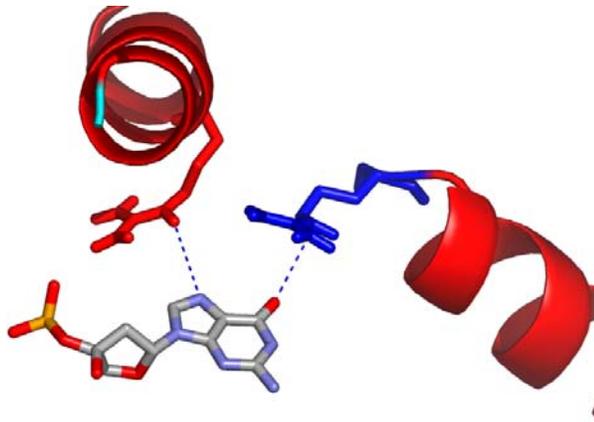
the completed structure of the homology model insight can be gained as to the probability of the protein making the necessary DNA interactions for binding. The completed structure for GD-AA is the best in terms of potential protein - DNA interactions (Figure 23). As a generalization, the negative one residue for all three domains has very strong interaction potential. Although the exact orientation is not ideal, the distance for interaction is in the scope needed for short range forces to take effect, less than 3 Å, in two of the three domains. The first domain was the only one of the three to have the key residues oriented in a



**Figure 24. Potential Protein DNA Interactions for AA-GD.** The DNA helix from the crystal structure for 1AAY was superimposed upon the completed structure for AA-GD. For each figure the three residues in the  $\alpha$ -helix are the -1, 3 and 6 positions from top to bottom. (a) Domain 1: Residues R, E and R. (b) Domain 2: Residues R, H and T. (c) Domain 3: Residues R, E and R.

manner to create a possible interaction between the protein and the DNA. The second and third domains are each missing an interaction. Domain two has the second residue as a glycine, which will not have a side chain to make a viable interaction, and the arginine in the third domain is oriented away from the DNA and thus no interaction is possible in this state.

The structure for AA-GD has more possible interactions between the key residues and DNA than the previous structure of GD-AA. In this model, each domain has at least three possible interactions between the protein and DNA due to the orientation of the side chains (Figure

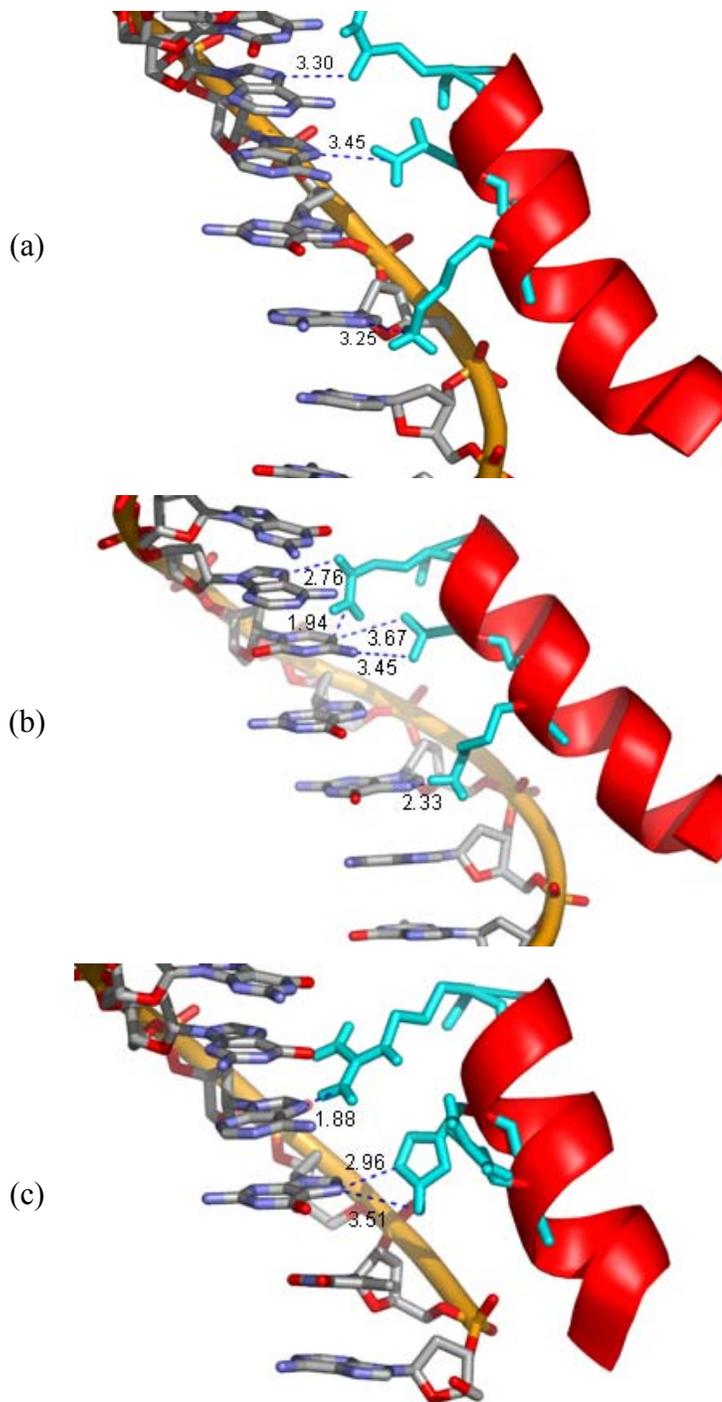


**Figure 25. Favorable Coordination of a Guanine by Residues in Two Domains.** The arginine from domain one is shown in red, and the arginine from domain two is shown in blue. Although both residues are coordinating the same nucleotide, a guanine, they are interacting with different areas of the nucleotide.

has each residue making contact with only one base, but the arginine in the negative one position is not interacting with the expected nucleotide. The arginine is actually interacting with the guanine that the arginine in the sixth position of the  $\alpha$ -helix in domain one is interacting, but on different atoms. Interestingly, the two residues do not conflict with one another as would be expected since the two domains were modeled independent of one another (Figure 25). Both the histidine and threonine are in favorable orientations to make the proper interactions. While there is potential for making correct contact with the DNA for domain three, the distance of the side chains is not close enough to create hydrogen bonds or short range forces.

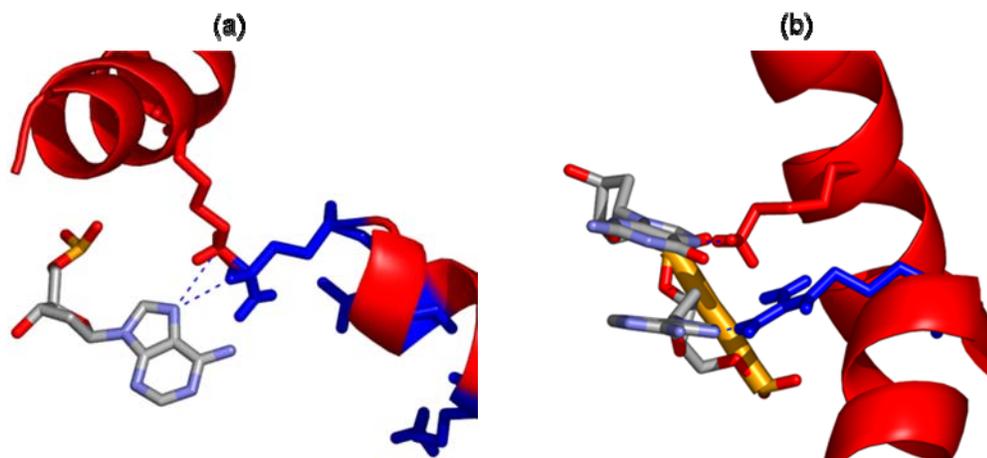
The model of ME-AA has some interesting side chain conformations that complicate the binding of DNA. ME-AA is the only structure that contains lysine residues in the key residues for binding to DNA. All three domains for position six of the  $\alpha$ -helix do not have a desirable conformation for binding to DNA in the proper and experimentally defined manner (Figure 26). For positions negative one and three in the  $\alpha$ -helix in domain one, the orientation of the residues allows for potential interaction between the protein and DNA. The lysine in position six is

24). The first domain has four potential interactions, with the first residue, an arginine, potentially interacting with two different bases in the DNA. Each  $\text{NH}_2$  group of the arginine could make contact with the DNA. Both of these interactions are more favorable than either of the two interactions for the third and sixth residues in the  $\alpha$ -helix, which are no closer than 4 Å to the DNA. The second domain



**Figure 26. Potential Protein DNA Interactions for ME-AA.** The DNA helix from the crystal structure for 1MEY was superimposed upon the completed structure for ME-AA. For each figure the three residues in the  $\alpha$ -helix are the -1, 3 and 6 positions from top to bottom. (a) Domain 1: Residues Q, N and K. (b) Domain 2: Residues Q, D and K. (c) Domain 3: Residues R, H and R.

skipping a nucleotide and interacting with the next nucleotide, an adenine, which is also being coordinated by the glutamine in position negative one in domain two. The interactions on the adenine nucleotide are unfavorable due to the fact that the two arginine residues are attempting to occupy the same space (Figure 27a). The aspartate residue in domain two can have two possible interactions with same nucleotide, both in an advantageous orientation. Once again, the lysine residue in domain two is positioned so that it is skipping a nucleotide, however the way that domain three is



**Figure 27. Favorable and Unfavorable DNA Interactions within ME-AA.** (a) Unfavorable DNA interactions between a lysine (red) in domain one and a glutamine (blue) in domain two. Both residues are attempting to interact with the same atom in the adenine and are occupying the same space. (b) Favorable interactions with a lysine (red) in domain two and an arginine (blue) in domain three.

positioned allows for there to be no conflict in the interactions between the arginine in the negative one position and the DNA (Figure 27b). The histidine in position three of domain three is not properly oriented to make a strong interaction with the DNA, but is within proximity to have both short and long range interactions. The arginine in position six is not in a desirable conformation for binding to the DNA.

Through close examination of the individual amino acids with respect to the nucleotides available, one can deduce that these structures are of sufficient quality to be able to perform protein – DNA interaction prediction algorithms. Some of the residues discussed will need to be further examined in terms of the clustering results to find a better orientation, but to do so in an automated fashion would require the need for a robust algorithm for scoring clusters.

Additionally, the DNA shown in these examples has been slightly distorted due to the fact that it originated from a crystal structure with a bound EGR protein, thus precise measurements are biased towards interactions occurring after the protein has bound to the DNA.

### *Protein – DNA Proposed Binding Mechanism*

The residues of interest from analysis against the crystal structure shed light into a possible mechanism for EGR proteins binding to DNA. An EGR protein will not bind to DNA in the absence of zinc ions (Huang and Adamson 1993). Until the point of binding, the structure of the protein is able to fluctuate and is seen throughout the MD simulation. Through analysis of the simulations, certain residues have sparked heightened interest, specifically position six in the alpha helix. Generally occupied by an arginine or lysine, position six is the only residue in the helix that consistently has a highly sampled state that is removed from the state of the crystal structure. Additionally, the negative one position of the  $\alpha$ -helix appears to fluctuate between two states throughout the simulation while taking time to sample the available space in between. The third position of the  $\alpha$ -helix is the one residue that tends to reach an equilibrium state and is generally an amino acid with a smaller side chain.

It is possible that the information gathered on the states of the positions of the individual side chains could be related to the binding mechanism of EGR proteins, which is still yet to be empirically determined. One possible mechanism that comes to light from this research is that the third position of the  $\alpha$ -helix remains steady, while the negative one position of the  $\alpha$ -helix samples the conformational space. When the negative one position residue finds a desirable ligand it would rotate to a position to bring the third position of the  $\alpha$ -helix closer to the DNA and increase the binding potential. After that point, the sixth position of the helix would come into play and reach out to make the final interaction between the protein and the ligand, in this case double stranded DNA.

