

**UCSF**

**UC San Francisco Electronic Theses and Dissertations**

**Title**

3D Motifs as Signatures of Protein Function and Evolution

**Permalink**

<https://escholarship.org/uc/item/1061t12w>

**Author**

Polacco, Benjamin John

**Publication Date**

2007-07-24

Peer reviewed|Thesis/dissertation

3D Motifs as Signatures of Protein Function and Evolution

by

Benjamin John Polacco

DISSERTATION

Submitted in partial satisfaction of the requirements for the degree of

DOCTOR OF PHILOSOPHY

in

Biological and Medical Informatics

in the

GRADUATE DIVISION

of the

UNIVERSITY OF CALIFORNIA, SAN FRANCISCO

Copyright (2007)

by

Benjamin Polacco

## Acknowledgments

I would like to thank my research advisor Patricia Babbitt for her support and insight; both were essential for the completion of this work. The members of the Babbitt lab were always ready to help when asked, and offered much feedback through formal group meetings as well as informal discussions.

Thanks are also due to my orals and thesis committee for their guidance and criticisms of my ideas when I presented them. Several good ideas contained here can be attributed to the insights of my committee, while all the questionable ideas are strictly my own.

While I often think this could have been completed more quickly without the responsibilities of being a father to my family, these years would have been much less fulfilling without them. I especially thank my wife for her support of me emotionally, and the family financially throughout this work.

This dissertation is divided into an introduction, four chapters, and a conclusion. Three chapters are based on work that is either published already or expected to be published soon. Chapter 1 is entirely my work, though some of its findings were published together with the work of Elaine Meng in the journal *Proteins*. The text of Chapter 2 is a reprint of the material as it appears in the journal *Bioinformatics*. The coauthor listed in that publication directed and supervised the research that forms the basis for this chapter. Chapter 4, is a manuscript that will be submitted for publication in a peer-reviewed journal.

Meng, E. C., B. J. Polacco and P.C. Babbitt. (2004). "Superfamily active site templates." *Proteins* 55(4): 962-76.

Polacco, B. J. and P. C. Babbitt (2006). "Automated discovery of 3D motifs for protein function annotation." *Bioinformatics* 22(6): 723-30.

## **Abstract: 3D Motifs as Signatures of Protein Function and Evolution**

**Benjamin Polacco**

The ability to predict a protein's function from its structure is becoming more important with the increasing pace at which international structural genomics projects make structures available for proteins with no known function. The function of a protein is frequently determined by relatively small regions in an overall structure. This dissertation investigates signature 3D motifs, or small subsets of a protein's residues, that capture the critical structural determinants of function shared by an entire group of proteins. First, with an investigation of randomly selected 3D motifs I show that motifs built from important functional residues are better at identifying proteins to a superfamily with a common functional mechanism than any other motifs. Next I develop a genetic algorithm, named GASPS, that chooses a motif based on its ability to identify a group of proteins. I demonstrate its effectiveness on four divergent superfamilies, and a convergent group of serine proteases. Again, I demonstrate that the best motifs, as chosen by GASPS this time, contain known functional residues. Chapter 3 investigates the use of a geometrical statistical model to predict the number of expected random matches to a motif. This simple geometrical model performs very well overall, but it under-predicts matches to motifs that are the result of general physical and chemical characteristics of proteins, such as disulfide bridges and hydrophobic clusters. This model is rejected for its use in GASPS in favor of the original empirical method. Finally, I report a broad survey of signature 3D motifs, generated by applying GASPS to all available functionally similar and homologous groups of proteins. Motifs are mostly restricted to homologous groups, with a higher chance of a better motif in homologous and isofunctional groups. I

report on general trends in structural conservation and find that catalytic, ligand binding, disulfide, and stabilized charged residues are over-represented among conserved motifs. Additionally, I find that glycines appear to be the most frequently conserved residue, especially important in ligand binding sites. This collection of motifs is useful for identification of function in unknown proteins, as well as describing trends in protein evolution.

# Table of Contents

<b>Introduction</b> .....	<b>1</b>
Chapter 1 Summary.....	3
Chapter 2 Summary.....	4
Chapter 3 Summary.....	4
Chapter 4 Summary.....	5
Conclusions.....	7
References.....	8
<b>Chapter 1: Random Motifs and Superfamily Active Site Templates</b> .....	<b>10</b>
Introduction.....	10
Materials and Methods.....	11
Motif Searches.....	11
Structure Libraries .....	12
Calculating Conservation .....	12
Results.....	13
Scores of randomly generated motifs .....	13
Substitutions in random motifs.....	16
Limiting Residue Types in Random Motifs.....	18
Important Residues in Motifs .....	19
Discussion.....	22
References.....	24
<b>Introduction to Chapter 2</b> .....	<b>26</b>
<b>Chapter 2: Automated Discovery of 3D Motifs for Protein Function Annotation</b> ....	<b>27</b>
Abstract.....	27
Introduction.....	28
Methods.....	30
Motif Representation and Matching .....	30
GASPS .....	31
Structure Library.....	33
Positive and Negative Sets .....	33
Cross-Validation .....	34
PSI-BLAST and CE Libraries .....	35
Results.....	36
Validation of GASPS.....	36
Detection of Key Functional Residues.....	41
Discussion .....	47
Using GASPS for Function Identification .....	47
Location of Functional Information.....	48
Inference of Function for Diverse Groups .....	50
Future Applications.....	50
Acknowledgements .....	51
References.....	51

<b>Supplementary Materials for Chapter 2.....</b>	<b>54</b>
Significance of Optimized GASPS Scores .....	54
Sources of Variability .....	54
Detection of New Unidentified Structures.....	57
Allowing Substitutions in Motifs .....	57
References .....	58
<b>Chapter 3: An analysis of computed expectations for random matches to 3D motifs</b>	<b>63</b>
.....	63
Introduction .....	63
Results.....	64
Modifying the GASPS scoring function. ....	64
Expectation values compared. ....	67
GASPS scores compared.....	70
GASPS with $G_c$ on Random Groups.....	71
Composition of motifs .....	72
Accuracy of motifs at identifying homologous groups.....	74
Discussion .....	78
References .....	78
<b>Introduction to Chapter 4 .....</b>	<b>80</b>
<b>Chapter 4: An exhaustive survey of 3D motifs.....</b>	<b>81</b>
Abstract .....	81
Introduction .....	82
Methods.....	86
GASPS .....	86
Protein groups.....	86
Searching motif libraries with proteins.....	87
Results.....	89
Quality of Motifs .....	89
Patterns of conservation in 3D.....	94
Residue types in motifs.....	98
Annotation of protein structures .....	101
Homology models.....	103
Discussion .....	105
References .....	108
<b>Conclusion.....</b>	<b>111</b>
<b>Appendix 1: GASPS Package.....</b>	<b>114</b>
ReadMe .....	114
GASPS.py .....	116
polacco/BlastXML.py.....	143
polacco/Data.py .....	147
polacco/MultiAlign.py .....	149
polacco/Spasm.py .....	165
polacco/XML.py.....	178
polacco/utils.py.....	181



test/astreal_1.65_SF.lib (partial).....	183
test/d2mnr_1.fasta .....	184
test/d2mnr_1.fasta.psiblast.xml.faln (partial).....	184
test/d2mnr_1.pdb (partial).....	185
test/enolase.lib (partial).....	185
test/enolase.list .....	186
<b>Appendix 2: GASPSdb CGI scripts .....</b>	<b>187</b>
GASPSdb .....	187
jsonMotif.....	209
<b>Appendix 3: GASPSdb Web Interface.....</b>	<b>213</b>
GASPSdb Home Page .....	213
GASPSdb Search Page .....	214
GASPSdb Browsing Page.....	215
GASPSdb Group Description Page .....	216
GASPSdb Search Results Page .....	217
GASPSdb Help Page .....	218
About GASPSdb.....	218
GASPSdb References Page.....	220
References .....	220

# List of Figures

## Chapter 1

Figure 1. Cumulative histograms of scores of randomly generated motifs. ....	14
Figure 2. Cumulative histograms of scores of conserved and close random motifs. ....	16
Figure 3. Cumulative histograms of scores of random motifs with allowed substitutions. .....	18
Figure 4. Cumulative histograms of scores of random motifs with only polar residues. ...	19
Figure 5. Residues that contribute to motif scores. ....	21
Figure 6. Scores of partial motifs based on the functional site. ....	21

## Chapter 2

Figure 1. Generality of GASPS motifs based on sensitivity from two experiments: cross-validation and detection of newer structures. ....	37
Figure 2. Sensitivity of GASPS motifs compared with other techniques. ....	39
Figure 3. Scores and functional significance of GASPS motifs. ....	43

## Chapter 2 Supplementary Materials

Figure i. Distributions of GASPS scores on artificial and real groups. ....	60
Figure ii. Stochasticity of GASPS results. ....	61
Figure iii. GASPS motifs for 2hlc, a trypsin-like serine protease. ....	62

## Chapter 3

Figure 1. Relation between "ROC Credit", P Values, and Expected False Positives. ....	67
Figure 2. Empirical counts of false positives versus computed expectation values. ....	70
Figure 3. GASPS scores (G) compared between empirical and computed methods. ....	71
Figure 4. Distributions of motifs by GASPS with $G_c$ on random groups. ....	72
Figure 5. Composition of motifs generated by GASPS with computed G scores. ....	73
Figure 6. SCOP superfamilies identified by motifs generated by empirical G scores compared to computed G scores. ....	76
Figure 7. SCOP families identified by motifs generated by empirical G scores compared to computed G scores. ....	77

## Chapter 4

Figure 1. Distribution of motif G-scores on SCOP groups. ....	91
Figure 2. Distribution of motif G-scores on Gene Ontology and SCOP groups. ....	91
Figure 3. Number of distinct EC classes at first position in each SCOP group. ....	94

Figure 4. Number of distinct EC classes at first two positions in each SCOP group.....	94
Figure 5. Residue interactions captured by motifs. ....	97
Figure 6. Dominance of residue types, compared against background residue frequency, and at different G-scores. ....	101
Figure 7. Coverage of GASPSdb compared to other 3D motif libraries and PSI-BLAST. .....	103

### **Appendix 3**

Figure 1. Home page of GASPSdb.....	213
Figure 2. GASPSdb Search Page. ....	214
Figure 3. GASPSdb Browse Page. ....	215
Figure 4. GASPSdb Group description page, partial.....	216
Figure 5. Search results table for search of 1rvk against SCOP superfamily motifs. ....	217

## List of Tables

### Chapter 2

Table 1. Functionally Similar Protein Groups .....	35
--	----

### Chapter 2 Supplementary Materials

Table i. Improvements in GASPS by using substitutions on Crotonase and HAD superfamilies.....	59
---	----

### Chapter 3

Table 1. Overlap of significant motifs with catalytic sites in CSA.....	74
---	----

### Chapter 4

Table 1. Group and motif counts by classification.....	88
--	----

# Introduction

As proteins are the major gene products that act in living cells, understanding the functions of proteins is a critical step in translating genomic sequences into useful biological knowledge relevant to the health sciences. With today's efforts in structural genomics that aim to provide for each protein a model of its shape or structure in a cell (Blundell et al. 2000), knowledge of a protein's structure is becoming a more common starting point for determining a protein's function (Teichmann et al. 2001; Watson et al. 2007). As different functions can be performed by proteins that have very similar overall structures and folds (Chothia 1992; Todd et al. 1999), it is clear that we have to look at fine-scale details or local protein structure to accurately describe a protein's function. Over evolutionary time, identical proteins can diverge to have very different sequences by the accumulation of random neutral changes that do not change function (neutral drift), but these proteins will still share whatever structural components have been critical to their function. Additionally, as proteins evolve to perform new functions they can make use of existing local structural features that contribute the same partial function to both the new and old functions (Gerlt et al. 2001; Bartlett et al. 2003). This explains, for the most part, why all members of a diverse group of proteins often make use of the same configuration of a small number of amino acids that can be directly related to function. We can use these clusters of amino acids, called three dimensional (3D) motifs, as signatures of function. This work investigates these signature 3D motifs to show how the identification and understanding of protein function can be advanced through these repeated structural elements.

Because 3D motifs are closely tied to the evolution of function, a study of 3D motifs also describes the manner in which protein function evolves. The evolution of new function proceeds through one of two paths tied to the existence of 3D motifs. First, as descendants of a single protein diverge in function, existing functional components can be entirely replaced by new functional components so that no 3D motif will persist between modern day proteins. On the other hand, as new functions evolve, proteins can make use of existing functional components to perform one or more components in the overall function. If across these different functions, the same functional component is reused, a 3D motif will persist in modern day proteins. 3D motifs can also be present in convergent proteins, those that perform the same function but have no common ancestor (Dodson et al. 1998). If we observe frequent cases of convergent motifs this is evidence that the possible ways any proteins can evolve to perform a single function are limited.

Much work has been done by others examining 3D motifs. Studies have shown their effectiveness on a handful of cases (Wallace et al. 1996; Fetrow et al. 1998; Russell 1998), tools have been developed that can search protein structures for matches to motifs (Artymiuk et al. 1994; Kleywegt 1999; Barker et al. 2003), and motifs are being collected from literature descriptions of enzyme active sites (Torrance et al. 2005). Still, no study has yet systematically asked on how many and on what types of various protein groups can we use signature 3D motifs. This dissertation extends our knowledge of 3D motifs by inventing a novel method for discovering signature 3D motifs and applying this method to a large set of protein groups. This generates a set of motifs that are not only useful for protein annotation, but because they were systematically generated, provide an even and unbiased picture of the distribution of 3D motifs and patterns within them. Specifically,

we see that homology is the most important generator of signature 3D motifs, but functional diversity also plays a role. Though groups with diverse functions and signature 3D motifs are not uncommon, homologous groups with many varied functions are less likely to have a signature motif.

## ***Chapter 1 Summary***

The work I describe in Chapter 1 lays the foundation for my method, and demonstrates how significant findings and research paths are often stumbled upon by accident. I worked together with Elaine Meng, who was investigating signature 3D motifs in the active sites of enzyme superfamilies (Meng et al. 2004). While Elaine assembled motifs from residues known to be functionally important for the superfamilies, I performed the control study to show that motifs based on this expert knowledge identified the superfamily better than motifs assembled from randomly chosen residues. To make a more compelling comparison, I tested constraining the residues in the random motifs by distance from each other, then conservation, and then residue type. I added an automated system for allowing position specific substitutions based on a multiple sequence alignment. While these increased the quality of the randomly generated motifs, the published result from this work remained that the motifs built from functional knowledge always outperformed the automatically generated random motifs (Meng et al. 2004). This same result viewed from a slightly different angle would provide the inspiration that led to this entire dissertation: with a few simple constraints, a random guess could produce motifs that begin to approach the quality of expert derived motifs.

## ***Chapter 2 Summary***

This earliest work not only provided the insight that would lead to the development of my method, but also provided most of the software development. My method was given the acronym GASPS for Genetic Algorithm Search for Patterns in Structures. A genetic algorithm develops solutions to problems by choosing from among a set of guesses the best ones, then making new guesses by adding to, deleting from, or recombining the best guesses made so far. Using most of the random motif generation system I presented in Chapter 1 to create the first guesses, I added a method for measuring performance of motifs, a system to alter and recombine motifs, and an iterative process. This resulted in a version of GASPS that I described in a published manuscript (Polacco et al. 2006), included here as Chapter 2. As an alternative to building motifs from often-limited expert knowledge, GASPS identifies patterns of 3 to 10 residues that maximize function prediction. The unbiased approach of GASPS allowed us to test the assumption that residues that provide function are the most informative for predicting function. I applied GASPS to superfamilies with varied functions as well as the serine proteases, an example of convergent evolution of active sites (Dodson et al. 1998). The motifs found by GASPS are as good at function prediction as 3D motifs based on expert knowledge. The GASPS motifs with the greatest ability to predict protein function consist mainly of known functional residues.

## ***Chapter 3 Summary***

In an effort to improve GASPS, I investigated the theoretical statistics of random matches to 3D motifs, or false positives. GASPS seeks to find a motif for a group where all group



members match within a deviation threshold stringent enough to make random matches to unrelated proteins rare. The original GASPS determines this threshold empirically, by searching for matches to each candidate motif among all non-group proteins that it should not match. The work described in Chapter 3, answers whether this empirical distribution of matches is necessary or instead is a theoretical statistical model of matches to 3D motifs sufficient. Computing the empirical distribution takes more time than any other GASPS step, so if it could be replaced it would significantly reduce the computing time necessary to generate motifs with GASPS. I show how the scoring function that GASPS uses to rank motifs can be modified to use a statistical model of motif matches developed by Stark et al. (2003). Qualitatively, motifs generated by this faster GASPS are very similar to the original GASPS, with similar rates of overlap with functionally significant residues. However, these motifs fail to identify new structures to the appropriate group with the same accuracy. This decreased accuracy is due more to false positives than false negatives, indicating the motifs are not as unique as the model would predict. This results from the use of a solely geometrical model that cannot account for common physically favorable interactions frequently observed across various protein groups, such as salt bridges or disulfides. This makes the faster GASPS a useful tool for discovering a motif that is well conserved by a group, but not for generating motifs useful for annotation of new structures. This faster GASPS was not used for any other work described here.

### ***Chapter 4 Summary***

The pieces are now in place to apply GASPS across as much of the protein universe as possible in order to generate as many signature 3D motifs as possible. Doing so allows for an examination of the evolution of fine scale protein structure by determining how

widespread are conserved 3D motifs, and what structural features tend to be conserved. I apply GASPS to homologous superfamilies and families in the Structural Classification of Proteins (SCOP) (Murzin et al. 1995), as well as isofunctional groups defined by the Gene Ontology (GO) (Ashburner et al. 2000). I find that non-homologous but isofunctional groups do not commonly share a motif. This suggests that most protein functions, at least as they are commonly described, can be accomplished by very different means in unrelated proteins. Homologous groups more often share a conserved motif, with about one third of all SCOP groups showing a strongly conserved motif. Many of these superfamilies with strong motifs have very diverse functions, revealing where evolution has reused functional components to produce different overall reactions. The remaining two thirds of groups with less-conserved motifs reveal that evolution of new functions in homologous groups is not usually constrained to maintain the positions of a critical set of residues.

These motifs also allow us to examine what features are among the most conserved. Again, we see a strong relationship between motifs and function. The motifs frequently overlap with known catalytic, metal and other ligand binding sites. Additionally, disulfides as well as stabilized charged residue pairs are frequent components of the most conserved motifs. Residue distribution among the motifs is mostly as expected based on these common features: cysteine, histidine, aspartate and glutamate are among the most frequent. More surprisingly, glycine, leucine and proline are ranked first, fourth and seventh, respectively, among the most frequent motif residues. The dominant role of leucine can be attributed mostly to its high frequency among the entire proteins. Glycine is well conserved where its unique backbone angles and space allowances (Jornvall et al.

1984; Dym et al. 2001) are critical for function. The unique geometry afforded glycine seems especially important at binding sites: glycines in motifs show the greatest rate of non-metal ligand interaction among all residue types.

To maximize the impact of this work, I have made available the motifs generated in this broad study via a web resource named GASPSdb (<http://gaspsdb.rbvi.ucsf.edu>). The motifs at this site can be searched, browsed or downloaded. One search capability enables users to search for matches to the GASPS motifs among a protein structure they can provide or choose from the Protein Data Bank (PDB) (Berman et al. 2000). Because each motif is generated to be a signature motif for a functional or homologous group, a matching motif indicates that the new structure is a likely member of the group, and the matched residues are likely to be important for the protein's function. I show that the GASPSdb resource provides a greater coverage than other available 3D motif resources (Stark et al. 2003; Torrance et al. 2005). It also proves effective on low quality structural models computed from homology. This effectiveness on homology models is very important for the description of function in the homology models that structural genomics aims to make possible (Blundell et al. 2000).

## ***Conclusions***

This study has grown from its beginnings where its goal was to merely show the ineffectiveness of randomly chosen motifs, to show how when combined with an effective selection and recombination process those same random motifs can become useful signatures of protein function and evolution. While I generated a large number of signature motifs that will enable us to more accurately annotate structures, I find that not

all groups can be identified by a signature motif. The distributions of these signature motifs represent just a single but useful view on the evolution of protein structure and function at a fine scale. We observe that evolution has used both schemes I presented regarding function and local structure. I identify both the homologous groups that have re-used functional features for multiple different overall functions, as well as groups which keep no single functional feature as they evolve to perform new overall functions. There are many known (and probably unknown) protein groups with insufficient structures for GASPS to work on effectively. As new structures are solved, an automated process like GASPS is well suited to continue to analyze new groups and new structures.

## **References**

- Artymiuk, P. J., A. R. Poirrette, et al. (1994). "A graph-theoretic approach to the identification of three-dimensional patterns of amino acid side-chains in protein structures." J Mol Biol **243**(2): 327-44.
- Ashburner, M., C. A. Ball, et al. (2000). "Gene ontology: tool for the unification of biology. The Gene Ontology Consortium." Nat Genet **25**(1): 25-9.
- Barker, J. A. and J. M. Thornton (2003). "An algorithm for constraint-based structural template matching: application to 3D templates with statistical analysis." Bioinformatics **19**(13): 1644-9.
- Bartlett, G. J., N. Borkakoti, et al. (2003). "Catalysing new reactions during evolution: economy of residues and mechanism." J Mol Biol **331**(4): 829-60.
- Berman, H. M., J. Westbrook, et al. (2000). "The Protein Data Bank." Nucleic Acids Res **28**(1): 235-42.
- Blundell, T. L. and K. Mizuguchi (2000). "Structural genomics: an overview." Prog Biophys Mol Biol **73**(5): 289-95.
- Chothia, C. (1992). "Proteins. One thousand families for the molecular biologist." Nature **357**(6379): 543-4.
- Dodson, G. and A. Wlodawer (1998). "Catalytic triads and their relatives." Trends Biochem Sci **23**(9): 347-52.
- Dym, O. and D. Eisenberg (2001). "Sequence-structure analysis of FAD-containing proteins." Protein Sci **10**(9): 1712-28.

- Fetrow, J. S. and J. Skolnick (1998). "Method for prediction of protein function from sequence using the sequence-to-structure-to-function paradigm with application to glutaredoxins/thioredoxins and T1 ribonucleases." J Mol Biol **281**(5): 949-68.
- Gerlt, J. A. and P. C. Babbitt (2001). "Divergent evolution of enzymatic function: mechanistically diverse superfamilies and functionally distinct suprafamilies." Annu Rev Biochem **70**: 209-46.
- Jornvall, H., H. von Bahr-Lindstrom, et al. (1984). "Extensive variations and basic features in the alcohol dehydrogenase-sorbitol dehydrogenase family." Eur J Biochem **140**(1): 17-23.
- Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." J Mol Biol **285**(4): 1887-97.
- Meng, E. C., B. J. Polacco, et al. (2004). "Superfamily active site templates." Proteins **55**(4): 962-76.
- Murzin, A. G., S. E. Brenner, et al. (1995). "SCOP: a structural classification of proteins database for the investigation of sequences and structures." J Mol Biol **247**(4): 536-40.
- Polacco, B. J. and P. C. Babbitt (2006). "Automated discovery of 3D motifs for protein function annotation." Bioinformatics **22**(6): 723-30.
- Russell, R. B. (1998). "Detection of protein three-dimensional side-chain patterns: new examples of convergent evolution." J Mol Biol **279**(5): 1211-27.
- Stark, A. and R. B. Russell (2003). "Annotation in three dimensions. PINTS: Patterns in Non-homologous Tertiary Structures." Nucleic Acids Res **31**(13): 3341-4.
- Stark, A., S. Sunyaev, et al. (2003). "A model for statistical significance of local similarities in structure." J Mol Biol **326**(5): 1307-16.
- Teichmann, S. A., A. G. Murzin, et al. (2001). "Determination of protein function, evolution and interactions by structural genomics." Curr Opin Struct Biol **11**(3): 354-63.
- Todd, A. E., C. A. Orengo, et al. (1999). "Evolution of protein function, from a structural perspective." Curr Opin Chem Biol **3**(5): 548-56.
- Torrance, J. W., G. J. Bartlett, et al. (2005). "Using a Library of Structural Templates to Recognise Catalytic Sites and Explore their Evolution in Homologous Families." J Mol Biol **347**(3): 565-81.
- Wallace, A. C., R. A. Laskowski, et al. (1996). "Derivation of 3D coordinate templates for searching structural databases: application to Ser-His-Asp catalytic triads in the serine proteinases and lipases." Protein Sci **5**(6): 1001-13.
- Watson, J. D., S. Sanderson, et al. (2007). "Towards fully automated structure-based function prediction in structural genomics: a case study." J Mol Biol **367**(5): 1511-22.

# Chapter 1: Random Motifs and Superfamily Active Site Templates

## *Introduction*

When this work was started there were already multiple studies showing that specific active site 3D motifs could be used successfully to identify specific protein functions (Artymiuk et al. 1994; Wallace et al. 1997; Fetrow et al. 1998; Russell 1998; Kleywegt 1999). A major appeal of 3D motifs is that they provide a direct link between structural details and function in a way that sequence based or whole protein fold-based comparisons could not. In addition to a detailed view on structure, accurately describing the linkage between structure and function can benefit from a detailed view of protein function. Instead of treating an enzyme's function as a single unit, it can be broken down into smaller mechanistic steps, and superfamilies of enzymes can share one or more functional steps (Gerlt et al. 2001; Babbitt 2003). The work I present here was my part of a collaboration to show that superfamilies of enzymes, and therefore just the smaller element of function that they share can be identified by a single motif (Meng et al. 2004). Superfamily active site template was the name given to a 3D motif that is shared among members of a diverse superfamily that are responsible for the superfamily's shared function. While previous studies of 3D motifs have constructed motifs based on knowledge of functional residues, none looked specifically at the question of whether there were other informative residues—residues that uniquely identified the group of proteins. Investigating this question was my contribution to the study. While my collaborator did the traditional motif-building jobs of compiling lists of functional

residues and their similarities between related proteins, I constructed thousands of motifs at random to determine whether the functional residues were required or whether there were other residues that were conserved in three dimensions across a superfamily.

The work I present in this chapter is an important component of the published superfamily study. Additionally, it provides an analysis of random motifs that would guide the development of my technique, named GASPS, described in Chapter 2. With GASPS I make use of random guesses in a genetic algorithm, so the knowledge of which constraints can lead to better random guesses and the ways in which partial solutions score compared to an overall solution are important.

Portions of this work were published previously in the journal *Proteins* (Meng et al. 2004).

## ***Materials and Methods***

### **Motif Searches**

Active site template searching was performed with SPASM (Kleywegt 1999). A motif is supplied to SPASM as a file containing the atomic coordinates of the residues of interest. These coordinates are taken from the original Protein Data Bank (PDB) (Berman et al. 2000) file of each source structure. SPASM allows explicit specification of the residue types that can match each motif residue, referred to as substitutions later. The  $\alpha$ -carbon (CA) and computed side-chain centroid (SC) are used to describe each motif residue. The internal CA-CA and SC-SC distances of the motif and each candidate match are compared, and candidate matches are pruned if they exceed user-specified maximum deviations, in our case the maximum CA-CA distance deviation was set to 5.0 Å, and the

maximum SC-SC distance deviation was 3.8 Å. The remaining candidate matches are reoriented onto the motif and those fulfilling a user-specified RMSD cutoff are saved (3.2 Å). Thus, the input parameters include motif coordinates, allowed substitutions, a maximum CA-CA deviation cutoff, a maximum SC-SC deviation cutoff, a maximum RMSD cutoff, and what database to search. SPASM-searchable databases are derived directly from PDB files, but have been preprocessed down to the CA and computed SC coordinates for each residue. The preprocessing program, MKSPAZ, is available along with SPASM from the Uppsala Software Factory (<http://xray.bmc.uu.se/usf/index.html>).

## **Structure Libraries**

Motif sensitivity and specificity was evaluated by searching a sequence-unique subset of the PDB; this database, spasm100, can be downloaded from the Uppsala Software Factory (<http://xray.bmc.uu.se/usf/index.html>). The July 2002 version of spasm100 (8255 entries, including 22 true positive enolase superfamily members) was used.

## **Calculating Conservation**

Conservation of positions in a protein structure were calculated from a multiple sequence alignment generated by BLAST (Altschul et al. 1997) with default values against a non redundant protein sequence database, nrdb90 (Holm et al. 1998). Conservation was calculated from the multiple sequence alignment by a method that weights to reduce the effects of redundancy, considers conservative substitutions based on a substitution matrix, and penalizes gaps (Valdar 2002).