There are many possible sources that could undermine the binding mechanism hypothesis. The current method for analyzing the predicted structure of homology models is through RMS analysis, which has been criticized as being a measure not specific enough to yield in-depth information about the differences in structure. Referring back to the graphs in Figures 11 – 19, calculating the RMS is sufficient to tell how much difference there is between two

different states, but is not specific enough to be able to cluster the data on this information alone. Because the comparison in these graphs is against a steady state can there be information gained. RMS analysis is also used in the cluster determination for the completed structures of the proteins. Here each time step is compared to every other time step for a given residue, and the time step with the highest amount of neighbors is reported and used in the structure of the completed protein. Since RMS is not being used in a conventional hierarchical or non-hierarchical clustering method, some of the reservations in using RMS as a means of determining structural information can be circumvented.

Another source of concern stems from the neighbor clustering method instantiated in finding the best structure to make the completed structure. Although the neighbor clustering algorithm reports all possible clusters, there is no straightforward way to identify markedly different clusters of a given residue. It is possible to manually search through the high ranking clusters to find disparate clusters, but that is neither time efficient or able to be automated. It is not guaranteed that the second or third highest ranking cluster is any different than the highest ranking cluster, since clusters are not distinct, non-overlapping sets. This problem presents itself when looking at residues with a high amount of clusters accommodating less than a majority of the time steps but receiving similar amounts of neighbors.

There are different approaches that could have been taken for clustering the side chains of the amino acids. One way would have been to analyze the side chains against a rotamer library using RMS, and then using the rotamer with the highest amount of similar structures in the final structure reported (Kimura et al. 2001). Another method would be to find the centroid for each side chain in a given residue and use hierarchical clustering methods in three-dimensional space to find the largest cluster within a given radius. The disadvantage of this clustering method would be that valuable information about the side chain and the direction that it is facing would be lost due to reduction of dimensionality. To alleviate the hindrance in the previous method, one could find the centroid of the last carbon in the side chain and any atoms extending beyond that

point and perform hierarchical clustering. The advantage to this method would be that the subject that is of highest concern, the location of the tip of the side chain, would be forefront consideration of this method. Ultimately, the atoms of the side chain would be neglected, but are of little concern since they are not interacting with the DNA.

### *Future Research*

There are many areas of this project that can be further expanded in addition to ideas previously mentioned. There are two other areas for improvement within the current method for predicting the structure of EGR proteins. The first improvement would be to make the method completely automated so that any user interested in EGR proteins would be able to create their own homology model and perform MD simulations. Automating the procedure would be labor intensive as there are steps in the method that at this point in time require human interaction and require making knowledge-based decisions. The second major area for improvement would be to develop a scoring function for addressing the quality of a homology model, and the completed structure after the simulation and clustering is complete. A scoring algorithm would also be beneficial in the clustering process, potentially ranking rotamers seen throughout the simulation. Residues with a more energetically favorable state would appear higher in the list of clustered residues, thus not relying solely on the RMS calculation between structures.

The next step for the completed predicted structure for zinc finger proteins would involve the use of a derivative of the ClusPro program, initially developed by Comeau *et al.* for predicting interactions between proteins (2004). ClusPro uses two unbound structures and performs automated rigid body docking, and ranks the results based upon clustering of the resolved conformations on the evaluation of the properties of free energy, electrostatic and desolvation calculations. Since this method uses two unbound structures it can also be applied to predict protein – DNA interactions, which is under development in the Camacho laboratory. The initial clustered result for GD-AA with a 0.75Å water layer was tested to see if protein – DNA

interactions could be predicted. The structures for the protein domains were predicted to a degree that they were able to pick up the experimentally proven interactions between 1G2D and the given DNA sequence.

Aside from predicting protein – DNA interactions, there are many more experiments that can be performed using MD simulations. Since we have MD simulation results for the homology modeled proteins without the presence of ions or DNA, the next step would be to add these incrementally to see how the protein reacts to the new conditions. One starting point would be to add a zinc ion into each of the zinc finger domains. It would also be advantageous to simulate the protein in the presence of DNA or a negatively charged object representing DNA both with and without zinc ion to see how the protein thus reacts. Another interesting simulation would be with all three domains together, and allow for the backbone to be unconstrained in the loop regions between the zinc finger domains.

## **Conclusion**

By exploring the EGR family of transcription factors there have been many insights gained into structure prediction for this specific family. Since the EGR family is one of the largest families of transcription factors, the method developed can easily be used for any of the proteins within the family. Using molecular dynamic simulations on singular domains yields a final clustered protein with side chains that have near-native conformations. By analyzing the results and superimposing a DNA helix one can infer information about potential protein – DNA interactions before performing actual interaction algorithms. Analyzing the completed domains before using interaction algorithms can help determine if there are residues that should be further explored. By looking at the results, and the analysis against the crystal structure, a hypothesis for EGR proteins binding to DNA can be surmised: the negative one position of the  $\alpha$ -helix is responsible for initial contact with the DNA, and the other key residues have interactions that are initiated after initial binding.

## References

- Altschul, S.F.; Madden, T.L.; Schäffer A.A.; Zhang J; Zhang Z.; Miller W; Lipman D.J., "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", *Nucleic Acids Res.* 1997, **25**, 3389-3402.
- Bateman, A.; Coin, L.; Durbin, R.; Finn, R.D.; Hollich, V.; Griffiths-Jones, S.; Khanna, A.; Marshall, M.; Moxon, S.; Sonnhammer, E.L.L.; Studholme, D.J.; Yeats, C.; Eddy, S.R., "The Pfam Protein Families Database", *Nucleic Acids Res. Database Issue* 2004, **32**, D138-D141.
- Benos, P.V.; Lapedes, A.S.; Stormo, G.D., "Probabilistic Code for DNA Recognition by Proteins of the EGR Family", *J. Mol. Biol.* 2002, **323**, 701-727.
- Brooks, B.R.; Bruccoleri, R.E.; Olafson, B.D.; States, D.J.; Swaminathan, S.; Karplus, M., "CHARMM: A Program for Macromolecular Energy, Minimization, and Dynamic Calculations", *J. Comput. Chem.* 1983, **4**, 187-217.
- Camacho, C.J., "Modeling Side-Chains Using Molecular Dynamics Improve Recognition of Binding Region in CAPRI Targets", *Proteins* 2005, **60**, 245-251.
- Christy, B.; Nathans, D., "DNA Binding Site of the Growth Factor-Inducible Protein Zif268", *Proc. Natl. Acad. Sci. USA.* 1989, **86**, 8737-8741.
- Comeau, S.R.; Gatchell, D.W.; Vajda, S.; Camacho, C.J., "ClusPro: an automated docking and discrimination method for the prediction of protein complexes", *Bioinformatics* 2004, **20**, 45-50.
- Davis, S.; Bozon, B.; Laroche, S., "How Necessary is the Activation of the Immediate Early Gene *zif268* in Synaptic Plasticity and Learning?", *Behav. Brain Res.* 2003, **142**, 17-30.
- Elrod-Erickson, M.; Rould, M.A.; Nekludova, L.; Pabo, C.O., "Zif268 Protein-DNA Complex Refined at 1.6 Å: A Model System for Understanding Zinc Finger – DNA Interactions", *Structure*, 1996, **4**, 1171-1180.
- Gashler, A.; Sukhatme, V.P., "Early Growth Response Protein 1 (Egr-1): Prototype of a Zinc-finger Family of Transcription Factors", *Prog. Nucleic Acid Res. Mol.Biol.*, 1995, **50**, 191-224.
- Huang, R.P.; Adamson, E.D., "Characterization of the DNA-binding Properties of the Early Growth Response-1 (Egr-1) Transcription Factor: Evidence for Modulation by a Redox Mechanism", *DNA Cell Biol.* 1993, **12(3)**, 265-273.
- Hulo, N.; Sigrist, C.J.A.; Le Saux, V.; Langendijk-Genevaux, P.S.; Bordoli, L.; Gattiker, A.; De Castro, E.; Bucher, P.; Bairoch, A., "Recent improvements to the PROSITE database", *Nucleic Acids Res.* 2004, **32**, 134-137.
- Kim, C.A.; Berg, J.M., "A 2.2 Å Resolution Crystal Structure of a Designed Zinc Finger Bound to DNA", *Nat. Struct. Biol.*, 1996, **3**, 940-945.
- Kimura, S.R.; Brower, R.C.; Vajda, S.; Camacho, C.J., "Dynamical view of the Positions of Key Side Chains in Protein – Protein Recognition", *Biophys J.* 2001, **80**, 635-642.

- Klug, A; Rhodes, D., “ ‘Zinc Fingers’: A Novel Protein Motif for Nucleic Acid Recognition”, *Trends Biochem. Sci.* 1987, **12**, 464 – 469.
- Kozakov, D.; Clodfelter, K.H.; Vajda, S.; Camacho, C.J., “Optimal Clustering for Detecting Near-Native Conformation in Protein Docking”, *Biophys. J.* 2005, **89**, 867-875.
- Lindahl, E.; Hess, B.; Van Der Spoel, D., “GROMACS 3.0: A Package for Molecular Simulation and Trajectory Analysis”, *J. Mol. Model.* 2001, **7**, 306-317.
- Paillard, G.; Deremble, C.; Lavery, R., “Looking into DNA recognition: zinc finger binding specificity”, *Nucleic Acids Res.* 2004, **32**, 6673-6682.
- Prasad, J.C.; Comeau, S.R.; Vajda, S.; and Camacho, C.J., “Consensus alignment for reliable framework prediction in homology modeling”, *Bioinformatics*, 2003, **19**, 1682-1691.
- Sukhatme, V.P.; Cao, X.; Chang, L.C.; Tsai-Morris, C.; Stamenkovich, D.; Ferreira, P.C.P; Cohen, D.R.; Edwards, S.A.; Shows, T.B.; Curran, T.; Le Beau, M.M.; Adamson, E.D., “A Zinc Finger-Encoding Gene Coregulated with *c-fos* during Growth and Differentiation, and after Cellular Depolarization”, *Cell* 1988, **53**, 37-43.
- Thompson, J.D.; Higgins, D.G.; Gibson, T.J., “CLUSTAL W: Improving the Sensitivity of Progressive Multiple Sequence Alignment Through Sequence Weighting, Position-specific Gap Penalties and Weight Matrix Choice”, *Nucleic Acids Res.* 1994, **22**, 4673 – 80.
- Wolf, S.A.; Grant, R.A., Elrod-Erickson, M.; Pabo, C.O., “Beyond the ‘Recognition Code’: Structures of Two Cys<sup>2</sup>Hys<sup>2</sup> Zinc Finger / TATA Box Complexes”, *Structure* 2001, **9**, 717-723.

## Appendix

<b>Appendix A.</b>	Multiple Sequence Alignment of EGR Protein Domains	Page 50
<b>Appendix B.</b>	Consensus Output	Page 63
<b>Appendix C.</b>	Domain Information for Each Homology Modeled EGR Protein	Page 69
<b>Appendix D.</b>	RMS versus Crystal Structure for Clustered Models	Page 70



Q9NWX2\_HUMAN\_81-103  
AIOL\_HUMAN\_174-196  
Q69BL7\_HUMAN\_174-196  
Q69BM3\_HUMAN\_140-162  
IKAR\_HUMAN\_173-195  
Q69BM4\_HUMAN\_173-195  
HELI\_HUMAN\_168-190  
Q6PQC6\_HUMAN\_142-164  
Q6PQC8\_HUMAN\_168-190  
Q6PQD0\_HUMAN\_142-164  
Q96JP3\_HUMAN\_175-197  
Q9H2S9\_HUMAN\_113-135  
Q71UL5\_HUMAN\_4-26  
Q9Y2Y4\_HUMAN\_401-423  
Q6ZMZ9\_HUMAN\_134-156  
Q6ZTI2\_HUMAN\_93-117  
ZIC4\_HUMAN\_220-244  
ZIC1\_HUMAN\_332-356  
Q96IQ9\_HUMAN\_176-198  
Q6PIF0\_HUMAN\_243-265  
Q6DKI8\_HUMAN\_56-80  
O95276\_HUMAN\_64-86  
PRDM6\_HUMAN\_542-564  
Q8WXE2\_HUMAN\_275-299  
Q96MM3\_HUMAN\_275-299  
PRD12\_HUMAN\_271-293  
Q96AM6\_HUMAN\_172-194  
Q96LB6\_HUMAN\_172-194  
Q8TAX0\_HUMAN\_175-197  
ZIC1\_HUMAN\_302-326  
ZIC4\_HUMAN\_190-214  
Q6PJS0\_HUMAN\_212-234  
Q96MU6\_HUMAN\_365-387  
Q6ZS14\_HUMAN\_108-131  
Q7Z7H7\_HUMAN\_473-495  
Q9BU03\_HUMAN\_71-93  
Q9H609\_HUMAN\_71-93  
Q6ZN24\_HUMAN\_186-208  
Q6ZN24\_HUMAN\_370-392  
Q7LDZ1\_HUMAN\_178-200  
Q7LDZ1\_HUMAN\_362-384  
Q7LDZ1\_HUMAN\_5-27  
Q6ZN24\_HUMAN\_80-102  
Q7Z6T2\_HUMAN\_769-793  
Q9NUN9\_HUMAN\_99-123  
Q9NQ73\_HUMAN\_769-793  
Q96PN7\_HUMAN\_1013-1037  
OVOL1\_HUMAN\_118-140  
Q86XL8\_HUMAN\_95-117  
Q5R2W1\_HUMAN\_708-731  
Q6ZMW8\_HUMAN\_890-913  
Q5T149\_HUMAN\_335-357  
Q8IZ20\_HUMAN\_335-357  
Q5T149\_HUMAN\_411-433  
Q8IZ20\_HUMAN\_411-433  
Q9H9H3\_HUMAN\_125-147

FKCS--ICQRHFKNLKT--FVKHQQ--LH---  
FKC--HLCNYACQRRDA--LTGH--LRTH---  
FKC--HLCNYACQRRDA--LTGH--LRTH---  
FKC--HLCNYACQRRDA--LTGH--LRTH---  
FKC--HLCNYACRRRDA--LTGH--LRTH---  
FKC--HLCNYACRRRDA--LTGH--LRTH---  
FKC--PFCSYACRRRDA--LTGH--LRTH---  
FKC--PFCSYACRRRDA--LTGH--LRTH---  
FKC--PFCSYACRRRDA--LTGH--LRTH---  
FKC--PFCSYACRRRDA--LTGH--LRTH---  
FKC--PFCNYACRRRDA--LTGH--LRTH---  
FKC--PFCNYACRRRDA--LTGH--LRTH---  
FEC--KLCHQRSRDYSA--MIKH--LRTH---  
FSC--SLCPQRSRDFSA--MTKH--LRTH---  
FEC--KVCQQAQRQSA--LTVH--KQCH---  
FRCEFEGCERRFANSSD--RKKH--SHVH---  
FRCEFEGCERRFANSSD--RKKH--SHVH---  
FKCEFEGCDRRFANSSD--RKKH--MHVH---  
FKC--ENCLLRFRTHRS--LFKH--LHVC---  
FEC--PNCHERFARNST--LKCH--LTAC---  
FECFVEGCCARFSARSS--LYIH--SKKH---  
FKCSQ--CGRGFVSAGV--LKAH--IRTH---  
FKCGY--CGRAFAGATT--LNNH--IRTH---  
FVCPFQGCNRRFIQSNN--LKAH--ILTH---  
FVCPFQGCNRRFIQSNN--LKAH--ILTH---  
FVCRF--CNRRFSQSST--LRNH--VRLH---  
FICKF--CGRHFTKSYN--LLIH--ERTH---  
FICKF--CGRHFTKSYN--LLIH--ERTH---  
FVCKF--CGRHFTKSYN--LLIH--ERTH---  
FPCFPFGCGKVFARSEN--LKI--KRTH---  
FPCFPFGCGKVFARSEN--LKI--KRTH---  
FECKE--CGRSFRNSSC--LNDH--IQIH---  
FACVV--CGKYFRNSSC--LNNH--VRIH---  
FQCKF--CVRYFRSKNL--LIEHT-RKVH---  
FECKD--CGETFNKSAA--LAEH--RKIH---  
FICFT--CARSFLSSKA--LITH--QRSH---  
FICFT--CARSFPSSKA--LITH--QRSH---  
IRCEF--CGEFFENRKG--LSSH--ARSH---  
IRCEF--CGEFFENRKG--LSSH--ARSH---  
IRCEF--CGEFFENRKG--LSSH--ARSH---  
IRCEF--CGEFFENRKG--LSSH--ARSH---  
TTCEV--CGACFETRKG--LSSH--ARSH---  
MRCDF--CGAGFDTRAG--LSSH--ARAH---  
FICEMPNCGAVFSSRQA--LNGH--ARIH---  
FICEMPNCGAVFSSRQA--LNGH--ARIH---  
FICEMPNCGAVFSSRQA--LNGH--ARIH---  
FICEMPNCGAVFSSRQA--LNGH--ARIH---  
FTCRV--CQKAFTYQRM--LNRH-MKC-H---  
FTCRV--CQKAFTYQRM--LNRH-MKC-H---  
WPCEK--CGKMFTVHKQ--LERH-QEL-LC--  
FICRK--CQMMFTDEDA--AVNH-QKS-FC--  
FQCAL--CQKSFTQLAH--LQKH-HLV-H---  
FQCAL--CQKSFTQLAH--LQKH-HLV-H---  
FQCSV--CRSRFTQHIH--LKLH-HRL-H---  
FQCSV--CRSRFTQHIH--LKLH-HRL-H---  
FSCDI--CGKLFTRREH--VKRH-SLV-H---

\* \* :



Q8N2J5\_HUMAN\_379-401  
Q8N782\_HUMAN\_73-95  
Q8N246\_HUMAN\_335-357  
O14859\_HUMAN\_43-65  
ZNF27\_HUMAN\_1-23  
Q8N4W9\_HUMAN\_53-75  
ZNF18\_HUMAN\_29-51  
Q13580\_HUMAN\_30-52  
Q8N782\_HUMAN\_45-67  
Q9H963\_HUMAN\_62-84  
ZNF79\_HUMAN\_70-92  
O14889\_HUMAN\_12-34  
O43694\_HUMAN\_47-69  
Q15923\_HUMAN\_1-21  
ZN253\_HUMAN\_200-222  
O14891\_HUMAN\_44-66  
Q96MU6\_HUMAN\_279-303  
Q15929\_HUMAN\_9-31  
ZN253\_HUMAN\_172-194  
O43698\_HUMAN\_35-57  
O43698\_HUMAN\_63-85  
Q9H3U2\_HUMAN\_32-54  
Q8N229\_HUMAN\_6-28  
Q8N229\_HUMAN\_34-56  
ZNF15\_HUMAN\_29-51  
ZN253\_HUMAN\_228-250  
Q9H3U2\_HUMAN\_4-26  
Q15929\_HUMAN\_37-59  
Q9H3U2\_HUMAN\_60-82  
ZNF15\_HUMAN\_1-23  
Q13580\_HUMAN\_2-24  
Q13580\_HUMAN\_86-108  
Q8N782\_HUMAN\_17-39  
ZNF27\_HUMAN\_29-51  
Q7Z7K7\_HUMAN\_195-217  
Q8TBC5\_HUMAN\_413-435  
Q9H9A0\_HUMAN\_413-435  
Q9BRK7\_HUMAN\_310-332  
O14860\_HUMAN\_44-66  
Q8N2C6\_HUMAN\_45-67  
O14852\_HUMAN\_2-24  
Q16524\_HUMAN\_11-33  
Q16524\_HUMAN\_67-89  
Q16524\_HUMAN\_39-61  
Q8N2C6\_HUMAN\_17-39  
Q15933\_HUMAN\_33-55  
Q15936\_HUMAN\_39-61  
O14885\_HUMAN\_33-55  
Q8N8C0\_HUMAN\_33-55  
Q6ZMZ9\_HUMAN\_106-128  
Q8N8C0\_HUMAN\_89-111  
Q8N246\_HUMAN\_279-301  
ZN212\_HUMAN\_455-477  
Q8N2J5\_HUMAN\_317-339  
ZN212\_HUMAN\_316-338  
ZNF32\_HUMAN\_1-23

YKCP--ECDSSFSSHKSS--LTKHQITH-----  
YECE--ECDKAFSFKSN--LESHRITH-----  
HECG--ECRKTFSYKSN--LIRHRRVH-----  
YKCK--QCGKGFSSRSA--VNVHCKVH-----  
FKCV--ECGKGFSSRSA--LNVHKKLH-----  
CKCH--QCGKVFSPRSL--LAEHEKIH-----  
CKCD--YCGKGFSDFSG--LRHHEKIH-----  
YKCN--ECGKTFGQNSD--LLIHKSIH-----  
HKCN--ECGKTFSQKSY--LACHRSIH-----  
YKCN--QCGKTFSYKSS--LVIHKAIH-----  
YKCN--ECGKFFSESSA--LIRHHIIH-----  
YRCN--ECGKAFSVRSS--LTTHQAIH-----  
YQCH--ECRKPLVSVSS--LTTHQTIH-----  
--CN--ECGKAFNQSAC--LMQHQRH-----  
YRCE--ECGKAFNQSAN--LTTHKRIH-----  
YRCE--ECGKAFGQSSS--LIHHQRH-----  
KPCELEECGKASPVSSS--LTQHVRH-----  
FKCI--VCGKAFNSSSN--LTTHKKIH-----  
FKCI--ICGKAFKRSST--LTTHKKIH-----  
FKCE--ECGKAFNHPSA--LTTHKFIH-----  
YKCE--ECEKAFNRFSY--LTKHKIIH-----  
YKCE--ECGKAFNRFST--LTKHKRIH-----  
YKCE--ECGKAFSVFST--LTKHKIIH-----  
YKCE--ECGKAFNQSSI--FTKHKIIH-----  
YKCK--ECGKAFNQSST--LMKHKIIH-----  
YKCE--ECGKAFNWSSD--LNKHKKIH-----  
YKCE--ECDRAFSQSSN--LTEHKKIH-----  
YRCE--ECGKAFKRSSH--LTVHKIVH-----  
YKCE--ECGKAFNQSYQ--LTRHKIVH-----  
YKCE--ECGKSFILSSH--LTTHKIIH-----  
YKCK--VCDKAFANNSH--LVRHTRH-----  
YKCK--VCDKAFANNSH--LVSHTSIH-----  
YKCE--ECDKAFRHNSA--LQRHRIH-----  
YNCE--ECGKAFIHDSQ--LQEHQRH-----  
YKCK--ECGKAFFFHSY--LVKHQRH-----  
YACG--ECGEAFAWLSH--LMEHHSSH-----  
YACG--ECGEAFAWLSH--LMEHHSSH-----  
YACG--ECGEAFAWLSH--LMEHHSSH-----  
IECI--ECGKAFNRRSY--RTWHQQRH-----  
YECV--ECGKAFNRRSP--LTRHQRH-----  
HECN--QCGKAFNRSSN--HIHHQKVH-----  
YECN--ACGEAFIRSKS--LARHQVLH-----  
YECS--ECGKAFSRSKC--LIRHQSLH-----  
YKCN--ACGRAFCSNRN--LIDHQRH-----  
YVCI--QCGKAFCRTTN--LIRHFSIH-----  
YECS--QCGKAFRQSTH--LTQHQRH-----  
YECS--QCGKAFRQSTH--LTQHQRH-----  
YECH--QCGKAFSQRAH--LTIHQRIH-----  
YECQ--ICGKPFKRAH--LTQHNRIH-----  
YECS--ECGKVIRKAW--FDQHQRH-----  
YECL--ECRKTFRRSAH--LIRHQRIH-----  
FECS--ECEESFSKKCH--LILHKIIH-----  
YSCT--ECEKSFVQKQH--LLQHQKIH-----  
YECA--ECEISFRHKQQ--LTLHQRIH-----  
YECS--ECEITFRYKQQ--LATHLRSH-----  
YECQ--ECGKSFRQKGS--LTLHERIH-----