## **Results**

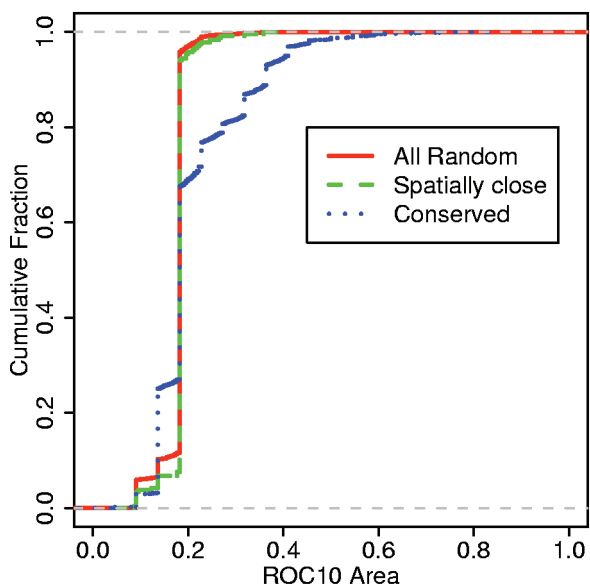
### **Scores of randomly generated motifs**

The superfamily motif derived from knowledge of functionally important residues from mandelate racemase (PDB id 2mnr) performed better than motifs from other available structures at identifying superfamily members with high sensitivity and specificity (Meng et al. 2004). Did the functional information identify the best residues for a motif, or could other motifs perform as well? To answer this, I generated motifs at random by selecting five residues entirely at random from the mandelate racemase structure. Each motif was scored by calculating the area under an ROC plot to 10 false positives based on the root mean squared deviation (RMSD) between the motif and its match in the structure (ROC10). All scores were normalized so that the maximum allowable ROC10 score was set to 1.0, the score that implies all superfamily structures match at a lower RMSD than any false positive. The vast majority of about 500 randomly generated motifs do not have an ROC10 area greater than 0.18 (Figure 1). This score corresponds to matching only the four superfamily structures that are most similar to 2mnr.

Most of the above randomly generated motifs appear very different from what commonly used 3D motifs look like. In an effort to make a more compelling comparison between randomly generated motifs, and those based on expert knowledge, I tested applying constraints on the generation of motifs. 3D motifs are typically composed of residues that are known to interact, so they must be close in space. The first constraint I applied to make the random motifs look more like typical motifs was therefore to restrict the residues to a 7.5 Å neighborhood, measured at their  $\alpha$ -carbon, of an initial chosen

residue. By itself, this restriction offered only a slight improvement to the ROC10 scores of the generated motifs.

A good 3D motif is maintained by evolution in all group structures, therefore its residues cannot be among the most variable in close relatives. The next constraint I tested was then to eliminate the most variable or least conserved residues observed in close homologs. Just eliminating residues with conservation below 0.6 showed a significant improvement in ROC10 areas, with no other constraints.

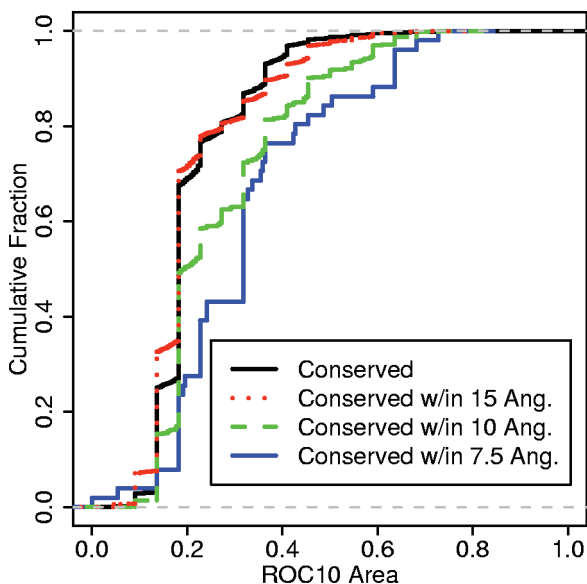


**Figure 1. Cumulative histograms of scores of randomly generated motifs.**

The red line, “All Random”, represents 5331 motifs of five residues chosen entirely at random from a single mandelate racemase structure. The green line, “Spatially Close”, represents 237 similar motifs with the only constraint that residues within a single motif are restricted to lie within a 7.5 Å neighborhood. The blue line, “Conserved”, represents 2825 random motifs with the only constraint that all residues must have a conservation score (see Methods) greater than 0.6.

While the spatial constraint by itself showed little effect, putting the conservation and spatial constraints together resulted in an even greater improvement in ROC10 areas (Figure 2). This effect is strongest when the residues are maintained within a 7.5 Å

neighborhood, compared with larger neighborhoods. By enforcing such a small neighborhood, I significantly decrease the number of possible random motifs because each residue only has a very limited set of residues it could build a motif with. This greatly minimizes the number of motifs that many conserved residues could be a part of because they do not cluster spatially with large numbers of other conserved residues. Others have actually used clusters of sequence-conserved residues on a protein structure to identify functionally important residues (Lichtarge et al. 1996). These constraints used here significantly enrich the available residues with functional residues, which can explain the increase in ROC10 areas for the smaller neighborhoods.

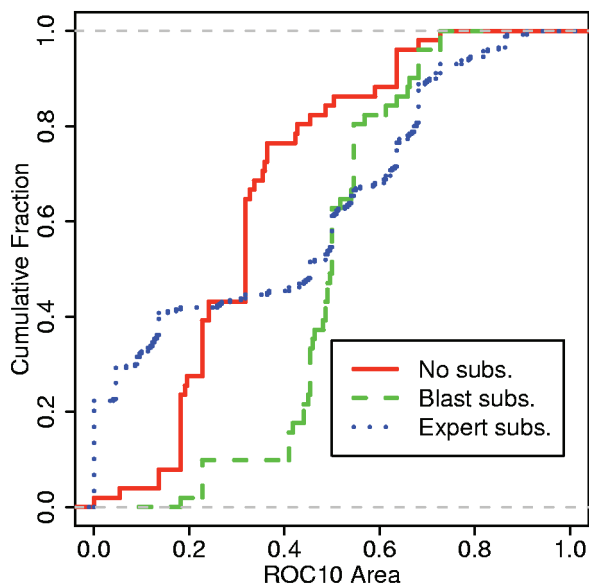


**Figure 2. Cumulative histograms of scores of conserved and close random motifs.** The black line, labeled “Conserved”, is identical to the same-labeled data shown in Figure 1. The remaining lines show the effects of adding an additional constraint that all residues must lie in a 7.5, 10, or 15 Å neighborhood. These lines represent 51, 496 and 3093 motifs, respectively.

### Substitutions in random motifs

While the simple constraints show improvement in the scores of random motifs, none performed as well as the motif based on expert knowledge. This difference is, in large part, due to the position-specific substitutions allowed in the expert motif. Three of the five residues in the expert derived motif allow a specific list of substitutions, and these substitutions are important for the high score of the motif. Not allowing these substitutions lowers its ROC10 area from 0.97 to 0.27. To provide randomly generated motifs this same flexibility, I allowed for position-specific substitutions chosen from the same multiple sequence alignment I used to measure conservation. Positions with poor conservation would have a very long allowable substitution list and could match most any residue in any protein, so it only makes sense to use conserved residues with this

substitution scheme. Choosing substitutions from the BLAST-generated multiple sequence alignments showed a large shift in the middle of the distribution to higher ROC10 areas, but it did not change the maximum score (Figure 3). The sequences in the BLAST alignment are all much more similar than the most distant relatives of mandelate racemase in the superfamily. These substitutions allow more frequent matching of the relatively close structures, but not the more distant ones. Capturing the substitutions necessary to match more distant relations will require a multiple sequence alignment that includes sequences that are more distant. Alignments of entire superfamilies are not accurately generated by automatic methods, but the use of a manually curated multiple sequence alignment (Babbitt et al. 1996) shows an increase in the maximum scores achieved by randomly generated motifs.



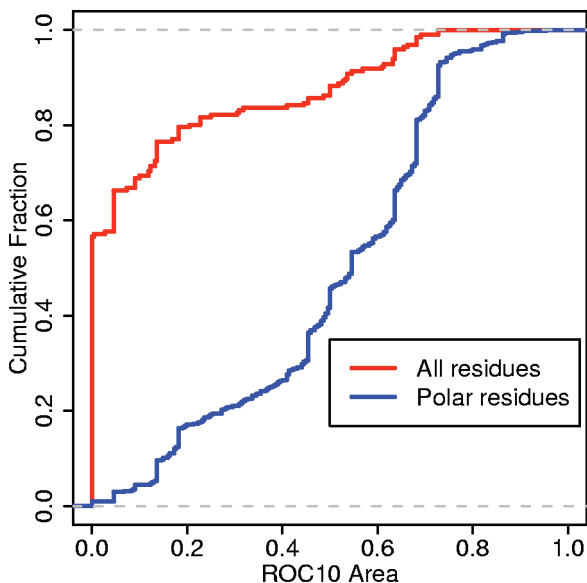
**Figure 3. Cumulative histograms of scores of random motifs with allowed substitutions.**

Red, “No subs.” line is identical to blue “Conserved w/in 7.5 Ång.” line in Figure 2. All motifs shown are constrained by conservation and spatial proximity (7.5 Å neighborhood). Green, “Blast subs.” line is the identical set of motifs with substitutions chosen from a BLAST-generated alignment. Blue, “Expert subs.” line represents 260 motifs generated identically except that the alignment is a manually curated superfamily alignment.

### Limiting Residue Types in Random Motifs

One notable feature of the distributions of motifs based on very diverse sequence alignments is the number of motifs with ROC10 areas at 0. Most of these motifs are sensitive enough to match the four structures that are very similar to 2mnr with low RMSD, but they also match many false positives at equivalently low or lower RMSD. Inspection of these motifs shows that these are composed mostly of hydrophobic residues that are freely substituted by other hydrophobic residues, especially at great evolutionary distances. It appears that matching a hydrophobic cluster is very easy among unrelated proteins. Furthermore, most previously described motifs and catalytic sites are composed of polar residues. As a final constraint to test, I restricted the motifs to use only the polar

residues. This eliminated the large number of motifs that score at 0.0, and shifted the entire distribution to the higher ROC10 areas.



**Figure 4. Cumulative histograms of scores of random motifs with only polar residues.**

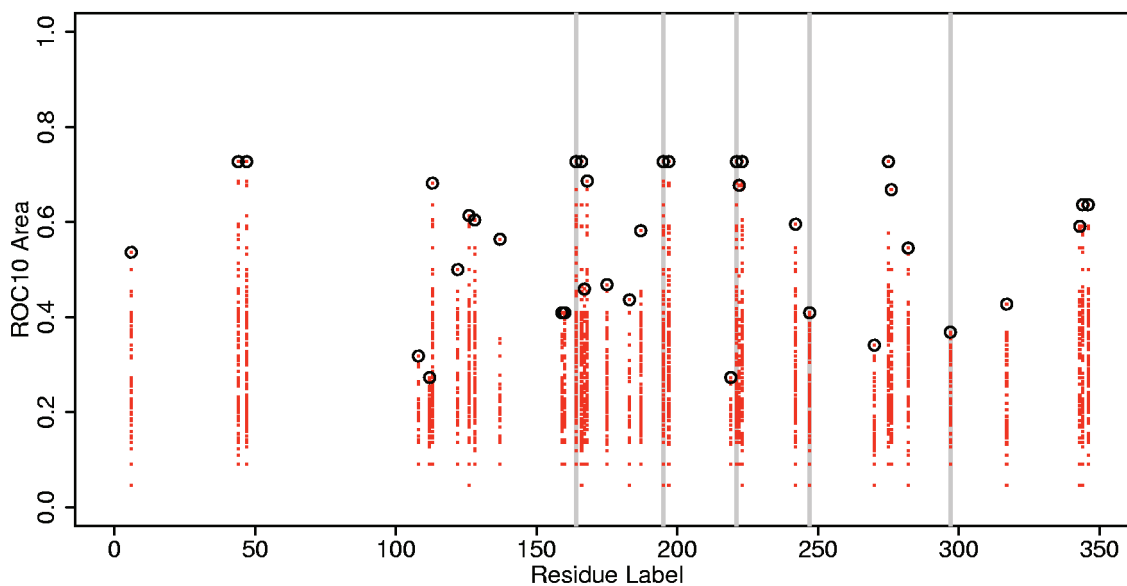
All motifs were generated with conservation and spatial proximity constraints (12 Å), as well as substitutions chosen from a manually curated superfamily multiple sequence alignment. The motifs represented by the blue, “Polar residues” line were further constrained to include only polar residues.

### Important Residues in Motifs

The analysis so far has focused mostly on the scores of motifs and less on the features of the motifs that contribute to the score. I have shown that no randomly generated motifs have classification ability as high as a motif based on the functional site, but many still have high classification ability. Are there other regions of the protein with high classification ability. Are there other regions of the protein with high classification ability that are not in the active site? I examined the residues used in motifs constrained by only conservation to find other residues that could contribute to high scores. Figure 5 shows each residue and an ROC10 area for each motif that contains it. While a functional residue is not sufficient for a high scoring motif—the spread of scores

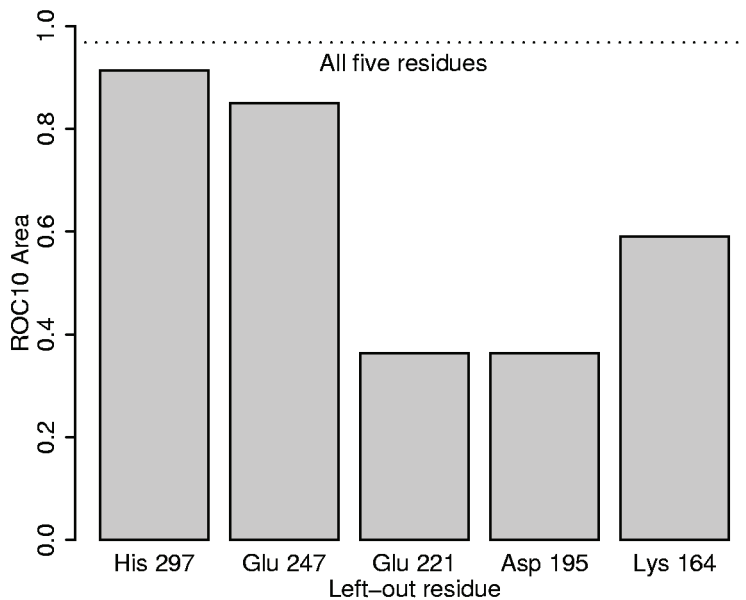
for all positions goes from close to zero to near the maximum for that position, at least some functional residues appear necessary for the highest scoring motifs. The highest-scoring randomly generated templates were very similar to the manually chosen template; all included at least one of the metal ligands, and most included two. Three of the five residues in the expert motif appear among the top scoring motifs, two metal binding ligands and a base. Six other residues appear in top-scoring motifs. Half are close sequence neighbors of these three functional residues. The remaining three top scoring residues are glycines that are more distant from the active site, though they always occur in top scoring motifs with at least one of the expert motif residues. As a final demonstration of the importance of these three functional residues to high scoring motifs, the ROC10 score for motifs built with only four of the expert motif residues show the greatest decline when the left-out residue is one of these three (Figure 6).