\* \* :

Q7Z7K7\_HUMAN\_279-301  
 Q7Z7K7\_HUMAN\_307-329  
 Q8N246\_HUMAN\_251-273  
 O14861\_HUMAN\_41-63  
 ZN154\_HUMAN\_63-85  
 Q7Z3Q9\_HUMAN\_161-183  
 Q15933\_HUMAN\_61-83  
 Q15936\_HUMAN\_67-89  
 Q7KZ25\_HUMAN\_73-95  
 ZK23\_HUMAN\_1-23  
 ZN154\_HUMAN\_35-57  
 ZK23\_HUMAN\_29-51  
 O14885\_HUMAN\_5-27  
 ZN154\_HUMAN\_7-29  
 ZNF17\_HUMAN\_29-51  
 ZNF17\_HUMAN\_1-23  
 Q5T149\_HUMAN\_307-329  
 Q8IZ20\_HUMAN\_307-329  
 ZNF12\_HUMAN\_29-51  
 ZN126\_HUMAN\_29-51  
 ZN126\_HUMAN\_57-79  
 ZN126\_HUMAN\_1-23  
 O14892\_HUMAN\_1-21  
 ZNF21\_HUMAN\_29-51  
 ZNF12\_HUMAN\_1-23  
 O14892\_HUMAN\_27-49  
 O14892\_HUMAN\_55-77  
 ZN174\_HUMAN\_354-376  
 Q5MPB1\_HUMAN\_54-76  
 Q82133\_HUMAN\_41-63  
 ZEP2\_HUMAN\_1186-1208  
 Q99302\_HUMAN\_3-25  
 Q8N8C0\_HUMAN\_61-83  
 HKR2\_HUMAN\_59-81  
 ZN174\_HUMAN\_382-404  
 ZBT12\_HUMAN\_387-409  
 Q8TAX0\_HUMAN\_231-253  
 Q96AM6\_HUMAN\_228-250  
 Q96LB6\_HUMAN\_228-250  
 Q9H963\_HUMAN\_34-56  
 Q9H963\_HUMAN\_90-112  
 Q15921\_HUMAN\_36-58  
 Q96K08\_HUMAN\_183-205  
 O14852\_HUMAN\_30-53  
 Q7KZ25\_HUMAN\_45-67  
 ZN396\_HUMAN\_251-273  
 Q5U5Z4\_HUMAN\_242-264  
 O14887\_HUMAN\_11-33  
 O14888\_HUMAN\_11-33  
 O14890\_HUMAN\_12-34  
 Q15921\_HUMAN\_8-30  
 O14891\_HUMAN\_16-38  
 Q6RFR8\_HUMAN\_138-160  
 Q8N4W9\_HUMAN\_25-47  
 Q15923\_HUMAN\_27-49  
 ZNF26\_HUMAN\_29-51

YECK--ECGKSFTSGST--LNQHQQIH-----  
 YHCK--QCGKSFTVGST--LIRHQQIH-----  
 YECS--KCEKAFTCKNT--LVQHQQIH-----  
 YECT--ECGKSFSVKGK--LIQHQRH-----  
 YECT--ECGKSFSHNSS--LIKHQRH-----  
 YGCT--DCGKAFSHKST--LIKHQRH-----  
 YECN--DCGKAFSHSSS--LTKHQRH-----  
 YECN--DCGNPFSHSSS--LTKHQRH-----  
 YVCN--DCGKAFSQSSS--LIYHQRH-----  
 YECS--ECGKSFRQRSG--LIQHRRH-----  
 YECS--ECGKSFTQNSG--LIKHRRVH-----  
 YECS--ECGKSFSQSAS--LIQHQRVH-----  
 YECE--ECGKEFRHISS--LIAHQRMH-----  
 YECS--ECGKFFPYSSS--LRKHQRVH-----  
 YECS--ECGKFFVDSCT--LKSHQRVH-----  
 YECN--KCGKFFRYCFT--LNRHQRVH-----  
 YECN--ICGKSFGQLSN--LKVHLRVH-----  
 YECN--ICGKSFGQLSN--LKVHLRVH-----  
 YECN--ECGKFFSRLSY--LTVHYRTH-----  
 YECN--QCGKAFSKSHS--LQCHKRTH-----  
 YECN--QCGKAFSQHGL--LQRHKRTH-----  
 YECN--QCGKAFAQHSS--LKCHYRTH-----  
 --CP--ECGKSFCQKVT--STQHQRTH-----  
 YACT--ECGKAFFREKST--FTVHQRTH-----  
 YKCS--ECGKCFCKRST--LTTHLRTH-----  
 YECN--ACGKTFYHKS--LTRHQI IH-----  
 YECY--ECGKTFCLKSD--LTIHQRS-----  
 YTCG--ECGNCFGRQST--LKLHQRH-----  
 YICE--ECGIRCKKPSM--LKKHIRTH-----  
 YICE--ECGIRCKKPSM--LKKHIRTH-----  
 YICE--ECGIRCKKPSM--LKKHIRTH-----  
 YVCE--ECGIRCKKPSM--LKKHIRTH-----  
 YECK--ECGKVFICST--LIQHKRTH-----  
 YMCG--HCGKCFRESSS--LAKHQRVH-----  
 YQCG--QCGKSFRQSSN--LHQHRLH-----  
 HSCG--ICGKCFTQKST--LHDHLNLH-----  
 FKCG--ECGKGFCQSRT--LAVHKT LH-----  
 FKCG--ECGKGFCQSRT--LAVHKT LH-----  
 FKCG--ECGKGFCQSRT--LAVHKT LH-----  
 YKCD--ICGKVFNQKRY--LAYHHRCH-----  
 HKCN--ECGKVFNQKAY--FASHHRLH-----  
 YECE--ECGKSFSRSSH--LAQHQRTH-----  
 YTCE--ECGKAFSRSSF--LVQHQRH-----  
 YTCV--ECGKAFSQSSH--LIQH--IHPH---  
 FKCD--ECGKGFVQSSH--LIQHQRH-----  
 QKCD--ECGKIFSQSSA--LILHQRH-----  
 FKCK--ECLKAFSQSSA--LIQHQRTH-----  
 HTCD--ECGKSFCYISA--LHIHQRVH-----  
 HTCD--ECGKSFCYISA--LHIHQRVH-----  
 HTCD--ECGKSFCYISA--QHIHQRVH-----  
 HKCN--ECGKSFCRLSH--LIQHQRTH-----  
 FKCI--ECGKAFLSSK--LIQHQRH-----  
 HECM--ICGKAFLHSH--LIQHQRH-----  
 HKCD--DCGKAFTSHSH--LVGHQRH-----  
 YTCT--ECGKAFTQNSS--LVEHERTH-----  
 FKCS--ECGKAFTQKSS--LSEHQRVH-----

\* \* :

Q15933\_HUMAN\_5-27  
 Q15936\_HUMAN\_11-33  
 Q96K08\_HUMAN\_155-177  
 O14889\_HUMAN\_40-62  
 O43694\_HUMAN\_19-41  
 Q92669\_HUMAN\_44-66  
 O14857\_HUMAN\_30-52  
 Z297B\_HUMAN\_373-394  
 ZN174\_HUMAN\_326-348  
 O14887\_HUMAN\_39-61  
 O14890\_HUMAN\_40-62  
 O14888\_HUMAN\_39-61  
 O14859\_HUMAN\_15-37  
 O14899\_HUMAN\_19-41  
 Q9UFH1\_HUMAN\_26-48  
 ZNF13\_HUMAN\_29-51  
 ZN229\_HUMAN\_377-399  
 O14899\_HUMAN\_47-69  
 O14899\_HUMAN\_75-97  
 Q9UFH1\_HUMAN\_54-76  
 ZN229\_HUMAN\_349-371  
 ZNF13\_HUMAN\_1-23  
 GLIS1\_HUMAN\_322-346  
 Q5VTL4\_HUMAN\_322-346  
 GLIS3\_HUMAN\_471-495  
 Q5VZV9\_HUMAN\_472-496  
 Q9BZE0\_HUMAN\_293-317  
 Q6ZTI2\_HUMAN\_123-147  
 ZIC4\_HUMAN\_250-274  
 ZIC1\_HUMAN\_362-384  
 Q5U5Z4\_HUMAN\_214-236  
 Q6ZN79\_HUMAN\_172-194  
 Q6ZN79\_HUMAN\_200-222  
 Q8TBC5\_HUMAN\_441-463  
 Q9H9A0\_HUMAN\_441-463  
 Q9BRK7\_HUMAN\_338-360  
 Q96MU6\_HUMAN\_337-359  
 Y0352\_HUMAN\_633-655  
 Q5U5Z4\_HUMAN\_270-292  
 ZN396\_HUMAN\_279-301  
 Q6PJS0\_HUMAN\_240-262  
 Q8NEI5\_HUMAN\_209-231  
 Q6Q7C8\_HUMAN\_195-219  
 Q7RTV3\_HUMAN\_195-219  
 ZNF21\_HUMAN\_1-23  
 Q6ZN79\_HUMAN\_228-250  
 ZN396\_HUMAN\_307-329  
 Q92669\_HUMAN\_16-38  
 ZNF79\_HUMAN\_42-64  
 Q7KZ25\_HUMAN\_17-39  
 Q8TB80\_HUMAN\_18-40  
 O14858\_HUMAN\_11-33  
 O14858\_HUMAN\_39-61  
 Q15930\_HUMAN\_48-70  
 Q8NAJ8\_HUMAN\_36-58  
 Q6RFS2\_HUMAN\_30-52

YGCN--ECGKTFSHSSS--LSQHERTH-----  
 YGCN--ECGKTFSHSSS--LSQHERTH-----  
 WKC�--ECGKTFTQSSS--LTQHQRTH-----  
 YKCN--ECGKVFTQNAH--LANHRRIH-----  
 LLCH--ECGKVFTQNSH--LVRHRGIH-----  
 FTCH--ECGKKFSQNSH--LIKHRRTH-----  
 YECS--DCGKSFTSKSQ--LLVHQPVH-----  
 YPC---QCGKSFTTHKSQ--RDRHMSMH-----  
 YKCD--DCGKSFTWNSE--LKRHKRVH-----  
 IKCD--VCGKEFSQSSH--LQTHQRVH-----  
 YKCD--VCGKEFSQSSR--LQTHQRVH-----  
 YKCY--VCGKEFSQSSH--LQTHQRVH-----  
 LKCD--ECGKEFSQGAH--LQTHQKVH-----  
 FKCE--ECGKEFSWSAG--LSAHQRVH-----  
 YKCE--VCGKGFQRSN--LQAHQRVH-----  
 YKCE--ECGKGFSRASN--LLAHQRGH-----  
 YKCE--ECGKAFGRSSN--LLVHQRVH-----  
 YTCQ--QCGKGFQASH--FHTHQRVH-----  
 YICD--VCKKGFQRSH--LIYHQRVH-----  
 YKCD--ACGKGFRWSSG--LLIHQRVH-----  
 YRCD--VCGKGFQYKSV--LLIHQGVH-----  
 YQCD--ACGKGFSRSSD--FNIHFRVH-----  
 YACQIPGCSKRYTDPSS--LRKHVKAH-----  
 YACQIPGCSKRYTDPSS--LRKHVKAH-----  
 YACQIPGCTKRYTDPSS--LRKHVKAH-----  
 YACQIPGCTKRYTDPSS--LRKHVKAH-----  
 YYCKMPGCHKRYTDPSS--LRKHIAH-----  
 YTCKVRGCDKCYTHPSS--LRKHMKVH-----  
 YTCKVRGCDKCYTHPSS--LRKHMKVH-----  
 YLCKM--CDKSYTHPSS--VRKHMKVH-----  
 YKCSH--CEKAFIHNSH--LRRHQKNH-----  
 YQCNL--CEKAYTNCFR--LRRHKMTH-----  
 YACHL--CGKAFTQCSH--LRRHEKTH-----  
 YACQG--CWKTFHFSLA--LAEHQKTH-----  
 YACQG--CWKTFHFSLA--LAEHQKTH-----  
 YACQG--CWKTFHFSLA--LAEHQKTH-----  
 YECKD--CGKACGGFYL--LNEHGKTH-----  
 YQCKV--CHKFFRGRST--IKCHLKTH-----  
 YIC--KECGKAFSHSAS--LCKHLRTH-----  
 YAC--DECAKAFSRSI--LIQHRRTH-----  
 HKC--TYCGKAFTRSTQ--LTEHVRTH-----  
 YTC--RYCGKIFPRSAN--LTRHLRTH-----  
 YLCDYPDCGKAFVQSGQ--LKTHQRLH-----  
 YLCDYPDCGKAFVQSGQ--LKTHQRLH-----  
 YEC--PVCWKAFSQKSQ--LIHQHQRTH-----  
 YKC--HQCGKAFIQSFN--LRRHERTH-----  
 YKC--HDCGKAFSQSSN--LFRHRKRH-----  
 YKC--SWCGKSFSQNTN--LHTHQHQRTH-----  
 YKC--SECCKAFSOSTN--LIHQKTH-----  
 YKC--NKCTKAFGCSSR--LIRHQHQRTH-----  
 YSC--SACGKCFGGSGD--LRRHVRTH-----  
 YECKQ--CSKAFPVYSS--YLRHEKIH-----  
 YECKQ--CSKAFPVYSS--YLRHERTH-----  
 YECKE--CGKAFSSFKY--FCRHERTH-----  
 YECKQ--CGKAFVSFTS--FRYHERTH-----  
 FECKR--CGKAFRSSSS--FRLHERTH-----