**Figure 5. Residues that contribute to motif scores.**

For each motif among the 2825 motifs labeled “Conserved” in Figure 1, their ROC10 area score is plotted against each motif residue (red dots). The gray vertical lines show the residues that make up the motif based on functional knowledge, Lys 164, Asp 195, Glu 221, Glu 247 and His 297. The top score for each residue is circled in black.



**Figure 6. Scores of partial motifs based on the functional site.**

Shown are the scores of five motifs made by leaving one residue out of the motif based on functional knowledge. The ROC10 area for the entire five-residue motif is shown by the dotted line labeled “All five residues”.

## ***Discussion***

The original purpose of this work was to show that expert knowledge of functionally important residues generated unique and concise 3D motifs in enzyme superfamilies. I have shown that very few randomly generated motifs approach the ability of active site motifs to uniquely identify a whole superfamily. Even when motifs are chosen based on constraints to make them more like active site motifs, randomly generated motifs only approach the ability of the expert motif when they contain most of the residues in an expert motif. It should be pointed out that these results are from a very limited data set: all analysis is based on motifs from just one structure in one group of proteins. While preliminary analysis of another structure (1ebh) in the enolase superfamily revealed similar trends (data not shown here), other proteins in other groups could potentially generate very different conclusions.

A major factor explaining the inability of any of the random motifs to match the effectiveness of the expert motif is the allowed substitutions at each position. While I tried different multiple sequence alignments for choosing allowed substitutions, I could not recreate the list of allowed substitutions that can be generated by detailed manual analysis. Part of the problem is inherent to residue-based motifs. Because multiple amino acids can provide the same chemical capabilities, and a single amino acid can provide multiple different chemical capabilities, residue based motifs are a simplified model of functional protein elements. On the other hand, the residues are the ‘atoms’ of protein evolution: individual changes are at the level of individual residues, not chemical groups. Nevertheless, the requirements of function often allow residues to be replaced based on the chemical requirements. For example, when a glutamate is required for its carboxylate

group, often an aspartate can serve as well in the same location. In this simple case, the use of substituting residue types can adequately describe the system, but there are cases that are more complex. For example, two residues in one structure might interact to perform the role of a single residue in a different structure, as a His-Asp dyad serves the role of a lysine in some members of the enolase superfamily (Babbitt et al. 1996). Active site descriptions that use chemical groups or physical and chemical descriptions could avoid this issue.

These results suggest that an automated method to choose a motif could easily be developed based on optimization of random guesses. It is promising that motifs that contain only a fraction of the most important residues, together with other less important residues, produce an intermediate score (see Figures Figure 5 and Figure 6). This could allow an optimization to incrementally learn the best motif through small changes. In my later work, I chose to optimize these motifs through a genetic algorithm (described in Chapter 2), and these results helped guide the development of that genetic algorithm. As a genetic algorithm relies on random guesses, making better random guesses could help reach a solution faster. It is important to recognize a balance though between applying constraints that provide useful limits as opposed to restrictive limits. As functional residues or simply residues that are useful classifiers should always be less variable than residues under no selective pressure, requiring a minimal level of conservation for motif residues was a useful limit. Requiring close spatial proximity of residues, on the other hand, could be an overly restrictive limit, even though it provides for a higher rate of high scoring guesses. While I see the best enrichment in high scoring guesses by restricting motif residues to a 7.5 Å neighborhood, this would eliminate the best observed motifs in

some cases. For my genetic algorithm, I only restrict initial guesses to a 12 Å neighborhood and allow the optimizations to ignore residue distances. Likewise, for the final version of my genetic algorithm the restriction by residue type to polar residues was not used.

## References

- Altschul, S. F., T. L. Madden, et al. (1997). "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res **25**(17): 3389-402.
- Artymiuk, P. J., A. R. Poirrette, et al. (1994). "A graph-theoretic approach to the identification of three-dimensional patterns of amino acid side-chains in protein structures." J Mol Biol **243**(2): 327-44.
- Babbitt, P. C. (2003). "Definitions of enzyme function for the structural genomics era." Curr Opin Chem Biol **7**(2): 230-7.
- Babbitt, P. C., M. S. Hasson, et al. (1996). "The enolase superfamily: a general strategy for enzyme-catalyzed abstraction of the alpha-protons of carboxylic acids." Biochemistry **35**(51): 16489-501.
- Berman, H. M., J. Westbrook, et al. (2000). "The Protein Data Bank." Nucleic Acids Res **28**(1): 235-42.
- Fetrow, J. S. and J. Skolnick (1998). "Method for prediction of protein function from sequence using the sequence-to-structure-to-function paradigm with application to glutaredoxins/thioredoxins and T1 ribonucleases." J Mol Biol **281**(5): 949-68.
- Gerlt, J. A. and P. C. Babbitt (2001). "Divergent evolution of enzymatic function: mechanistically diverse superfamilies and functionally distinct suprafamilies." Annu Rev Biochem **70**: 209-46.
- Holm, L. and C. Sander (1998). "Removing near-neighbour redundancy from large protein sequence collections." Bioinformatics **14**(5): 423-9.
- Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." J Mol Biol **285**(4): 1887-97.
- Lichtarge, O., H. R. Bourne, et al. (1996). "An evolutionary trace method defines binding surfaces common to protein families." J Mol Biol **257**(2): 342-58.
- Meng, E. C., B. J. Polacco, et al. (2004). "Superfamily active site templates." Proteins **55**(4): 962-76.
- Ponomarenko, J. V., P. E. Bourne, et al. (2005). "Assigning new GO annotations to protein data bank sequences by combining structure and sequence homology." Proteins **58**(4): 855-65.

- Russell, R. B. (1998). "Detection of protein three-dimensional side-chain patterns: new examples of convergent evolution." J Mol Biol **279**(5): 1211-27.
- Valdar, W. S. (2002). "Scoring residue conservation." Proteins **48**(2): 227-41.
- Wallace, A. C., N. Borkakoti, et al. (1997). "TESS: a geometric hashing algorithm for deriving 3D coordinate templates for searching structural databases. Application to enzyme active sites." Protein Sci **6**(11): 2308-23.

## Introduction to Chapter 2

The work described in Chapter 2 is a natural progression on the work of Chapter 1. With the knowledge of how to make the best randomly derived motifs, I next sought to optimize the best random motifs to build towards motifs as good as or better than expert-built motifs at identifying a group of proteins with similar functions. For the optimization, I chose to use a genetic algorithm. Genetic algorithms discover optimized solutions to problems by randomly modifying and combining the best observed guesses in an iterative process. In my case, solutions are 3D motifs, or simply a collection of residue coordinates chosen from a protein structure. They can be modified by adding or removing residues, and recombined by taking a random subset of the combination of two motifs. It was a straightforward step to take my programming work used in Chapter 1 to generate random guesses, and put it together with these simple modification and recombination routines together with a slightly modified scoring function to generate a simple but effective genetic algorithm. The results, as described in more detail in the body of the chapter, show that while the expert built motifs still showed advantages for certain structures in certain groups, the genetic algorithm can build useful and functionally related motifs that are often as good as the expert built motifs.

The remainder of this chapter represents a verbatim copy of this manuscript, including the supplementary materials, in the journal *Bioinformatics* (Polacco et al. 2006).

Polacco, B. J. and P. C. Babbitt (2006). "Automated discovery of 3D motifs for protein function annotation." *Bioinformatics* **22**(6): 723-30.

## Chapter 2: Automated Discovery of 3D Motifs for Protein Function Annotation

### ***Abstract***

**Motivation:** Function inference from structure is facilitated by the use of patterns of residues (3D motifs), normally identified by expert knowledge, that correlate with function. As an alternative to often limited expert knowledge, we use machine-learning techniques to identify patterns of 3 to 10 residues that maximize function prediction. This approach allows us to test the assumption that residues that provide function are the most informative for predicting function.

**Results:** We apply our method, GASPS, to the haloacid dehalogenase, enolase, amidohydrolase and crotonase superfamilies and to the serine proteases. The motifs found by GASPS are as good at function prediction as 3D motifs based on expert knowledge. The GASPS motifs with the greatest ability to predict protein function consist mainly of known functional residues. However, several residues with no known functional role are equally predictive. For four groups, we show that the predictive power of our 3D motifs is comparable to or better than that of approaches that use the entire fold (CE) or sequence profiles (PSI-BLAST).

## ***Introduction***

The increasing availability of structural data for proteins of unknown function creates a demand for *in silico* methods to infer the function of these proteins using structural information (Teichmann et al. 2001). But while comparison of overall structures can extend homology detection to evolutionary distances where sequence similarity is undetectable (Chothia et al. 1986), fold comparison often does not identify similarities among functionally significant residues or atoms involved in a protein function's mechanism. Together, the coordinates of these residues or atoms can define a 3D motif. There are many available motif-matching methods that can be used to identify a protein with a matching motif and thus a similar function and mechanism (for example, Artymiuk et al. 1994; Fetrow et al. 1998; Barker et al. 2003). Such methods offer useful complements to fold-based homology comparisons, especially in cases where homologs have diverged in function.

In earlier studies, 3D motifs have typically been chosen based on expert knowledge of functionally important residues in enzyme active sites such as the catalytic triad of the serine proteases. These motifs have been successful at identifying specific enzymatic activities (Torrance et al. 2005), binding relationships (Artymiuk et al. 1994), and superfamily membership (Meng et al. 2004). However, in the absence of a large data source of functional information, accumulation of motifs is slow. The catalytic site atlas (CSA) is a new effort to create a comprehensive database of functional information gleaned from the literature (Porter et al. 2004). It currently provides 147 non-redundant active site motifs for enzymes (Torrance et al. 2005). Similarly, Arakaki et al. (2004) presented an automated method that used the functional information in feature records of



the Swiss-Prot database to construct 3D motifs for 162 different enzymes. Even this method is limited by the shortage of functional information in Swiss-Prot. There are numerous other examples of computational approaches to predict functionally important residues (for example, Zvelebil et al. 1988; Elcock 2001; Wangikar et al. 2003), but these may not be accurate enough to translate to useful motifs (see Discussion).

An alternative is the use of automated 3D motif detection methods. These have shown some success, though none has mapped motifs to specific protein functions with the design goal of characterizing novel proteins with high accuracy. PINTS detects repeated patterns of sidechains between pairwise comparisons of diverse structures, and has generated a large set of repeated motifs (Russell 1998). A similar data-mining approach that compares all patterns across an entire library of structures finds the catalytic triads of proteases along with metal binding sites, salt bridges and similar structural features (Oldfield 2002). Although such general structural features do not provide much specific functional information, they dominate the databases of motifs generated by these types of methods.

We present here a new approach for automated 3D motif generation named GASPS (Genetic Algorithm Search for Patterns in Structures). GASPS was developed with two basic design goals. First, for any specified group of proteins, GASPS should find the motif most useful for identifying the group. Second, GASPS should rely as little as possible on knowledge about what is likely a predictive or functionally important residue. We validate the effectiveness of GASPS on four highly divergent groups of enzymes: the convergent serine proteases (SP), the amidohydrolase superfamily (AHS), the enolase superfamily (ES), and the haloacid dehalogenase superfamily (HADS). These motifs

verify that many, but not all of the previously known functionally important residues are the best predictive residues (along with additional unexpected residues). We describe the crotonase superfamily (CS) as an example of a group that is not well suited for characterization by 3D motifs as they are typically defined.

## ***Methods***

### **Motif Representation and Matching**

As an initial test of principle, we adopted the motif model and matching algorithm of SPASM (Kleywegt 1999), although GASPS can be adapted for use with other motif matching algorithms as well. A motif is a small set of residues (<10 for this study) taken from a single chain, here called the query chain. For each position, SPASM requires a matching residue to be of the identical type with no substitutions. Alternatively, a unique set of residues at each position may be specified that can be substituted with no penalty, though in the course of our study we were unable to use this feature effectively (see Supplementary Materials). SPASM models each residue with just two points, backbone Ca and the sidechain geometrical center. SPASM computes a superposition root-mean-squared deviation (RMSD) for each match it finds within user-defined thresholds of RMSD, sidechain distance deviations (SCD), and C $\alpha$  distance deviations (C $\alpha$ D). For this study, thresholds were set to RMSD=3.2Å, SCD=3.8Å, and C $\alpha$ D=5.0Å. SPASM allows the use of several additional constraints that were not used for this study. Only the match with the best RMSD is considered from each structure.

## **GASPS**

GASPS generates motifs by selecting residues from a single query chain. Here, functional sites and motifs that span more than one chain are not directly addressed. These motifs are scored for their ability to accurately discriminate the positive from the negative sets. There are four main components to a GASPS run: query processing, initial guesses, scoring, and refined guesses.

**Query Processing** To limit the search space, only the 100 most conserved residues in the query chain are considered for inclusion in a motif. Conservation is calculated from a multiple sequence alignment by weighting sequences to reduce the effects of redundancy, considering conservative substitutions based on a substitution matrix, and to penalize gaps (Valdar 2002). All multiple sequence alignments were generated by a two-iteration PSI-BLAST (Altschul et al. 1997) search against nrdb90 (Holm et al. 1997) built in February, 2004.

**Initial Guesses** Fifty candidate motifs are initially chosen spread equally across the linear sequence of the query chain to provide coverage of all regions. For each random guess, a first residue is selected from the query chain and then four other residues are randomly chosen such that each  $\alpha$ -carbon is within 12Å of the first  $\alpha$ -carbon.

**Scoring Function** Candidate motifs are scored for their ability to discriminate between the positive and negative proteins based on the best RMSD matches from a SPASM search. The query structure, which is always a perfect match to the motif, is excluded from the positive set. The scoring function is primarily the normalized area under a receiver-operator characteristic (ROC) plot to five false positives (a false positive

rate of  $\sim 0.001$ ). If the sorted RMSD scores for structures in the negative set are ( $f_1, f_2, f_3, \dots, f_n$ ), then this area, called R, can be computed explicitly as:

$$R = \frac{1}{5} \sum_{i=1}^5 \frac{T(f_i)}{T_{\max}}$$

where  $T(f)$  is the number of true positives with a better RMSD match than a given false positive and  $T_{\max}$  is the size of the positive set. R ranges from 0 to 1. Because R is based on discrete counts, different motifs will frequently have identical R scores. To avoid ties, we include an additional term in the GASPS scoring function. This term, S, is the normalized difference in median RMSD between the true positives and false positives, only considering those that score better than the fifth false positive ( $f_5$ ). This can be explicitly defined as:

$$S = \frac{\text{median}(f_{1-5}) - \text{median}(t_{1-m})}{\text{median}(f_{1-5})}$$

where  $t_{1-m}$  is the set of RMSDs from the true positive matches that are better (less) than the fifth false positive ( $f_5$ ). When no true positives are hit ( $R=0$ ), S is set to zero. The overall GASPS score (G) is then the sum of S and R weighted to emphasize the ROC score, and is composed:

$$G = 1.0R + 0.1S$$

**Refined Guesses** The 16 highest scoring motifs of any round are included in the next round and used as parents for constructing 36 novel motifs via one random process: deletion, insertion, mutation or recombination. The only restriction on the new motifs is that they contain at least 3 residues and at most 10. Deletions and insertions generate a new motif by removing or adding a residue to a single parent motif. A mutation is a

combination of a deletion and an insertion. A recombination is a random subset of the combination of two parent motifs. The top-scoring motif after fifty rounds of refinement is considered the final winner. Most GASPS runs in this study took between 12 and 18 hours on a single 2.667 GHz Intel Xeon processor. Most of this time was spent completing the SPASM searches against the negative set, which time scales directly with the number of proteins in the negative set.