\* \* \*







KLF1\_HUMAN\_339-361  
Q6PIJ5\_HUMAN\_339-361  
Q5T3J9\_HUMAN\_406-428  
Q8N717\_HUMAN\_481-503  
KLF4\_HUMAN\_447-469  
KLF2\_HUMAN\_332-354  
KLF5\_HUMAN\_433-455  
KLF6\_HUMAN\_260-282  
Q7Z3W8\_HUMAN\_35-57  
KLF7\_HUMAN\_279-301  
Q7Z3H8\_HUMAN\_246-268  
Q5XG86\_HUMAN\_466-490  
ZBTB7\_HUMAN\_466-490  
O73453\_HUMAN\_448-472  
KLF15\_HUMAN\_381-403  
Q8NDX6\_HUMAN\_157-179  
Q8TB80\_HUMAN\_46-68  
Q5XG86\_HUMAN\_410-432  
ZBTB7\_HUMAN\_410-432  
O73453\_HUMAN\_392-414  
Z297B\_HUMAN\_428-450  
Q6ZN18\_HUMAN\_45-70  
Q96BG3\_HUMAN\_53-78  
Q8TAX0\_HUMAN\_203-225  
Q96AM6\_HUMAN\_200-222  
Q96LB6\_HUMAN\_200-222  
Q8IUL5\_HUMAN\_310-333  
Q8N2Y5\_HUMAN\_120-143  
Q96BR9\_HUMAN\_310-333  
Q96BX0\_HUMAN\_105-128  
ZBTB8\_HUMAN\_369-392  
GLIS1\_HUMAN\_292-316  
GLIS3\_HUMAN\_441-465  
Q5VTL4\_HUMAN\_292-316  
Q5VZV9\_HUMAN\_442-466  
Q9BZE0\_HUMAN\_263-287  
Q6N043\_HUMAN\_418-441  
Q9H740\_HUMAN\_254-277  
Q6MZM6\_HUMAN\_464-487  
Q6N085\_HUMAN\_418-441  
Q9HCI8\_HUMAN\_434-457  
Q9NXS0\_HUMAN\_122-145  
ZNF18\_HUMAN\_1-23  
Q6PIF0\_HUMAN\_215-237  
Q6PIF0\_HUMAN\_279-301  
Q9HD72\_HUMAN\_383-405  
Q9Y2Y4\_HUMAN\_373-395  
Y0352\_HUMAN\_605-627  
O14861\_HUMAN\_13-35  
Q6PJS0\_HUMAN\_128-150  
Q5T967\_HUMAN\_241-263  
Q9BYX9\_HUMAN\_4-26  
Z297B\_HUMAN\_400-422  
Q5T7W1\_HUMAN\_115-137  
Q5T7W2\_HUMAN\_115-137  
Q6ZT53\_HUMAN\_115-137

FRCQ--LCPRAFSRSDH--LALHMKR-H----  
FRCQ--LCPRAFSRSDH--LALHMKR-H----  
FQCQ--KCDRAFSRSDH--LALHMKR-H----  
FQCQ--KCDRAFSRSDH--LALHMKR-H----  
FQCQ--KCDRAFSRSDH--LALHMKR-H----  
FQCH--LCDRAFSRSDH--LALHMKR-H----  
FQCG--VCNRSFSRSDH--LALHMKR-H----  
FKCS--HCDRCFSRSDH--LALHMKR-H----  
FKCS--HCDRCFSRSDH--LALHMKR-H----  
FKCN--HCDRCFSRSDH--LALHMKR-H----  
FKCN--HCDRCFSRSDH--LALHMKR-H----  
YQCD--SCCKTFVRSRDH--LHRHLKK-DGC--  
YQCD--SCCKTFVRSRDH--LHRHLKK-DGC--  
YQCE--FCYKSFTRSDH--LHRHIKR-QSC--  
YQCP--VCEKKFARSRDH--LSKHIV-H----  
YQCE--RCHQCFSTRDR--LLRHKRM-C----  
YTCE--ICNKCFTRSAV--LRRHKKM-H----  
YECN--ICKVRFTRQDK--LKVHMRK-H----  
YECN--ICKVRFTRQDK--LKVHMRK-H----  
YMCT--ICEVRFTRQDK--LKIHMVK-H----  
YECN--ICAKRFMWRDS--FHRHVTS-C----  
YNCCWDQCQACFNSSPD--LADHIRS-IH---  
YNCCWDQCQACFNSSPD--LADHIRS-IH---  
YTC--DICHKAFRRQDH--LRDH-RY-IH---  
YTC--DICHKAFRRQDH--LRDH-RY-IH---  
YTC--DICHKAFRRQDH--LRDH-RY-IH---  
YPC--QACGKRFSRLDH--LSSHFR-IH---  
YPC--QACGKRFSRLDH--LSSHFR-IH---  
YPC--QACGKRFSRLDH--LSSHFR-IH---  
YPC--QACGKRFSRLDH--LSSHFR-IH---  
YPC--ETCGKRFRQEH--LRSHALS-VH---  
YLCQHPGCQKAFSNSSD--RAKHQRT-H----  
YLCQHPGCQKAFSNSSD--RAKHQRT-H----  
YLCQHPGCQKAFSNSSD--RAKHQRT-H----  
YLCQHPGCQKAFSNSSD--RAKHQRT-H----  
YVCPYEGCNKRYSNSSD--RFKHTRT-H----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
YVCQV--CNYRSSFSD--VETHFRSH----  
FQCTI--CKKAFLRSSD--FVKHQRT-H----  
HVCQY--CEKQFDHFGH--FKEHLRK-H----  
YECQV--CNSVFNSWDQ--FKDHLVI-H----  
YSCPV--CGLRFKRD--MSYHVS-H----  
YACSV--CGKRFSLKHQ--METHYRV-H----  
YSCKV--CGKRFHTSE--FNYHRI-H----  
YQCK--ECGKGFNNNTK--LIQHORI-H----  
YKCK--ECGKGFKYFAS--LDNHMGI-H----  
FKCT--ECGKAFKYKHH--LKEHLRI-H----  
YPCM--ICGKKFKSRGF--LKRHMKN-H----  
YCGG--VCGKKFKMKHH--LVGHMKI-H----  
YECG--ICGKKYKYNYNC--FQTHVRA-H----  
YECG--ICGKKYKYNYNC--FQTHVRA-H----  
YECG--ICGKKYKYNYNC--FQTHVRA-H----

\* \* :

Q5T7W1\_HUMAN\_156-178  
Q5T7W2\_HUMAN\_156-178  
Q6ZT53\_HUMAN\_156-178  
Q5T7W1\_HUMAN\_244-266  
Q5T7W2\_HUMAN\_232-254  
Q6ZT53\_HUMAN\_224-246  
Q6ZMZ8\_HUMAN\_41-63  
Q8WXE2\_HUMAN\_217-239  
Q96MM3\_HUMAN\_217-239  
Q9NWX2\_HUMAN\_27-49  
ZN469\_HUMAN\_156-178  
ZBT12\_HUMAN\_359-381  
ZN580\_HUMAN\_120-142  
Q658W5\_HUMAN\_197-219  
ZBTB2\_HUMAN\_254-276  
ZBT25\_HUMAN\_236-258  
PRD12\_HUMAN\_243-265  
Q658W5\_HUMAN\_306-328  
ZBTB2\_HUMAN\_363-385  
ZBT25\_HUMAN\_347-369  
Q8WV14\_HUMAN\_183-205  
SP2\_HUMAN\_518-542  
O00110\_HUMAN\_192-214  
ZFHX2\_HUMAN\_624-646  
ZN580\_HUMAN\_150-172  
Q68DQ8\_HUMAN\_993-1016  
Q9H937\_HUMAN\_581-604  
Q5VWB8\_HUMAN\_581-604  
ZN592\_HUMAN\_1013-1036  
Q9UPR8\_HUMAN\_711-734  
Q8NEI5\_HUMAN\_237-260  
PRDM6\_HUMAN\_514-536  
Q8TBE5\_HUMAN\_138-161  
Q9H5V7\_HUMAN\_138-161  
Q9H2T0\_HUMAN\_138-161  
Q9NZF0\_HUMAN\_158-181  
ZN644\_HUMAN\_586-609  
Q7Z5E7\_HUMAN\_102-124  
Q5MPB1\_HUMAN\_82-106  
Q82133\_HUMAN\_69-93  
Q99302\_HUMAN\_31-55  
ZEP2\_HUMAN\_1214-1238  
Q9NV14\_HUMAN\_229-253  
ZN644\_HUMAN\_410-432  
Q63HK5\_HUMAN\_1041-1064  
Q9H0G6\_HUMAN\_858-881  
Q9P254\_HUMAN\_1039-1062  
Q6DKI8\_HUMAN\_89-114  
ZBT12\_HUMAN\_415-438  
Q6ZN29\_HUMAN\_3-26  
Q9UBK3\_HUMAN\_85-108  
Q5VWB8\_HUMAN\_415-437  
Q9H937\_HUMAN\_415-437  
Q68DQ8\_HUMAN\_827-849  
Q8N8K7\_HUMAN\_291-313  
Q96T92\_HUMAN\_292-314