## **Structure Library**

Most analyses in this study used a set of structures selected from the Protein Data Bank (PDB) (Berman et al. 2000) to represent all families in The Structural Classification of Proteins (SCOP) version 1.65 (Murzin et al. 1995). The selection algorithm treats each SCOP family individually and has three main goals: 1) mutant removal based on text matching PDB fields, 2) sequence redundancy filter to 40% identity, and 3) favoring the highest quality structures based on resolution. No distinction is made between apo and holo structures. The entire corresponding PDB chains for each of the SCOP domains are included, regardless of similarities at other domains. On SCOP version 1.65, this selection results in 5440 unique domains on 4243 unique chains.

## **Positive and Negative Sets**

We chose five well-characterized positive groups so that all members within each group share a similar function, and this shared function is dependent on known functional residues (Table 1). Definitions for the four superfamilies were taken from the Structure-Function Linkage Database (SFLD) (Pegg et al. 2005). However, the SFLD as yet contains only a few superfamilies, so to mimic a more typical usage of GASPS on less

than perfect classifications, for all five groups of proteins studied here, a positive set of structures was selected based on SCOP superfamily and family classifications (given in Table 1). Each positive set is a subset of the structure library with the modification that all chains within a PDB structure file are included. Sequence identities between all pairs of homologous chains used as query chains range from 14% to 40%. The negative set is the entire structure library, excluding all chains that contain at least one domain meeting the criteria for inclusion in the positive set.

### **Cross-Validation**

Complete rounds of leave-one-out cross-validation were performed for several query structures in each group. For the smaller groups, each structure in the positive set was used once as a query structure. For the larger groups, AHS and SP, a randomly selected subset of the structures was used. For each query structure, all possible positive training sets were produced by excluding one other (non-query) positive structure. The corresponding positive test sets each contained just the excluded structure. Similarly, the negative set was equally divided to produce as many negative test sets as positive test sets. The corresponding negative training sets are simply the entire negative set excluding a negative test set. Using ES as an example, this procedure required 42 runs of GASPS (7 query structures multiplied by 6 left-out positive structures). The reported sensitivity on the test sets is the portion of GASPS runs where the final GASPS motif from each training run was able to discriminate the left-out positive member from the left-out negative test set at an RMSD threshold equal to the RMSD of the fifth-best false positive match on the training sets. Those runs for which the final trained GASPS motif did not score significantly on the training set were excluded.

**Table 1. Functionally Similar Protein Groups**

Group	SCOP groups	N <sup>1</sup>	Known Functional Residues
Amidohydrolase Superfamily	c.1.9	16	(1a4m) H15 H17 H214 H238 D295
Enolase Superfamily	c.1.11	7	(2mnr) K160 D191 E219 D244 K268
Crotonase Superfamily	c.14.1.3	7	(1mj3) backbone A98 G141
Haloacid Dehalogenase Superfamily	c.108.1	12	(1fez) D12 T126 R160 D186 D190
Serine Proteases	b.47.1.1, b.47.1.2, b.47.1.3, c.41.1.1	38	(2hlc) H57 D102 S195

1) Number of non-redundant structures in positive set.

## PSI-BLAST and CE Libraries

For BLAST (Altschul et al. 1997) and PSI-BLAST comparisons with GASPS the libraries were the set of unique chains from the same PDB files used in the positive and negative sets for GASPS (described above). For the Combinatorial-Extension algorithm (CE) (Shindyalov et al. 1998), to avoid computing all-by-all pairwise comparisons, the negative sets were reduced to the likely high-scoring members for each positive group. For most groups, this meant limiting the negative set to those chains with the same SCOP fold as a catalytic domain in the positive set. However, according to SCOP, HADS is the only superfamily of the HAD-like fold, so its negative set for CE was chosen based on CATH (Orengo et al. 1997) instead. For SP, there were an insufficient number of same-fold structures that were not serine proteases to provide negative sets for both the subtilisins and trypsins. An additional SCOP fold (b.43: Reductase/isomerase/elongation factor common domain) was included in the library that commonly scored highly against SP folds according to the CE internet database (<http://cl.sdsc.edu/ce.html>).

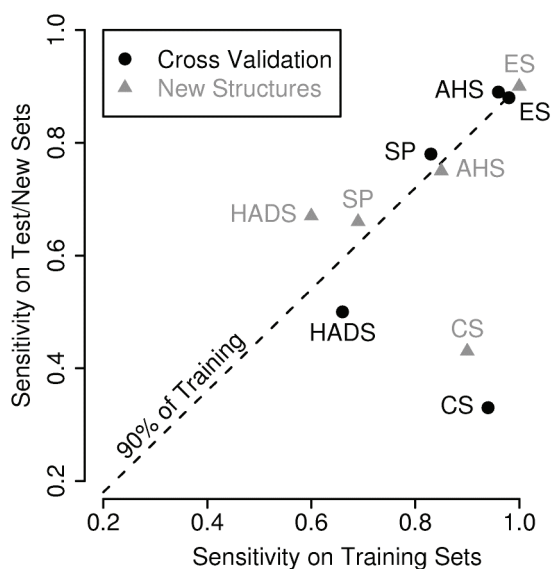
## **Results**

### **Validation of GASPS**

**Significance of Optimized GASPS Scores.** To determine whether any GASPS motif likely represents more than a chance co-occurrence of residues, we computed significance cutoffs from empirical distributions of GASPS motifs due to chance alone. To ensure that any motif was due to chance, artificial positive groups were generated by randomly selecting structures from the structure library, each with a different fold. Based on these distributions, provided in Supplementary Materials, we can set GASPS score thresholds for moderate significance ( $p < 0.01$ ): for groups of approximately 5 structures motifs must score greater than 0.55 and in larger groups of 10 or more structures they must score above 0.4.

**Cross-Validation Studies.** To estimate the performance of GASPS on new proteins, leave-one-out cross-validation studies were completed on each of the groups in Table 1. RMSD thresholds were chosen for each top GASPS motif to give a false positive rate of approximately 0.0013 (5 false positives) on the training set. With the exception of CS, sensitivity is high and there is a close correspondence between the training and test sets (Figure 1). Sensitivity on the test sets for most cases is approximately 90% of that on the training cases. The false positive rate (and its complement, sensitivity) shows an even tighter correspondence with an average rate of 0.0014 on the test cases. The fact that CS is one of the smallest groups and also that it lacks highly conserved sidechains in the active site, as described below, likely contribute to the poor performance of GASPS for this superfamily.





**Figure 1. Generality of GASPS motifs based on sensitivity from two experiments: cross-validation and detection of newer structures.**

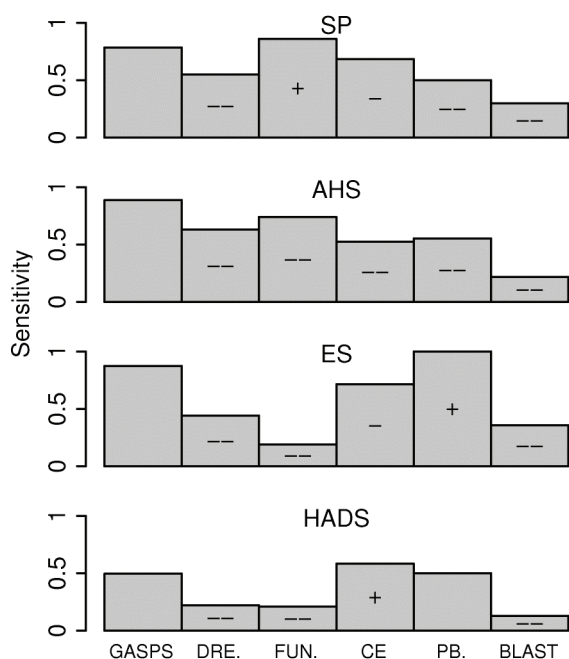
Black filled circles show average sensitivities of motifs from leave-one-out runs on the cross-validation training (x-axis) and test (y-axis) sets. Gray triangles show sensitivities of motifs generated on the full training sets (all motifs in Figure 3) when used to detect structures in the full training set (x-axis) compared with novel structures solved after the training set was established (y-axis).

**Detection of New Structures.** Across all groups, we identified 12 new structures in the PDB that were not yet classified by SCOP (as of version 1.65, December 2003), by a combination of searches based on literature, annotation and sequence similarity, along with communications with collaborators. These 12 proteins all share less than 40% sequence identity with each other or with any protein in the original training set. Motifs generated on the full training set, one for each query structure (shown in Figure 3), were tested for their ability to match the appropriate new structures within the RMSD thresholds determined on the full training set. For these 12 structures, the group-based average rate of matches is 68% compared with 81% on the structures included in the full training set. If CS is excluded, the group-based average rate of matches is 75%, compared to 79% on the training set (Figure 1). This is an average across all motifs in each group

including those with insignificant scores and very poor match rates. The expected match rate for any given motif appears linked to its original GASPS score. Excluding CS, no top-scoring motifs in any one group missed any of the new structures, and only 1 of 9 insignificantly scoring motifs matched any new structures. No new structures, CS included, failed to match any motif in their group.

**Comparisons with other 3D motif methods.** A key benefit of GASPS is that it requires no knowledge of functionally important residues. However, even on groups where functionally important residues are known, GASPS is still useful if it is able to select a more sensitive motif. We constructed motifs built from the functionally important residues (see Table 1) for all possible query structures and compared their sensitivity to GASPS motifs. For all groups except SP, the GASPS motifs have higher sensitivity than simply using these functional residues (Figure 2, “FUN”).

Of other available techniques, the closest to GASPS in principle, is DRESPAT (Wangikar et al. 2003) that detects similar patterns of residues within a group of structures. We used DRESPAT with previously published parameters and a pattern size of four residues to generate patterns for the groups in our data set. The resulting top ranked patterns identify some functionally important residues for all groups in this study except for CS (not shown). However, they fail to identify superfamily members with similar specificity and sensitivities to those of GASPS motifs (Figure 2). It may be possible to adjust the parameters and desired pattern size to improve the performance on a case-by-case basis, but the DRESPAT technique is not designed to automate or aid such a procedure.



**Figure 2. Sensitivity of GASPS motifs compared with other techniques.**

Sensitivity shown for GASPS is measured by cross validation (Fig. 1). For all other techniques, the sensitivity is measured at the threshold of the fifth false positive. Other techniques are DRESPAT (DRE.) motifs, motifs built from functional residues (FUN.), CE, PSI-BLAST (PB.) and BLAST. Plus signs (+) indicate significantly better sensitivity than GASPS within the protein group, and minus signs indicate significantly worse performance at  $p < 0.05$ . Double signs (++) or (--) indicate a greater degree of statistical significance ( $p < 0.0001$ ). CS results are not shown because no 3D motif methods were able to characterize it effectively

Two other 3D motif libraries have recently been used to identify functional or homologous relationships, the motifs used by PINTS and the CSA. As these libraries were not specifically constructed to identify members of the groups in this study, it is impossible to run a parallel experiment for a direct comparison with the techniques shown in Figure 2. PINTS has been used to confirm superfamily membership and binding relationships of structural genomics proteins by finding matches to motifs derived from proximity to ligands or SITE annotations in PDB records (Stark et al. 2004). We tested this same technique (made available at <http://www.russell.embl.de/pints/>) by asking

whether the structures used in our study matched with high specificity those motifs that came from other non-redundant (<40% sequence identity) group members. The measured sensitivities of GASPS motifs (Figures 1 and 2) greatly outperform PINTS for all five groups at similar or even much lower rates of specificity. To be generous, PINTS could adequately serve its purpose if for any query structure it only detects a single true positive motif and ranks it highest among all matches. Even using this much less stringent definition of sensitivity for PINTS than used for GASPS, only for SP does PINTS score better than our reported sensitivities for GASPS motifs.

The motifs derived from functional knowledge in the CSA are available for searching by the program Jess (Barker et al. 2003) at their website (<http://www.ebi.ac.uk/thornton-srv/databases/CSA/>). We used each of the structures in our positive sets to search the CSA with Jess and scored true positives according to whether the motif originated from the same group (defined in Table 1) as the original query. Maintaining similar specificity as required for GASPS, we should require that JESS, with only 147 motifs, identify true positives with greater E-value than any false positive. Only structures from SP reliably matched any true positives. Outside of SP, only three structures (one each from AHS, ES and HADS) matched any CSA motif, but all three of these motifs came from structures that shared more than 40% sequence identity with the query. Relaxing the specificity to five false positives only allowed two other HADS structures to match the haloacid dehalogenase motif. Even though several of the false positive matches had E-values that suggested significance ( $\sim 10^{-4}$ ), none correctly predicted the function or group membership of the original query. These high quality motifs in the CSA are useful for

detecting certain specific functions, but they cannot adequately detect the diverse functions or distant relationships covered by the superfamilies studied here.