YTCD--ICGKKYKYYS--FQEHRDL-H----  
YTCD--ICGKKYKYYS--FQEHRDL-H----  
YTCD--ICGKKYKYYS--FQEHRDL-H----  
YTCE--FCGKQYKYYP--YQEHAL-H----  
YTCE--FCGKQYKYYP--YQEHAL-H----  
YTCE--FCGKQYKYYP--YQEHAL-H----  
YTCE--FCGKQYKYYP--YQEHAL-H----  
YKCD--QCGYLSKTANK--LIEHVRV-H----  
HVCA--ECGKAFVESSK--LKRHFLVH-----  
HVCA--ECGKAFVESSK--LKRHFLVH-----  
YVCN--ICFKHFETPSK--LARHYLIH-----  
RDCH--HCGKRFPKPKF--LQRHLAVH-----  
FMCP--RCGKQFNHSSN--LNRHMNVH-----  
FTCG--ACGKAFKRSSH--LSRHRATH-----  
YACH--LCGRRFTLRSS--LREHLQIH-----  
YACH--LCGRRFTLRSS--LREHLQIH-----  
HLCH--YCGERFDSRSN--LRQHLHTH-----  
MRCV--ICHRGFNSRSN--LRSHMRIH-----  
YECT--ICGRKFIQKSH--WREHMYIH-----  
YECT--ICGRKFIQKSH--WREHMYIH-----  
MSCT--ICGHKFPKRSQ--LLEHMYTH-----  
HVC--TDCGRRFTYPSL--LVSHRR-MH----  
HVCHIPDCGKTFRKTSL--LRAHVR-LH----  
HVC--EDCGFTSSRPDT--YAQHRA-LH----  
HTC--DQCAISFSSQDL--LTSHRR-LH----  
HTC--PLCPRRFQDAAE--LAQHVR-LH----  
FPC--RLCERSFCSAPS--LRRHVRVNH----  
FPC--RLCERSFCSAPS--LRRHVRVNH----  
FPC--RLCERSFCSAPS--LRRHVRVNH----  
YPC--RQCEQSFTHPNS--LRKHIRNNH----  
YQC--KQCEESFHYKSQ--LRNHEREQH----  
YRC--KYCDRSFSISSN--LQRHVRNIH----  
YQC--GHCSQSFSQPSE--LRNHV-VTH----  
YKC--ELCSFRCSDRSN--LSHHRRRKH----  
YKC--ELCSFRCSDRSN--LSHHRRRKH----  
YKC--ELCSFRCSDRSN--LSHHRRRKH----  
YIC--KMCPFTTSAKSV--LKKHTEYLH----  
YIC--KMCPFTTSAKSV--LKKHTEYLH----  
LKC--KVCLRPFGDPSN--LNKHIR-LH----  
YHCTY--CNFSFKTKGN--LTKHMKSKAH---  
YHCTY--CNFSFKTKGN--LTKHMKSKAH---  
YVCKH--CHFAFKTKGN--LTKHMKSKAH---  
YVCKL--CNFAFKTKGN--LTKHMKSKAH---  
YFCLHFNCNESFKLPFQ--LAQHTKS--H---  
YPCTK--CNVNFREKKH--LHRHMMY--H---  
YQCKL--CNRTFASKHA--VKLHLSKT-H---  
YQCKL--CNRTFASKHA--VKLHLSKT-H---  
YQCKL--CNRTFASKHA--VKLHLSKT-H---  
SRCPVSTCNRLFTSKHS--MKAHMVRQ-H---  
YRCSY--CDVRFAHKPA--IRRHLKEQ-H---  
YQCKK--CNVVFPRIFD--LITHQKKQ-C---  
YRCED--CDQLFESKAE--LADHQKFP-C---  
YKCAM--CDTVFTHKPL--LSSHFDQH-----  
YKCAM--CDTVFTHKPL--LSSHFDQH-----  
YKCAM--CDTVFTHKPL--LSSHFDQH-----  
YRCPE--CDKVFSPAN--LASHRRWH-----  
YRCPE--CDKVFSPAN--LASHRRWH-----

\* \* :

Q96Q84\_HUMAN\_291-313  
 ZN580\_HUMAN\_92-114  
 Q96JV0\_HUMAN\_737-759  
 Q96JV0\_HUMAN\_896-918  
 Q6Q7C8\_HUMAN\_167-189  
 Q7RTV3\_HUMAN\_167-189  
 ZNF26\_HUMAN\_1-23  
 Q71UL5\_HUMAN\_32-54  
 Q9Y2Y4\_HUMAN\_428-450  
 Q6ZN29\_HUMAN\_184-206  
 FOG1\_HUMAN\_290-314  
 FOG2\_HUMAN\_296-320  
 Q6ZMW8\_HUMAN\_162-185  
 Q6ZN29\_HUMAN\_366-389  
 Q8TBE5\_HUMAN\_110-132  
 Q9H2T0\_HUMAN\_110-132  
 Q9H5V7\_HUMAN\_110-132  
 POGZ\_HUMAN\_494-516  
 Q5SZS4\_HUMAN\_494-516  
 Q5SZS3\_HUMAN\_439-461  
 Q5SZS1\_HUMAN\_399-421  
 Q5SZS2\_HUMAN\_441-463  
 Q9UPR8\_HUMAN\_683-705  
 POGZ\_HUMAN\_619-641  
 Q5SZS1\_HUMAN\_524-546  
 Q5SZS2\_HUMAN\_575-597  
 Q5SZS3\_HUMAN\_564-586  
 Q5SZS4\_HUMAN\_619-641  
 Q7Z6T2\_HUMAN\_351-373  
 Q7Z6T3\_HUMAN\_351-373  
 Q9NQ73\_HUMAN\_351-373  
 Q9NQ72\_HUMAN\_351-373  
 Q96PN7\_HUMAN\_512-534  
 Q9H0G6\_HUMAN\_92-116  
 Q9P254\_HUMAN\_273-297  
 Q63HK5\_HUMAN\_275-299  
 Q96MY0\_HUMAN\_186-208  
 Q9UBK3\_HUMAN\_139-161  
 Y0352\_HUMAN\_508-530  
 AIOL\_HUMAN\_202-224  
 Q69BM1\_HUMAN\_59-81  
 Q69BL7\_HUMAN\_202-224  
 HELI\_HUMAN\_196-219  
 Q6PQC8\_HUMAN\_196-219  
 Q6PQD0\_HUMAN\_170-193  
 IKAR\_HUMAN\_201-224  
 Q96JP3\_HUMAN\_208-231  
 Q9H2S9\_HUMAN\_146-169  
 Q6ZN20\_HUMAN\_244-267  
 ZN409\_HUMAN\_453-476  
 Q5T967\_HUMAN\_171-194  
 Q8IUL5\_HUMAN\_282-304  
 Q8N2Y5\_HUMAN\_92-114  
 Q96BR9\_HUMAN\_282-304  
 Q96BX0\_HUMAN\_77-99  
 Q9H9H3\_HUMAN\_97-119

YRCPE--CDKVFSPAN--LASHRRWH-----  
 YSCPE--CARVFASPLR--LQSHRVSH-----  
 YQCKH--CDSKLQSTAE--LTSHLNIH-----  
 YHCSQ--CDRVLMSMQG--LRSHERSH-----  
 IRCNI--CNRVFPREKS--LQAHKRTH-----  
 IRCNI--CNRVFPREKS--LQAHKRTH-----  
 YECNE--CEKAYPRKAS--LQIHQKTH-----  
 YQCTI--CTEYCPSSLSS--MQKHMKGH-----  
 YRCSL--CGAGCPSLAS--MQAHMRGH-----  
 YQCDQ--CTVAFPTLEL--WQEHQHVH-----  
 RVCPPFPQCRKSCPSASS--LEIHMSR--H---  
 SLCPFPQCTKSFSNARA--LEMHLNS--H---  
 KRCPF--CRALFKAKSA--LESHIRSR-H---  
 KRCPF--CRALFKAKSA--LESHIRSR-H---  
 HRCHL--CPFASAYERH--LEAHMRS--H---  
 HRCHL--CPFASAYERH--LEAHMRS--H---  
 HRCHL--CPFASAYERH--LEAHMRS--H---  
 FRCPH--CTKRLKNNIR--FMNHMKH--H---  
 FRCPH--CTKRLKNNIR--FMNHMKH--H---  
 FRCPH--CTKRLKNNIR--FMNHMKH--H---  
 FRCPH--CTKRLKNNIR--FMNHMKH--H---  
 YQCPH--CEHIADNSKD--LESHMIH--H---  
 LLCPY--CLKVFKNGNA--FQQHYMR--H---  
 LLCPY--CLKVFKNGNA--FQQHYMR--H---  
 LLCPY--CLKVFKNGNA--FQQHYMR--H---  
 LLCPY--CLKVFKNGNA--FQQHYMR--H---  
 LLCPY--CLKVFKNGNA--FQQHYMR--H---  
 LTCSI--CLKEFKNLPA--LNGHMRS--H---  
 LTCSI--CLKEFKNLPA--LNGHMRS--H---  
 LTCSI--CLKEFKNLPA--LNGHMRS--H---  
 LTCSI--CLKEFKNLPA--LNGHMRS--H---  
 LKCMY--CGHSFESLQD--LSVHMIKTKH---  
 LKCMY--CGHSFESLQD--LSVHMIKTKH---  
 LKCMY--CGHSFESLQD--LSVHMIKTKH---  
 LKCSI--CGHLFSSCS--LEKHAES--H---  
 QECKE--CDQVFPDLQS--LEKHMLS--H---  
 FSCSV--CANSFVDWHL--LEKHMAV--H---  
 YKCE--FCGRSYQRSS--LEEHEKER-C----  
 YKCE--FCGRSYQRSS--LEEHEKER-C----  
 YKCE--FCGRSYQRSS--LEEHEKER-C----  
 HKCN--YCGRSYQRSS--LEEHEKER-CH---  
 HKCN--YCGRSYQRSS--LEEHEKER-CH---  
 HKCN--YCGRSYQRSS--LEEHEKER-CH---  
 HKCG--YCGRSYQRSS--LEEHEKER-CH---  
 YKCN--YCGRSYKQOST--LEEHEKER-CH---  
 YKCN--YCGRSYKQOST--LEEHEKER-CH---  
 LKCP--KCNWHYKYQQT--LEAHMKE-KH---  
 LKCP--KCNWHYKYQQT--LDVHMQE-KH---  
 LTCP--YCDRGYKRFTS--LKEHIKY-RH---  
 FKCP--YCTHVVKRKAD--LKRHLR--CH---  
 FKCP--YCTHVVKRKAD--LKRHLR--CH---  
 FKCP--YCTHVVKRKAD--LKRHLR--CH---  
 FKCP--YCTHVVKRKAD--LKRHLR--CH---  
 LKCP--HCSYVAKYRRT--LKRHLL--IH---

\* \* \*

ZBT10\_HUMAN\_698-720  
 ZBTB8\_HUMAN\_341-363  
 ZFHX2\_HUMAN\_103-127  
 ZFHX2\_HUMAN\_46-70  
 PRD12\_HUMAN\_299-323  
 Q6ZN20\_HUMAN\_299-323  
 ZN409\_HUMAN\_508-532  
 Q6ZN20\_HUMAN\_367-391  
 Q9BYX9\_HUMAN\_35-57  
 Q6ZS14\_HUMAN\_600-623  
 Q9UPR8\_HUMAN\_655-677  
 Q5T1J7\_HUMAN\_442-465  
 Q8NAR2\_HUMAN\_13-36  
 Q6T3A3\_HUMAN\_363-386  
 Q6ZN30\_HUMAN\_441-464  
 BNC1\_HUMAN\_357-380  
 Q6T3A3\_HUMAN\_957-980  
 Q9NXV0\_HUMAN\_325-348  
 Q9H6J0\_HUMAN\_251-274  
 Q6ZN30\_HUMAN\_1035-1058  
 BNC1\_HUMAN\_928-951  
 Q5T1J7\_HUMAN\_834-857  
 Q9NXV0\_HUMAN\_123-146  
 Q6T3A3\_HUMAN\_755-778  
 Q6ZN30\_HUMAN\_833-856  
 Q8NAR2\_HUMAN\_405-428  
 Q9H6J0\_HUMAN\_49-72  
 BNC1\_HUMAN\_720-743  
 Q9C0D4\_HUMAN\_136-159  
 OVOL1\_HUMAN\_174-197  
 Q86XL8\_HUMAN\_151-174  
 O00110\_HUMAN\_153-176  
 Q7M4M1\_HUMAN\_36-59  
 Q7Z5E7\_HUMAN\_131-154  
 Q96MY0\_HUMAN\_528-551  
 FOG1\_HUMAN\_241-264  
 ZN592\_HUMAN\_1153-1176  
 Q86XK6\_HUMAN\_120-143  
 ZN512\_HUMAN\_197-220  
 Q86XK6\_HUMAN\_210-233  
 ZN512\_HUMAN\_287-310  
 Q86XK6\_HUMAN\_363-386  
 ZN512\_HUMAN\_440-463  
 Q8N8K7\_HUMAN\_525-548  
 Q96Q84\_HUMAN\_525-548  
 Q96T92\_HUMAN\_526-549  
 Q96MY0\_HUMAN\_705-728  
 Q8NDX6\_HUMAN\_101-123  
 Q96JV0\_HUMAN\_959-982