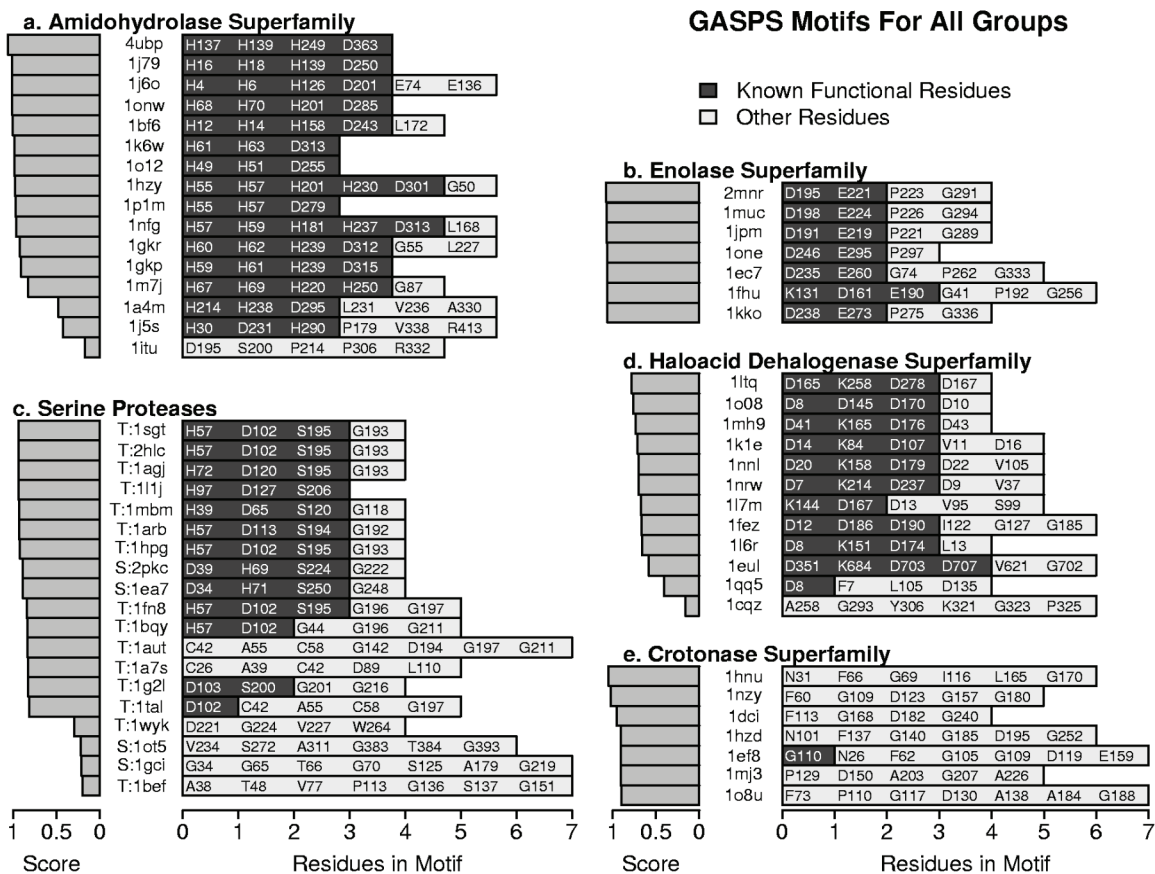
**Prediction Ability Compared to Whole-Chain Tools.** We compared the sensitivity of GASPS motifs (as estimated by cross validation) with other tools that use the whole protein chain including the sequence-based tools BLAST and PSI-BLAST (Altschul et al. 1997) and the fold comparison tool CE (Shindyalov et al. 1998). All members of all positive groups were used as queries for each of the methods, and these were searched against the appropriate library as described in Methods. All sensitivities were measured by counting the fraction of positives that scored better than the fifth best-scoring negative for each query (Figure 2). No single method is better than all other methods for all of the groups in this study. CS is easily grouped by most methods with the exception of GASPS motifs. The fold comparison tool CE performs well on groups with unique folds such as HADS. AHS and ES, on the other hand, share the common  $(\beta/\alpha)_8$  fold with many other superfamilies, which may help explain why CE performs worse than GASPS in these cases. PSI-BLAST performs better than GASPS only for the least divergent superfamilies considered, ES and CS, where PSI-BLAST performs perfectly. With the exception of CS, GASPS motifs outperform BLAST on all groups.

### **Detection of Key Functional Residues**

An advantage of our method is that the selection of the residues in a motif is unbiased towards any preconceived notions of functionally important residues except indirectly via our exclusion of the least conserved residues. This allows us to ask if there is a relationship between the residues that discriminate proteins of related function and the

residues that we know from experimental studies provide function. Table 1 shows the residues that are known to be directly involved in shared function or used in previous functional motif studies. These are used as our gold standard of key functional residues. Every positive structure was used once as a query structure except for SP from which only a diverse subset of structures was chosen. The best motifs from each of these runs are presented in detail in Figure 3. There is a clear trend for the proportion of functional residues in a motif to increase as the motif score rises.

As a stochastic search method, GASPS can be expected to produce different motifs in identically configured runs, and we expect several of the lower scoring motifs presented in Figure 3 are not the best motif a given query can provide. The results of repeated runs for several configurations are presented in Supplementary Materials. Clearly, multiple GASPS runs per group are necessary to ensure that an optimal motif is found for any group. For some single query structures, however, repeated runs suggest there is no combination of residues that provide a useful motif. Meanwhile, the optimal motifs for the majority of other query structures are highly similar. Taken together, these results suggest that to generate a set of the most useful and inclusive motifs for any group of proteins, limited resources are better spent on running GASPS on many different query structures than on running GASPS multiple times on the same structure.



**Figure 3. Scores and functional significance of GASPS motifs.**

The results of a single GASPS run are presented for each named query structure. Residues in the motif that correspond to previously identified functional residues or known active-site motif residues are darkly shaded. All other residues are lightly shaded regardless of subsequent determination of their functional significance. For the serine proteases, query structures are labeled “T:” to denote trypsin-like folds or “S:” for subtilisin-like folds.

**Amidohydrolase Superfamily** The amidohydrolase superfamily (AHS) is a functionally diverse superfamily composed of homologs with a  $(\beta/\alpha)_8$  (TIM) barrel fold that share a conserved mechanistic step mediated by a conserved set of active site residues (Holm et al. 1997; Gerlt et al. 2003). All known members of the superfamily are metal-dependent and require either one or two divalent metal ions. Five conserved metal ligands comprising four histidines and an aspartic acid have been identified as functionally important in all groups within the superfamily. Only one GASPS run on this

superfamily resulted in a motif with an insignificant score and no overlap with any of these metal ligands (Figure 3a). The remaining runs all resulted in motifs that contained at least three of the five conserved ligands. The other residues in the significant motifs are all distant from the metal ligands and thus probably not directly involved in the enzyme's active site.

**Enolase Superfamily** Like the amidohydrolase superfamily, the enolase superfamily (ES) is made up of homologs with a C-terminal  $(\beta/\alpha)_8$  barrel fold plus an N-terminal domain representing a unique fold. All functionally diverse members share a common mechanistic step (Babbitt et al. 1996; Gerlt et al. 2005). Past studies have carefully documented conserved elements responsible for the shared aspects of mechanism, and motifs based on this functional information have been generated with success (Meng et al. 2004). In this study, the conserved residues considered to play a functional role consist of three divalent metal ligands and two basic residues. All motifs resulting from GASPS runs contained at least the same two metal ligands, and one run contained one of the basic residues (Figure 3b). The remaining metal ligand and both basic residues are known to have variable residue types across members of the superfamily, possibly explaining why GASPS has trouble locating them. A highly conserved residue among the GASPS motifs that has not been identified as functionally important is a proline that is two positions downstream from the second metal ligand. Here called the “downstream proline”, it appears in all ES motifs.

**Haloacid Dehalogenase Superfamily** The haloacid dehalogenase superfamily (HADS) comprises enzymes with diverse functions, yet all members share a common mechanistic step associated with hydrolytic nucleophilic substitution via a conserved



aspartate and a few other residues (Allen et al. 2004). The HADS fold is unique according to SCOP, though CATH divides it into two domains: a common Rossmann fold domain and a domain unique to the superfamily. Our laboratory has previously developed motifs in a manual process based on expert knowledge (Meng et al. 2004), and the residues in these motifs are here considered the conserved functional residues. While the catalytic roles may be conserved at each of these positions, all but the obligate aspartate are substituted in diverse members of the superfamily, as apparently required to accommodate differences in their specific mechanisms and overall functions. Despite these substitutions, most GASPS runs still contain three of the five functional residues (Figure 3d). The nucleophilic aspartate appears in all significant motifs where possible. (The 117m structure contains two alternate conformations listed for this aspartate, D11, which precluded it from inclusion in a motif.) Nearly as frequent as the nucleophilic aspartate is another aspartate two positions downstream that has been implicated by others in binding and protonation of the substrate (Allen et al. 2004).

**Serine Protease Families** The serine proteases (SP) are a polyphyletic group consisting mainly of two non-homologous families: the subtilisins and trypsins. They are grouped together by virtue of their common functions and use of a structurally similar catalytic triad in their active sites that appear to be the result of convergent evolution (Dodson et al. 1998). Slightly more than half of the motifs and the highest scoring (10 of 19) included the entire triad (Figure 3c). Most triad-containing motifs included only one additional residue: a glycine involved in formation of the conserved oxyanion hole (for example, 2hlc G193) in trypsins. Though this glycine matches a conserved glycine in the subtilisins, the NH group in the subtilisins points away from the active site cavity. Of the

nine remaining motifs, four had insignificant scores, three included partial catalytic triads, and one was built from a heparin binding protein (1a7s) that, despite its homology to the trypsins, does not contain the catalytic triad or perform protease activity. Many significant runs seemed to be distracted by a disulfide bridge and neighboring alanine near the active site (C42-C58, A55, see Figure 3c and Supplementary Materials), which are well conserved among the trypsins but not the subtilisins.

**Crotonase Superfamily** Members of the crotonase superfamily (CS) display great catalytic diversity, yet all are unified by a common structure-based stabilization of an enolate anion intermediate of acyl-CoA substrates (Holden et al. 2001). Unlike the other groups given in Table 1, however, this shared chemistry is not performed by sidechains but by two structurally conserved NH groups of the peptide backbone that function as part of an oxyanion hole. The sidechains of these residues are not strictly conserved across the superfamily nor are there any other sidechains known or predicted to act in catalysis that are conserved across the entire superfamily. The crotonase superfamily therefore provides a test of GASPS and sidechain-based motifs on a group that may not contain a structural motif focused on sidechains. As expected, an insignificant number of residues in the motifs (1 of 33, for all motifs) is involved in the formation of the characteristic oxyanion hole (Figure 3e). The common residues in the motifs that do discriminate this superfamily seem unlikely to play a direct role in the enzyme's function, based on their distance from the active site. Examples include a conserved phenylalanine (1hnu F66) that is buried but lines an interior cavity and an aspartate (1hnu D135) involved in a conserved salt bridge.

## ***Discussion***

### **Using GASPS for Function Identification**

The performance of GASPS-generated motifs is comparable to that of 3D motifs generated based on expert knowledge of functional sites in other proteins (Artymiuk et al. 1994; Wallace et al. 1997; Fetrow et al. 1998; Kleywegt 1999; Torrance et al. 2005). Furthermore, GASPS motifs improve the coverage of protein functions offered by publicly available sources of 3D motifs (Stark et al. 2003; Porter et al. 2004). Searching with protein fragments in three-dimensional motifs developed by GASPS was also found to be comparable or better than commonly used methods of annotation transfer that use an entire protein chain such as PSI-BLAST or CE. Unlike these methods that use an entire protein or domain, GASPS is able to focus on the features of protein structure most likely to tell us the most about protein function. GASPS therefore provides a method of generating motifs useful for function or superfamily prediction in an automated fashion with no prior knowledge of mechanistic details. Such motifs can be used to verify similarity of active sites in proteins in which only similarity of fold has been previously identified. For example, GASPS motifs could be used for distinguishing functional differences among families of  $(\beta/\alpha)_8$  proteins.

GASPS requires only a prior grouping of related proteins, so GASPS is limited only to groups with sufficient available structures. We cannot say for certain how many structures are required, but it appears to depend on the variability among the available structures. In the current study, all structures shared less than 40% sequence identity, and GASPS still was able to find general motifs for groups with as few as seven structures.

While only 14% of superfamilies and 6% of families in the structure library used here have this many structures, these superfamilies and families make up the majority of protein structures (60% and 32%, respectively). Theoretically, it would seem possible for two highly diverged structures to share only their unique functional motif. However, for most proteins, even of different folds, it appears that sharing similar residues in 3D space occurs frequently enough by chance alone to require more than two structures to produce a trusted motif (Wangikar et al. 2003).

SPASM (and therefore GASPS, as used in this study) considers only a single point for each Ca and sidechain. With most catalysis carried out by sidechains (Bartlett et al. 2002), we believe the inclusion of the sidechains allows for better characterization of functional sites than if only the backbone placement were considered. Motifs could be represented with more precision by using the location of chemical groups, or even individual functional atoms. However, given the variability in placements of functional atoms in crystallographic structures, (DePristo et al. 2004; Torrance et al. 2005), approximating the entire sidechain by a single rigid point may be more appropriate.

## **Location of Functional Information**

GASPS makes no assumptions about the location of functional information except that it can be resolved to individual residues and that it will be relatively well conserved. The observed correspondence between information useful for classification and functionally significant residues is a result of the choice of positive sets based on shared chemical activities used in this study. The use of GASPS on sets based on other shared characteristics, such as homology, binding partners, or cofactors, may identify the residues most attributable to those shared characteristics.

It should be noted that the motifs generated by GASPS may not be the only, or even the most informative structural features. GASPS is expected to miss informative structural features if the features are either inconsistent between members of the group, such as the substituted residues in HADS, or if the features are not based on individual sidechains such as backbone interactions or helix dipoles. The CS results provide a case in point.

In addition to the previously identified functionally important positions, other positions occur with high frequency among the motifs for these groups. These positions may, for example, merely provide a simple geometric positioning constraint for the other motif residues that aids specificity. However, based on their conservation in 3D space, these positions are likely to play an important role for the protein, especially when located in the active site. For example, when the conserved “downstream proline” in the enolase superfamily is mutated to alanine in the muconate lactonizing enzyme from *Pseudomonas putida* (equivalent to structure 1muc) it results in an insoluble protein, (R. Nagatani and P. Babbitt, personal communication,) suggesting that this proline may be important for folding or stability of the soluble globular protein.

Based on its ability to identify at least a subset of the functionally important residues, GASPS appears similar to the fully automated DRESPAT, which successfully locates functionally important residues by identifying shared structural patterns in a set of functionally similar protein structures (Wangikar et al. 2003). The main differences between GASPS and DRESPAT are that GASPS compares patterns with a negative set and chooses patterns based on their predictive ability. Wangikar et al.(2003) suggest that DRESPAT patterns may represent useful 3D motifs. However, in the course of this study

we found that when DRESPAT patterns were converted to motifs for use by the search tool SPASM, they were not as accurate as those motifs generated by GASPS.

## **Inference of Function for Diverse Groups**

Four of the five groups in this study have been described as “mechanistically diverse superfamilies” (Gerlt et al. 2001) consisting of divergent enzymes that perform many different overall biochemical functions, but utilize a common mechanistic step such as a partial reaction. Any motifs that identify proteins to these groups will therefore identify the shared mechanistic step but not the overall biochemical function. By mapping a specific mechanistic step to specific structural elements, we are using a finer-resolution view of protein function than overall biochemical function, but one that is more appropriate for such diverse groups (Babbitt 2003).

## **Future Applications**

If applied to an exhaustive functional classification of proteins, GASPS has the potential to generate an unbiased set of 3D motifs that can aid in function prediction for novel proteins. In addition to aiding protein classification, the collection of 3D motifs can represent hypotheses about determinants of function shared among related proteins. In this regard, the high-scoring motifs can serve as starting points for studies attempting to link function to structure, especially in a superfamily context. Additionally such a study would systematically investigate the utility of 3D motifs at identification of functions other than catalysis, such as ligand binding.

For groups with few experimental structures available, especially those coming from structural genomics initiatives, GASPS would have insufficient structures without the use

of predicted structures, generated by homology modeling, for example. Past work has specifically demonstrated the effectiveness of predicted structures for matching previously determined functional motifs (Arakaki et al. 2004). It remains unclear whether predicted structures can be used reliably for generating motifs. Work is ongoing in our laboratory to investigate this issue.

## ***Acknowledgements***

We thank Elaine Meng for her useful discussions and careful reading of an earlier version of the manuscript. This work was supported by NSF grant DBI-0234768.

## ***References***

- Allen, K. N. and D. Dunaway-Mariano (2004). "Phosphoryl group transfer: evolution of a catalytic scaffold." Trends Biochem Sci **29**(9): 495-503.
- Altschul, S. F., T. L. Madden, et al. (1997). "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res **25**(17): 3389-402.
- Arakaki, A. K., Y. Zhang, et al. (2004). "Large-scale assessment of the utility of low-resolution protein structures for biochemical function assignment." Bioinformatics **20**(7): 1087-96.
- Artymiuk, P. J., A. R. Poirrette, et al. (1994). "A graph-theoretic approach to the identification of three-dimensional patterns of amino acid side-chains in protein structures." J Mol Biol **243**(2): 327-44.
- Babbitt, P. C. (2003). "Definitions of enzyme function for the structural genomics era." Curr Opin Chem Biol **7**(2): 230-7.
- Babbitt, P. C., M. S. Hasson, et al. (1996). "The enolase superfamily: a general strategy for enzyme-catalyzed abstraction of the alpha-protons of carboxylic acids." Biochemistry **35**(51): 16489-501.
- Barker, J. A. and J. M. Thornton (2003). "An algorithm for constraint-based structural template matching: application to 3D templates with statistical analysis." Bioinformatics **19**(13): 1644-9.
- Bartlett, G. J., C. T. Porter, et al. (2002). "Analysis of catalytic residues in enzyme active sites." J Mol Biol **324**(1): 105-21.

- Berman, H. M., J. Westbrook, et al. (2000). "The Protein Data Bank." Nucleic Acids Res **28**(1): 235-42.
- Chothia, C. and A. M. Lesk (1986). "The relation between the divergence of sequence and structure in proteins." Embo J **5**(4): 823-6.
- DePristo, M. A., P. I. de Bakker, et al. (2004). "Heterogeneity and inaccuracy in protein structures solved by x-ray crystallography." Structure (Camb) **12**(5): 831-8.
- Dodson, G. and A. Wlodawer (1998). "Catalytic triads and their relatives." Trends Biochem Sci **23**(9): 347-52.
- Elcock, A. H. (2001). "Prediction of functionally important residues based solely on the computed energetics of protein structure." J Mol Biol **312**(4): 885-96.
- Fetrow, J. S. and J. Skolnick (1998). "Method for prediction of protein function from sequence using the sequence-to-structure-to-function paradigm with application to glutaredoxins/thioredoxins and T1 ribonucleases." J Mol Biol **281**(5): 949-68.
- Gerlt, J. A. and P. C. Babbitt (2001). "Divergent evolution of enzymatic function: mechanistically diverse superfamilies and functionally distinct suprafamilies." Annu Rev Biochem **70**: 209-46.
- Gerlt, J. A., P. C. Babbitt, et al. (2005). "Divergent evolution in the enolase superfamily: the interplay of mechanism and specificity." Arch Biochem Biophys **433**(1): 59-70.
- Gerlt, J. A. and F. M. Raushel (2003). "Evolution of function in (beta/alpha)<sub>8</sub>-barrel enzymes." Curr Opin Chem Biol **7**(2): 252-64.
- Holden, H. M., M. M. Benning, et al. (2001). "The crotonase superfamily: divergently related enzymes that catalyze different reactions involving acyl coenzyme a thioesters." Acc Chem Res **34**(2): 145-57.
- Holm, L. and C. Sander (1997). "An evolutionary treasure: unification of a broad set of amidohydrolases related to urease." Proteins **28**(1): 72-82.
- Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." J Mol Biol **285**(4): 1887-97.
- Meng, E. C., B. J. Polacco, et al. (2004). "Superfamily active site templates." Proteins **55**(4): 962-76.
- Murzin, A. G., S. E. Brenner, et al. (1995). "SCOP: a structural classification of proteins database for the investigation of sequences and structures." J Mol Biol **247**(4): 536-40.
- Oldfield, T. J. (2002). "Data mining the protein data bank: residue interactions." Proteins **49**(4): 510-28.
- Orengo, C. A., A. D. Michie, et al. (1997). "CATH--a hierarchic classification of protein domain structures." Structure **5**(8): 1093-108.
- Pegg, S. C., S. Brown, et al. (2005). "Representing structure-function relationships in mechanistically diverse enzyme superfamilies." Pac Symp Biocomput: 358-69.



- Porter, C. T., G. J. Bartlett, et al. (2004). "The Catalytic Site Atlas: a resource of catalytic sites and residues identified in enzymes using structural data." Nucleic Acids Res **32 Database issue**: D129-33.
- Russell, R. B. (1998). "Detection of protein three-dimensional side-chain patterns: new examples of convergent evolution." J Mol Biol **279**(5): 1211-27.
- Shindyalov, I. N. and P. E. Bourne (1998). "Protein structure alignment by incremental combinatorial extension (CE) of the optimal path." Protein Eng **11**(9): 739-47.
- Stark, A. and R. B. Russell (2003). "Annotation in three dimensions. PINTS: Patterns in Non-homologous Tertiary Structures." Nucleic Acids Res **31**(13): 3341-4.
- Stark, A., A. Shkumatov, et al. (2004). "Finding functional sites in structural genomics proteins." Structure (Camb) **12**(8): 1405-12.
- Teichmann, S. A., A. G. Murzin, et al. (2001). "Determination of protein function, evolution and interactions by structural genomics." Curr Opin Struct Biol **11**(3): 354-63.
- Torrance, J. W., G. J. Bartlett, et al. (2005). "Using a Library of Structural Templates to Recognise Catalytic Sites and Explore their Evolution in Homologous Families." J Mol Biol **347**(3): 565-81.
- Valdar, W. S. (2002). "Scoring residue conservation." Proteins **48**(2): 227-41.
- Wallace, A. C., N. Borkakoti, et al. (1997). "TESS: a geometric hashing algorithm for deriving 3D coordinate templates for searching structural databases. Application to enzyme active sites." Protein Sci **6**(11): 2308-23.
- Wangikar, P. P., A. V. Tendulkar, et al. (2003). "Functional sites in protein families uncovered via an objective and automated graph theoretic approach." J Mol Biol **326**(3): 955-78.
- Zvelebil, M. J. and M. J. Sternberg (1988). "Analysis and prediction of the location of catalytic residues in enzymes." Protein Eng **2**(2): 127-38.

## Supplementary Materials for Chapter 2

### ***Significance of Optimized GASPS Scores***

As a way of determining whether any GASPS motif and its score are significant and likely represent more than a chance co-occurrence of residues, we computed empirical distributions of the scores of refined GASPS motifs due to chance alone. To ensure that the motifs were due to chance alone, artificial positive groups were generated by randomly selecting structures from the structure library, each with a different fold. The distributions of final GASPS scores on these artificial positive groups of size 5 and 10 structures are shown in the histograms in Figure i. As expected, these chance motifs do a better job of discriminating the smaller groups (have higher GASPS scores), though in both cases the discrimination is far from perfect. Using these distributions we can set GASPS score thresholds for moderate significance ( $p < 0.01$ ): motifs found in groups of approximately 5 structures must score greater than 0.55 and in larger groups of 10 or more structures they must score above 0.4.

### ***Sources of Variability***

As a stochastic search method, GASPS generates some variability in GASPS scores by producing different motifs in identically configured runs. However, it is clear from other observations that some variability stems from variation between query structures. Recently, Torrance *et al.* (2005) reported variation above 1.0Å RMSD at active site motifs for 20% of protein pairs in a set of proteins with less than 20% sequence identity. This variation causes corresponding motifs from different query structures to provide different

scores. For example, the GASPS score for the catalytic triad of trypsin-like structure 1hpg is 0.83, but the same motif from another trypsin-like protein, 1agj, is 0.92. To determine how much of the variation seen across GASPS results is due to the query structure or the stochastic nature of the genetic algorithm, the best, worst, and the two queries scoring closest to the median were chosen from each of these groups and run through GASPS five additional times.

Figure ii plots the GASPS scores for these runs (along with histograms of insignificant GASPS scores, as in Figure ii for reference). It is encouraging that the highest scores for each query structure seem to be the most frequent in most cases, especially for AHS and HADS. We also note that there are runs for which a high-scoring motif exists as indicated by other successful runs with the same query, but GASPS fails to locate it or any other significant motif; see for example 1hzy from AHS, and 1one and 1kko from ES in Figure ii. There are also structures that seem poorly suited for motif generation with GASPS, see 1itu (AHS), 1cqz(HADS) and 1bef(SP) in Figure ii. Manual examination of these structures suggests the active sites and other structurally conserved sites of these structures may be relatively deformed.

In many cases, while GASPS may fail to find the most significant motif, it can still find a significant motif. For example, three of the six GASPS runs with 2mnr (mandelate racemase from ES) as the query structure do not find the two metal ligands that all top-scoring 2mnr motifs contain, but they do find the conserved lysine (K164) and a nearby lysine (K166) in the active site, along with a distant pair of glycines (G44 and G47) located in the non-catalytic N-terminal domain. Similarly, two runs with 1one (enolase from ES) find another lysine (K345) important for catalysis in enolase(Babbitt et al.

1996) and a glutamine (Q167) that is very close to the active site. As an example from SP, the disulfide bridge near the active site occurred in all less optimal motifs for 1bqy and 2hlc (C42 and C58), usually along with one or more members of the catalytic triad (Figure iii b). Likewise, the less optimal motif for 1nrw in HADS still shares two functional aspartates in common with the best motif (D7 and D237), but matches a different third functional residue (D241 instead of K214) along with three other nearby residues (G42, G236, and A248). Thus, these suboptimal motifs may be useful for identifying additional residues that are potentially important for protein function.

In some cases, however, a suboptimal motif is best viewed as incomplete progression towards the optimal motif. For example, three different motifs that were found for 1ec7 (glucarate dehydratase from ES) all have similar scores and share the two most common metal ligands (D235 and E260). They differ by whether they include either a distant glycine (G115) located in the N-terminal domain, or the “downstream proline” (P262, see the *Enolase Superfamily* section above) with a glycine (G333) located at the N-terminal end of the barrel along with another distant N-terminal glycine (G74) in the third case. Similarly, the less than optimal motif for SP member 1fn8 (catalytic triad plus G196 and G197) is very close to the optimal motifs (catalytic triad plus G193), and GASPS likely would have settled on the optimal motif in this case with a few more rounds of optimization.

Taken together, these observations suggest that to generate a set of the most useful and inclusive motifs, limited resources are better spent on running GASPS on many different query structures than on running GASPS multiple times on the same structure.

## ***Detection of New Unidentified Structures***

The GASPS motifs in Figure 3 were used to search for additional group members among a non-redundant (<50% sequence identity) subset of the complete PDB in late September, 2004. Ten structures, previously unidentified by the authors, were identified by these GASPS motifs and later confirmed as group members by consulting the literature or expert collaborators: 7 from HADS, 2 from SP, and 1 from CS. The rate of new false positives at the chosen thresholds was also consistent with expectations. Thresholds were chosen based on a false positive rate of 0.0012 on the training set (5 false positives from among the original library of 4243 structures), and the rate on the newer subset of the PDB was identical within rounding errors (average 8.0 false positives from 6673 structures).

## ***Allowing Substitutions in Motifs***

In an attempt to improve the ability of GASPS to identify functionally important residues and identify motifs useful for classifying protein structures, we tested a scheme for allowing position-specific substitutions. Allowed substitutions were chosen based on the multiple sequence alignment used to measure conservation, and were fixed for a position throughout the GASPS run.

The introduction of substitutions allows GASPS to achieve much higher scores on the randomly generated, unrelated groups, though only for HAD and Crotonase superfamilies are statistically significant increases in GASPS scores observed (Wilcoxon ranked sum test (Hollander 1973),  $p < 0.001$  and  $p < 0.01$ , respectively. Data not shown.) However, the improvements did not hold up to a full set of cross-validation analyses using all

structures as queries and test structures (Table i). It appears that allowing substitutions raises GASPS scores for some families, especially when substitutions help to identify a few more functionally important residues. However, in its most simple current implementation, allowing substitutions in GASPS appears to be more prone to overfitting. Thus, we cannot say that allowing substitutions in the current scheme provides for more general motifs.

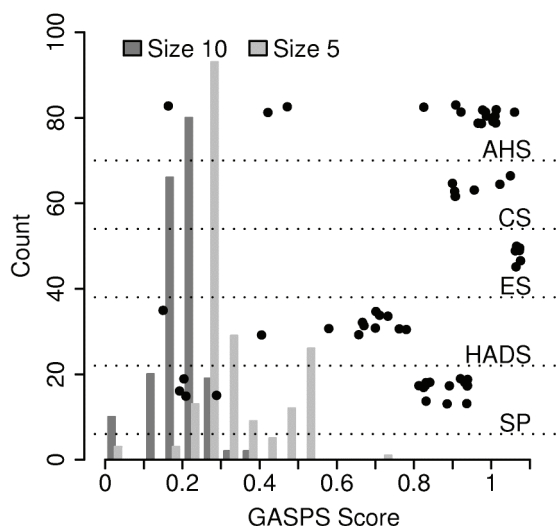
## **References**

- Babbitt, P. C., M. S. Hasson, et al. (1996). "The enolase superfamily: a general strategy for enzyme-catalyzed abstraction of the alpha-protons of carboxylic acids." Biochemistry **35**(51): 16489-501.
- Hollander, M. D. A. W. (1973). Nonparametric statistical inference. New York, John Wiley & Sons.
- Torrance, J. W., G. J. Bartlett, et al. (2005). "Using a Library of Structural Templates to Recognize Catalytic Sites and Explore their Evolution in Homologous Families." J Mol Biol **347**(3): 565-81.