LKCP--HCSYVAKYRRT--LKRHLL--IH---  
 HKCP--FCPYTAKQKGI--LKRHIR--SH---  
 FKCT--VCRVSYNQSST--LEIHMRSVLH---  
 YKCT--VCKESFTQKNI--LLVHYNSVSH---  
 YKCQ--VCQSAYSQLAG--LRAHQKSARH---  
 FRCE--VCNYSTTTTKGN--LSIHMQSDKH---  
 YRCD--VCNYSTTTTKGN--LSIHMQSDKH---  
 WRCE--VCDYETNVARN--LRIHMTSEKH---  
 YRCT--DCDYTTNKKIS--LHNHLES--H---  
 YKCT--MCNYSTTTTLKG--LRVHQQH-KH---  
 YRCR--LCHYTSNGN-KG--YIKQHLR-VH---  
 VFCN--ACGK-TFYDKG-TLKIHYNA-VH---  
 VFCN--ACGK-TFYDKG-TLKIHYNA-VH---  
 VFCN--ACGK-TFYDKG-TLKIHYNA-VH---  
 VFCN--ACGK-TFYDKG-TLKIHYNA-VH---  
 VFCT--ACEK-TFYDKG-TLKIHYNA-VH---  
 IMCN--ICLK-MYSNKG-TLRVHYKT-VH---  
 IMCN--ICLK-MYSNKG-TLRVHYKT-VH---  
 IMCN--ICLK-MYSNKG-TLRVHYKT-VH---  
 IMCN--ICLK-MYSNKG-TLRVHYKT-VH---  
 ITCH--LCQK-TYSNKG-TFRAHYKT-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 KICY--VCKK-SFKSSY-SVKLHYRN-VH---  
 FQCD--ICKK-TFKNAC-SVKIHHKN-MH---  
 YQCE--YCDYGAIRNDY-IVK-HTKR-VH---  
 YKCS--LCDK-AFTQRC-SLESHLKK-IH---  
 YKCS--LCDK-AFTQRC-SLESHLKK-IH---  
 FRCS--ACGK-AFTQRC-SLEAHLAK-VH---  
 YKCK--DCGN-AFIWRA-SLQYHVKK-VH---  
 YRCE--FCGK-VLVRRR-DLERHVKS-RH---  
 CACT--DCGQ-VATNRT-DLEIHVKR-CH---  
 FPCK--DCGIWYRSERN--LQAHLLY-YC---  
 AQCL--LCGLCYTSASS--LSRHLEFI-VH---  
 FTCH--HCGKQLRSLAG--MKYHVMA-NH---  
 FTCH--HCGKQLRSLAG--MKYHVMA-NH---  
 LKCH--HCGKPYRSKAG--LAYHLRS-EH---  
 LKCH--HCGKPYRSKAG--LAYHLRS-EH---  
 YKCL--LCQKEFVSESG--VKYHINS-VH---  
 YKCL--LCQKEFVSESG--VKYHINS-VH---  
 FSCK--HCPSTFFSSPG--LTRHINK-CH---  
 FSCK--HCPSTFFSSPG--LTRHINK-CH---  
 FSCK--HCPSTFFSSPG--LTRHINK-CH---  
 FQCK--KCFYKTRSTV--LTRHIKL-RH---  
 FVCE--HCFGAFRSSYH--LKRHIL--IH---  
 FRCK--LC--SFKSSYNSRLKTHILK-AH---

\* \* :

## Appendix B – Consensus Output

Target GD-AA : Template specified/identified is 1AAY Chain A

Columns 1-2 correspond to the target.

Columns 3-4 correspond to the template.

Column 5 is the confidence (0-9) for that pair of aligned residues according to sequence alignments

Column 6 is the selection status.

An 'S' indicates that the corresponding pair of residues has been confidently aligned and may be used in homology modeling.

A '.' indicates otherwise.

Please disregard column 7. This is only for diagnostic purposes and will soon be removed.

----- GD-AA.consensus -----

R 3	R 1	9	.	H
P 4	P 2	9	.	
Y 5	Y 3	9	.	
A 6	A 4	9	.	
C 7	C 5	9	.	
P 8	P 6	9	.	
V 9	V 7	9	.	
E 10	E 8	6	.	
S 11	S 9	7	S	D
C 12	C 10	9	S	D
D 13	D 11	9	S	D
R 14	R 12	9	S	D
R 15	R 13	9	S	D
F 16	F 14	9	S	D
S 17	S 15	9	S	D
Q 18	R 16	9	S	D
K 19	S 17	9	S	D
T 20	D 18	9	S	D
N 21	E 19	9	S	A
L 22	L 20	9	S	A
D 23	T 21	9	S	E
T 24	R 22	9	S	E
H 25	H 23	9	S	AE
I 26	I 24	9	S	AE
R 27	R 25	9	S	AE
I 28	I 26	9	S	E
H 29	H 27	9	S	E
T 30	T 28	9	S	E
G 31	G 29	9	S	E
Q 32	Q 30	9	S	E
K 33	K 31	9	S	E
P 34	P 32	9	S	E
F 35	F 33	9	S	E
Q 36	Q 34	9	S	E
C 37	C 35	9	S	E
R 38	R 36	9	S	E
I 39	I 37	9	S	E
C 40	C 38	9	S	E

M 41	M 39	9	S	E
R 42	R 40	9	S	E
N 43	N 41	9	S	E
F 44	F 42	9	S	AE
S 45	S 43	9	S	AE
Q 46	R 44	9	S	E
H 47	S 45	9	S	E
T 48	D 46	9	S	E
G 49	H 47	9	.	G
L 50	L 48	9	.	DG
N 51	T 49	9	.	DG
Q 52	T 50	9	.	AG
H 53	H 51	9	.	AG
I 54	I 52	9	.	EG
R 55	R 53	9	.	EG
T 56	T 54	9	.	EG
H 57	H 55	9	.	EG
T 58	T 56	9	.	EG
G 59	G 57	9	.	EG
E 60	E 58	9	.	EG
K 61	K 59	9	.	EG
P 62	P 60	9	.	EG
F 63	F 61	9	.	EG
A 64	A 62	9	.	EG
C 65	C 63	9	.	EG
D 66	D 64	9	.	EG
I 67	I 65	9	.	EG
C 68	C 66	9	.	EG
G 69	G 67	9	.	EG
R 70	R 68	9	.	EG
K 71	K 69	9	.	EG
F 72	F 70	9	.	AEG
A 73	A 71	9	.	AEG
T 74	R 72	9	.	EG
L 75	S 73	9	.	EG
H 76	D 74	9	.	EG
T 77	E 75	9	.	EG
R 78	R 76	9	.	EGI
D 79	K 77	9	.	EGI
R 80	R 78	9	.	EGI
H 81	H 79	9	.	EGI
T 82	T 80	9	.	EGI
K 83	K 81	9	.	EGI
I 84	I 82	9	.	EGI
H 85	H 83	9	.	EGI
L 86	L 84	9	.	EGI
R 87	R 85	9	.	EGI

Target 1AAY : Template specified/identified is 1G2D Chain C

Columns 1-2 correspond to the target.

Columns 3-4 correspond to the template.

Column 5 is the confidence (0-9) for that pair of aligned residues according to sequence alignments

Column 6 is the selection status.

An 'S' indicates that the corresponding pair of residues has been confidently aligned and may be used in homology modeling.

A '.' indicates otherwise.

Please disregard column 7. This is only for diagnostic purposes and will soon be removed.

----- AA-GD.consensus -----

M 1	M 1	5	.	H
E 2	E 2	7	.	DH
R 3	R 3	9	.	DH
P 4	P 4	9	S	D
Y 5	Y 5	9	S	D
A 6	A 6	9	S	D
C 7	C 7	9	S	D
P 8	P 8	9	S	D
V 9	V 9	9	S	D
E 10	E 10	9	S	D
S 11	S 11	9	S	D
C 12	C 12	9	S	D
D 13	D 13	9	S	D
R 14	R 14	9	S	D
R 15	R 15	9	S	D
F 16	F 16	9	S	D
S 17	S 17	9	S	D
R 18	Q 18	9	S	D
S 19	K 19	9	S	D
D 20	T 20	9	S	D
E 21	N 21	9	S	D
L 22	L 22	9	S	D
T 23	D 23	9	S	D
R 24	T 24	9	S	A
H 25	H 25	9	S	A
I 26	I 26	9	S	A
R 27	R 27	9	S	A
I 28	I 28	9	S	E
H 29	H 29	9	S	E
T 30	T 30	9	S	E
G 31	G 31	9	S	E
Q 32	Q 32	9	S	E
K 33	K 33	9	S	E
P 34	P 34	9	S	E
F 35	F 35	9	S	E
Q 36	Q 36	9	S	E
C 37	C 37	9	S	E
R 38	R 38	9	S	E
I 39	I 39	9	S	E
C 40	C 40	9	S	E
M 41	M 41	9	S	E

R 42	R 42	9	S	E
N 43	N 43	9	S	E
F 44	F 44	9	S	AE
S 45	S 45	9	S	AE
R 46	Q 46	9	S	E
S 47	H 47	9	S	E
D 48	T 48	9	S	E
H 49	G 49	9	S	A
L 50	L 50	9	S	A
T 51	N 51	9	S	E
T 52	Q 52	9	S	E
H 53	H 53	9	S	E
I 54	I 54	9	S	E
R 55	R 55	9	S	E
T 56	T 56	9	S	E
H 57	H 57	9	S	E
T 58	T 58	9	S	E
G 59	G 59	9	S	E
E 60	E 60	9	S	E
K 61	K 61	9	S	E
P 62	P 62	9	S	E
F 63	F 63	9	S	E
A 64	A 64	9	S	E
C 65	C 65	9	S	E
D 66	D 66	9	S	E
I 67	I 67	9	S	E
C 68	C 68	9	S	E
G 69	G 69	9	S	E
R 70	R 70	9	S	E
K 71	K 71	9	S	E
F 72	F 72	9	S	AE
A 73	A 73	9	S	AE
R 74	T 74	9	S	AE
S 75	L 75	9	S	E
D 76	H 76	9	S	E
E 77	T 77	9	S	E
R 78	R 78	9	S	E
K 79	D 79	9	.	EI
R 80	R 80	9	.	EI
H 81	H 81	9	.	AEI
T 82	T 82	9	.	AEI
K 83	K 83	9	.	EI
I 84	I 84	9	.	EI
H 85	H 85	9	.	EI
L 86	L 86	9	.	EI
R 87	R 87	9	.	EI
Q 88	Q 88	9	.	EI
K 89	K 89	9	.	EI

Target ME-AA : Template specified/identified is 1AAY Chain A

Columns 1-2 correspond to the target.

Columns 3-4 correspond to the template.

Column 5 is the confidence (0-9) for that pair of aligned residues according to sequence alignments

Column 6 is the selection status.

An 'S' indicates that the corresponding pair of residues has been confidently aligned and may be used in homology modeling.

A '.' indicates otherwise.

Please disregard column 7. This is only for diagnostic purposes and will soon be removed.

----- ME-AA.consensus -----

K 3	R 1	9	.	F
P 4	P 2	9	.	F
Y 5	Y 3	9	.	F
K 6	A 4	9	.	F
C 7	C 5	9	.	F
P 8	P 6	9	.	F
E 9	V 7	7	.	F
C 10	C 10	9	.	F
G 11	D 11	9	.	F
K 12	R 12	9	S	D
S 13	R 13	9	S	D
F 14	F 14	9	S	D
S 15	S 15	9	S	D
Q 16	R 16	9	S	D
S 17	S 17	9	S	D
S 18	D 18	9	S	D
N 19	E 19	9	S	A
L 20	L 20	9	S	A
Q 21	T 21	9	S	E
K 22	R 22	9	S	E
H 23	H 23	9	S	AE
Q 24	I 24	9	S	AE
R 25	R 25	9	S	AE
T 26	I 26	9	S	E
H 27	H 27	9	S	E
T 28	T 28	9	S	E
G 29	G 29	9	S	E
E 30	Q 30	9	S	E
K 31	K 31	9	S	E
P 32	P 32	9	S	E
Y 33	F 33	9	S	E
K 34	Q 34	9	S	E
C 35	C 35	9	S	E
P 36	R 36	9	S	E
E 37	I 37	9	S	E
C 38	C 38	9	S	E
G 39	M 39	9	.	
K 40	R 40	9	S	D
S 41	N 41	9	S	D
F 42	F 42	9	S	A

S 43	S 43	9	S	A
Q 44	R 44	9	S	E
S 45	S 45	9	S	E
S 46	D 46	9	S	E
D 47	H 47	9	S	E
L 48	L 48	9	S	E
Q 49	T 49	9	S	E
K 50	T 50	9	S	AE
H 51	H 51	9	S	AE
Q 52	I 52	9	S	E
R 53	R 53	9	S	E
T 54	T 54	9	S	E
H 55	H 55	9	S	E
T 56	T 56	9	S	E
G 57	G 57	9	S	E
E 58	E 58	9	S	E
K 59	K 59	9	S	E
P 60	P 60	9	S	E
Y 61	F 61	9	S	E
K 62	A 62	9	S	E
C 63	C 63	9	S	E
P 64	D 64	9	S	E
E 65	I 65	9	S	E
C 66	C 66	9	S	E
G 67	G 67	9	S	E
K 68	R 68	9	S	E
S 69	K 69	9	S	E
F 70	F 70	9	S	AE
S 71	A 71	9	S	AE
R 72	R 72	9	S	E
S 73	S 73	9	S	E
D 74	D 74	9	S	E
H 75	E 75	9	S	E
L 76	R 76	9	.	EI
S 77	K 77	9	.	EI
R 78	R 78	9	.	EI
H 79	H 79	9	.	EI
Q 80	T 80	9	.	EI
R 81	K 81	9	.	EI
T 82	I 82	9	.	EI
H 83	H 83	9	.	EI
Q 84	L 84	9	.	EI
N 85	R 85	9	.	EI

## Appendix C – Domain Information for Each Homology Modeled EGR Protein

Areas of overlap are underlined.

### AA-GD

RPYACPVESCDRRFSRSDDELTRHIRIHTGQKPFQCRICMRNFSRSDHLTTHIRTHTGEKPFA  
CDICGRKFARSDERKRHTKIHRLR

#### AA-GD Domain 1 (Residues 1 – 32)

RPYACPVESCDRRFSRSDDELTRHIRIHTGQKP

#### AA-GD Domain 2 (Residues 30 – 60)

QKPFQCRICMRNFSRSDHLTTHIRTHTGEKP

#### AA-GD Domain 3 (Residues 57 – 85)

EKPFACDICGRKFARSDERKRHTKIHRLR

### GD-AA

MERPYPVESCRRFSQKTNLDTHIRIHTGQKPFQCRICMRNFSQHTGLNQHIRTHTGE  
KPFACDICGRKFATLHTRDRHTKIHRLRQK

#### GD-AA Domain 1 (Residues 1 – 34)

MERPYPVESCRRFSQKTNLDTHIRIHTGQKP

#### GD-AA Domain 2 (Residues 32 – 61)

QKPFQCRICMRNFSQHTGLNQHIRTHTGEK

#### GD-AA Domain 3 (Residues 60 – 89)

EKPFACDICGRKFATLHTRDRHTKIHRLRQK

### ME-AA

KPYKCPECGKSFSQSSNLQKHQRTHTGEKPYKCPECGKSFSQSSDLQKHQRTHTGEKPY  
KCPECGKSFSRSDHLSRHQRTHQ

#### GD-AA Domain 1 (Residues 1 – 30)

KPYKCPECGKSFSQSSNLQKHQRTHTGEKP

#### ME-AA Domain 2 (Residues 28 – 58)

EKPYKCPECGKSFSQSSDLQKHQRTHTGEKP

#### ME-AA Domain 3 (Residues 56 – 82)

EKPYKCPECGKSFSRSDHLSRHQRTHQ

## Appendix D – RMS versus Crystal Structure for Clustered Models

### AA - GD

Residue Number	Residue Type	Time (ps)	RMS	Residue Number	Residue Type	Time (ps)	RMS
1	R	870	0.77	44	R	2970	2.3
2	P	485	0.82	45	S	1560	2.05
3	Y	2752.5	1.45	46	D	1557.5	1.27
4	A	3470	0.63	47	H	1087.5	0.82
5	C	1095	0.94	48	L	1037.5	0.47
6	P	1767.5	1.67	49	T	385	0.88
7	V	1555	1.01	50	T	837.5	0.67
8	E	1625	2.61	51	H	3847.5	1.52
9	S	1542.5	2.95	52	I	2415	0.57
10	C	280	2.05	53	R	1245	0.66
11	D	3592.5	2.48	54	T	1817.5	1.04
12	R	3400	3.04	55	H	2265	1.7
13	R	1600	0.66	56	T	1237.5	0.39
14	F	3425	1.72	57	G	317.5	0.29
15	S	1517.5	0.95	58	E	2300	1.74
16	R	3102.5	1.13	59	K	3112.5	0.65
17	S	2672.5	0.64	60	P	3182.5	0.41
18	D	1442.5	2.06	61	F	795	1.46
19	E	237.5	0.97	62	A	317.5	0.14
20	L	1115	0.57	63	C	2592.5	0.87
21	T	1090	1.18	64	D	1055	2.44
22	R	910	2.43	65	I	2675	1.09
23	H	882.5	1.77	66	C	1410	0.84
24	I	3047.5	1.85	67	G	3040	0.25
25	R	2610	1.76	68	R	830	3.68
26	I	1897.5	2.13	69	K	1265	1.23
27	H	2850	1.27	70	F	3562.5	1.5
28	T	837.5	335	71	A	2205	0.31
29	G	3010	0.87	72	R	3402.5	1.39
30	Q	3132.5	2.65	73	S	3367.5	0.33
31	K	3177.5	3.53	74	D	2490	1.91
32	P	2622.5	0.58	75	E	2592.5	1.07
33	F	3887.5	1.61	76	R	3342.5	2.14
34	Q	607.5	0.77	77	K	860	2.07
35	C	3927.5	0.9	78	R	1415	1.55
36	R	2102.5	4.02	79	H	1117.5	1.25
37	I	2365	0.94	80	T	1515	1.0
38	C	2437.5	0.49	81	K	3647.5	2.27
39	M	2717.5	0.51	82	I	3400	0.51
40	R	547.5	4.9	83	H	655	2.84
41	N	2235	0.49	84	L	272.5	1.17
42	F	200	0.46	85	R	1062.5	2.55
43	S	1745	0.78				

**GD-AA**

Residue Number	Residue Type	Time (ps)	RMS	Residue Number	Residue Type	Time (ps)	RMS
1	M	2410	1.3	45	S	525	0.84
2	E	460	2.74	46	Q	1120	0.97
3	R	3970	2.27	47	H	2235	0.87
4	P	2622.5	0.61	48	T	1080	1.16
5	Y	3580	1.49	49	G	2622.5	0.25
6	A	317.5	0.25	50	L	3827.5	0.6
7	C	2432.5	0.45	51	N	305	1.02
8	P	3597.5	1439.5	52	Q	1352.5	0.84
9	V	2622.5	0.41	53	H	1320	1.23
10	E	3105	2.73	54	I	1082.5	1.59
11	S	3760	0.98	55	R	3207.5	2.21
12	C	325	1.26	56	T	3087.5	2.07
13	D	890	1.47	57	H	3302.5	1.18
14	R	2645	1.92	58	T	1805	0.32
15	R	1107.5	1.63	59	G	1465	0.16
16	F	525	1.47	60	E	1200	1.52
17	S	1030	0.68	61	K	640	0.61
18	Q	377.5	0.86	62	P	3467.5	0.48
19	K	1097.5	1.34	63	F	2010	1.57
20	T	317.5	1.22	64	A	2622.5	0.49
21	N	3382.5	1.14	65	C	485	0.95
22	L	940	0.28	66	D	2325	0.92
23	D	3225	1.91	67	I	992.5	1.12
24	T	1680	0.75	68	C	2010	0.49
25	H	2105	0.64	69	G	2683.5	0.54
26	I	3972.5	1.54	70	R	2620	1.04
27	R	1837.5	1.13	71	K	1622.5	1.71
28	I	1350	0.68	72	F	1222.5	1.56
29	H	2285	1.24	73	A	2622.5	0.15
30	T	3507.5	0.29	74	T	2927.5	0.39
31	G	3207.5	0.15	75	L	3757.5	0.7
32	Q	1377.5	0.83	76	H	3087.5	1.36
33	K	2982.5	0.32	77	T	2820	0.71
34	P	317.5	0.32	78	R	1817.5	1.61
35	F	3660	1.56	79	D	862.5	1.96
36	Q	2680	0.38	80	R	1287.5	3.62
37	C	1820	1.08	81	H	1072.5	0.99
38	R	3945	2.19	82	T	3252.5	0.91
39	I	2137.5	0.87	83	K	2417.5	1.72
40	C	815	0.91	84	I	1045	0.65
41	M	2342.5	0.79	85	H	1665	2.8
42	R	2140	5.41	86	L	2012.5	1.07
43	N	2625	1.12	87	R	2017.5	3.46
44	F	3955	0.7	88	Q	2225	0.93

## ME – AA

Residue Number	Residue Type	Time (ps)	RMS	Residue Number	Residue Type	Time (ps)	RMS
1	K	3617.5	1.78	42	Q	1320	0.71
2	P	3702.5	0.55	43	S	1707.5	0.33
3	Y	1537.5	1.71	44	S	417.5	0.51
4	K	2650	0.94	45	D	3735	1.07
5	C	2385	0.5	46	L	2090	0.5
6	P	392.5	0.57	47	Q	1307.5	1.97
7	E	1435	2.12	48	K	957.5	2.4
8	C	3412.5	1.25	49	H	1412.5	0.79
9	G	532.5	0.67	50	Q	1592.5	0.76
10	K	1640	0.59	51	R	1710	2.28
11	S	992.5	0.75	52	T	2560	0.46
12	F	1135	1.45	53	H	1665	2.22
13	S	322.5	0.82	54	T	2622.5	0.56
14	Q	735	1.18	55	G	860	0.39
15	S	600	0.48	56	E	1752.5	2.14
16	S	1665	0.66	57	K	2640	2.56
17	N	2420	1.11	58	P	1732.5	0.35
18	L	3112.5	0.46	59	Y	3780	1.49
19	Q	2910	0.76	60	K	1090	1.86
20	K	3075	2.21	61	C	1707.5	0.6
21	H	2597.5	0.66	62	P	2187.5	0.61
22	Q	2297.5	1.04	63	E	3995	2.49
23	R	3812.5	2.01	64	C	2287.5	0.98
24	T	3247.5	0.46	65	G	3007.5	0.25
25	H	1170	2.3	66	K	817.5	0.71
26	T	1995	0.37	67	S	620	0.47
27	G	482.5	0.21	68	F	3277.5	1.46
28	E	3640	1.98	69	S	2882.5	0.45
29	K	2115	1.1	70	R	2680	1.98
30	P	1400	0.47	71	S	1120	0.44
31	Y	1105	0.66	72	D	2455	0.55
32	K	402.5	3.06	73	H	2262.5	2.08
33	C	3605	0.64	74	L	3957.5	0.51
34	P	1960	0.42	75	S	2935	0.84
35	E	2465	2.86	76	R	2005	4.51
36	C	1627.5	0.62	77	H	1265	0.93
37	G	2260	0.22	78	Q	2327.5	0.69
38	K	2097.5	1.81	79	R	2650	4.92
39	S	2397.5	0.9	80	T	1457.5	1.52
40	F	2985	0.95	81	H	660	1.88
41	S	922.5	0.79				