**Table i. Improvements in GASPS by using substitutions on Crotonase and HAD superfamilies.**

	Average Training Scores		Average Test Scores		Number of test structures matched		
	Subs	No Subs	Subs	No Subs	Subs	No Subs	
CS	1.06	0.98	0.35	0.33	14	14	
HADS	0.85	0.63	0.5	0.45	70	59	$p=0.22^a$

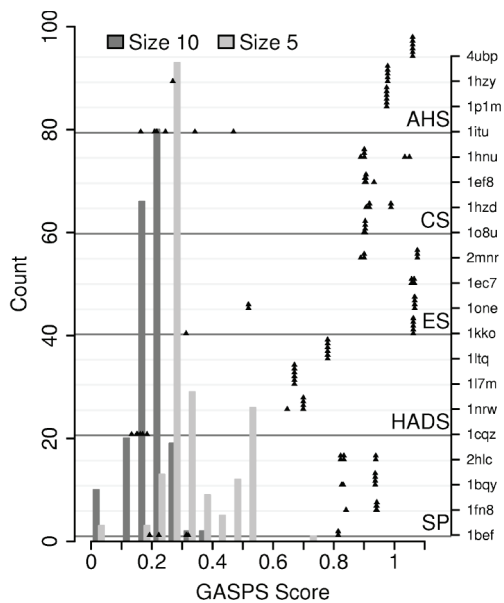
(a) Fisher's exact test on count data for HADS alone.



**Figure i. Distributions of GASPS scores on artificial and real groups.**

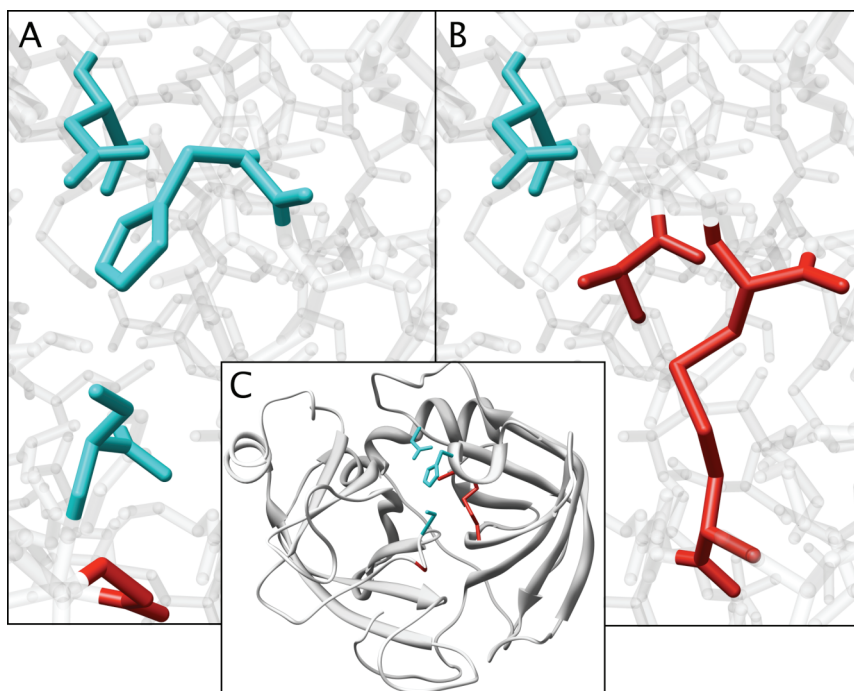
GASPS scores on randomly selected artificial groups of 10 and 5 structures are presented as histograms. GASPS scores on real groups (corresponding to motifs in Figure 2) are presented as scatter plots arranged by group; within groups the vertical placement is randomly chosen to avoid overlaps. Each point is a single GASPS run using a unique query chain. The counts on the y-axis are relevant only to the histograms.





**Figure ii. Stochasticity of GASPS results.**

For each of the five groups, data are presented from six repeated GASPS runs on four query chains indicated by their PDB identifier. Stacked points (triangles) represent identical GASPS results for the query chain. Histograms are redundant with **Error! Reference source not found.**, but included here for reference. The counts on the left y-axis are relevant only to the histograms. The right y-axis and its connecting horizontal lines identify the query structure used to generate the GASPS scores that stack on each line.



**Figure iii. GASPS motifs for 2hlc, a trypsin-like serine protease.**

Residues in motifs are highlighted in red or in cyan if the residue is also part of the catalytic triad. (a.) The top-scoring GASPS motif from among all runs includes the entire catalytic triad (H57, D102, S195) and a nearby glycine (G193). (b.) The top-scoring motif from an identically configured, but lower scoring, GASPS run includes D102 of the catalytic triad, and a nearby disulfide bridge and alanine (C42, C58, A55). (c.) Residues highlighted in panels (a) and (b) are shown relative to the entire domain.

## **Chapter 3: An analysis of computed expectations for random matches to 3D motifs**

### ***Introduction***

The scientific evaluation of any hypothesis requires a method to evaluate the likelihood that observations could be explained by an alternative, usually simpler and less interesting, hypothesis. For searches of biological databases, in our case protein structures or 3D motifs, the hypothesis of interest is that any match is the result of a meaningful biological relationship, such as shared ancestry or function, and the alternative that a match is the result of chance placements of residues within the physical and chemical constraints of protein structure. For GASPS motifs, the biological relationship of interest is a group of proteins, defined by homology or functional similarity (Polacco et al. 2006). GASPS seeks to find a motif where all group members match within a deviation threshold stringent enough to make matches to non-group proteins rare. The original GASPS determines this threshold empirically, by searching for matches to each candidate motif among all non-group proteins that it should not match. Each run of GASPS on a single modern processor can take as long as 20 hours, and the repeated searches of a structure database take approximately 98% of this time. To see if this step could be replaced by a quicker approach, I evaluated a more theoretical approach to calculating the expected number of random matches to any motif. This is not solely a practical, statistical inquiry. If motifs are adequately scored by a theoretical model based only on geometry, this indicates conserved elements are products solely of the unique functional constraints of a group. Instead, I show evidence for the alternative, that the conserved elements of

protein structures are often not unique to a group, so that general physical or chemical constraints of protein structure can lead to similar arrangements of residues in unrelated proteins.

## **Results**

### **Modifying the GASPS scoring function.**

GASPS was designed to find a motif that discriminates between a group of proteins that share a trait of interest (such as the serine proteases) and all other proteins. In effect, GASPS chooses a motif that is matched by all of the group proteins within an RMSD stringent enough to make random matches to other proteins rare. By searching each candidate motif against the background structure library, GASPS empirically computes an expectation for each true positive match based on its relative RMSD. This time-consuming step in GASPS can be replaced by the use of an accurate model that can be used to compute the likelihood of a match within any RMSD threshold. Stark et al. (2003) have developed a method for computing the expected number of matches to a motif based only on the number of residues in the motif, abundance of residues in the database, and number of atoms used per residue in the motif. We compare here the accuracy of using the computed expectation numbers ( $E_c$ ) with actual expectations, or counts of false positives, computed empirically ( $E_e$ ).

The majority of the GASPS score of a motif is based on the number of random matches that are expected at better or equal RMSD than to its matches to group proteins. This score is based on the area under an ROC curve to five false positives and so can be computed by summing the vertical “columns” on an ROC plot (Equation 1).

### Equation 1

$$Area = \frac{1}{5} \sum_{i=1}^5 \frac{T(f_i)}{T_{\max}}$$

$T(f_i)$  is the count of true positives with scores better than false positive  $f_i$ , and  $T_{\max}$  is the maximum number of true positives, i.e. the number of proteins in the positive set. For purposes of this explanation, it will be clearer to instead take the equivalent area by summing the “rows” (Equation 2).

### Equation 2

$$Area = \frac{1}{T_{\max}} \sum_{i=1}^{T_{\max}} \left( 1 - \frac{F(t_i)}{5} \right)$$

$F(t_i)$  is the number of false positive matches with an RMSD equal to or better than that of a given true positive match,  $t_i$ , and is assigned a maximum value of five for the equation above.  $F(t)$  is treated here as the empirical expectation value ( $E_e$ ). The summed term  $(1 - F(t_i)/5)$ , is the “credit” granted to each true positive match  $t_i$  by an ROC area calculation. If a true positive matches better than any false positive, it is given full credit  $(1-0/5)$ , with partial credit granted if it matches worse than only a few ( $x$ ) false positives  $(1-x/5)$ , and no credit  $(1-5/5)$  if it matches worse than 5 or more false positives. Traditionally, a probability ( $P$ ) value can be assigned to any match from an expectation value according to Equation 5. This ROC credit term is approximately equivalent to this  $P$  value subtracted from 1 (see Figure 1), so for explanatory purposes we define an “empirical probability”  $P_e$  by Equation 3.

### Equation 3

$$P_e(t) = \begin{cases} F(t)/5 & \text{for } F(t) \leq 5 \\ 1 & \text{for } F(t) > 5 \end{cases}$$

The ROC area component of the G score calculated by Equation 2 is then just the average of  $(1-P_e)$  for the best match for each protein in the positive set. With the GASPS scoring function now viewed in terms of P values, a theoretical system for computing P values can easily be substituted for the empirical system. Because a theoretical computation of P values should be continuous, the second term used in the empirical scoring function, which served mainly to break ties when the discrete  $P_e$  gave identical scores, is not necessary.

The computed expectation,  $E_c$ , is the expected number of matches in the negative set and can be computed at any RMSD (R) by Equation 4.

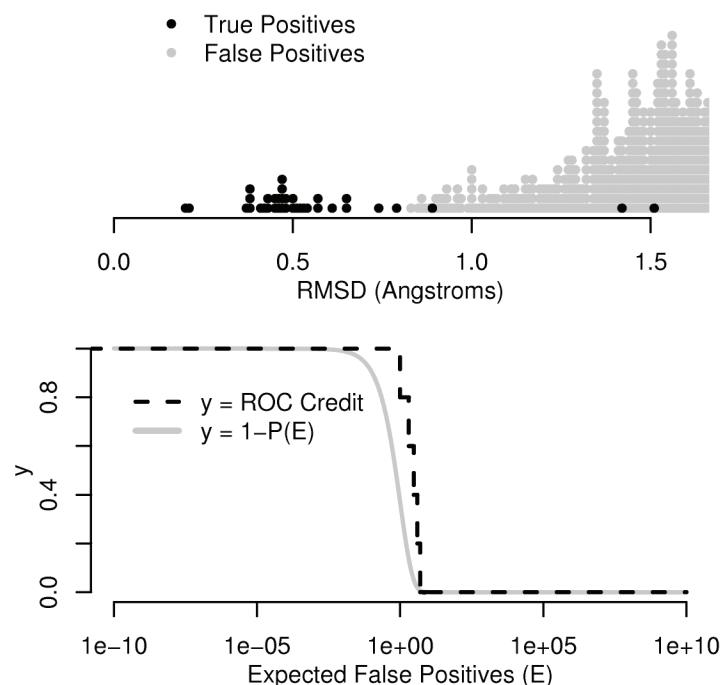
**Equation 4**

$$E_c(R) = DA\Phi a_3^N R^{2.93N-5.88} [c_2 R^2]^N$$

In this formula, taken from Stark et al. (2003), D is the number of proteins in the database, N is the number of residues in the motif,  $\Phi$  is the products of abundances (as percentages) of residue types, and the remaining parameters are empirically derived constants:  $A=3.70 \times 10^6$ ,  $a_3=1.79 \times 10^{-3}$ ,  $c=0.196$ . The expected number of matches is converted to a P value using Equation 5, which depends on a Poisson distribution of matches among the protein structures.

**Equation 5**

$$P = 1 - e^{-E}$$



**Figure 1. Relation between "ROC Credit", P Values, and Expected False Positives.**

These two plots are intended to demonstrate the similarity between calculating a P value based on a Poisson distribution and continuous expectation values, and 'P values' based on number of observed false positives (1- 'ROC Credit'). The top plot shows the distributions of true positive and false positive matches to an actual motif from a serine protease. It is drawn so that its RMSD values roughly correspond to the E values on the lower plot. For both systems of scoring true positives, the majority of true positives (black dots) are scored equivalently. Those on the left are given a score of 1, or approximately 1, and those on the right are given a score of 0, or approximately 0. For this case, only the three true positives near  $E=1$  will be treated differently.

### Expectation values compared.

Figure 1 shows that values of  $P_e$  and  $P_c$  are nearly identical given the same expectation values. Therefore, GASPS modified to use a computed score as opposed to an empirical score should provide a similar result as long as  $E_c$  accurately tracks  $E_e$ . One important difference between  $E_c$  (Stark et al. 2003) and  $E_e$  as defined here is that  $E_c$  predicts the

total number of matches, including multiple non-independent matches in the same protein, whereas  $E_e$  counts only the number of proteins with one or more matches<sup>1</sup>. Multiple matches to a motif frequently occur within the same protein and make use of some of the same residues, but only the best match is usually biologically significant (Meng et al. 2004). Coincident motifs occur more often than expected by chance: assuming 5 matches are assorted randomly across 4000 proteins, the likelihood of any coincident match in any protein is 0.002 (by ‘Birthday Paradox’ arguments). Forty percent of motifs with 6 residues, when searched against the negative set used in Chapter 2, have repeated proteins among their best 5 matches. For smaller motifs, this becomes less of a problem with the same percentages being 22%, 9%, and 5% for motifs with 5, 4 and 3 residues, respectively.

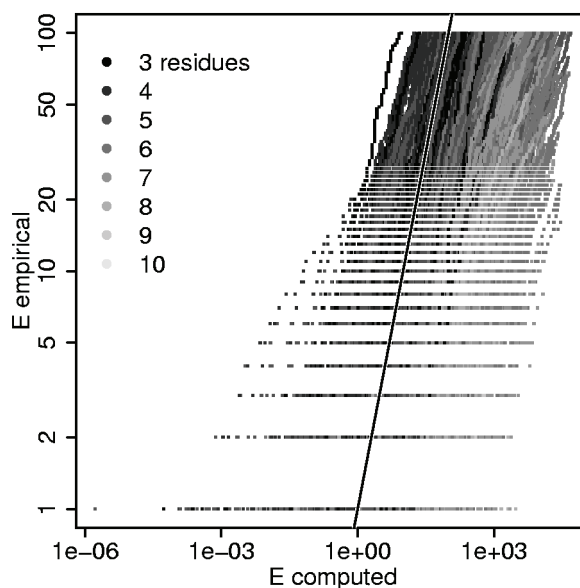
I examined the actual relation between  $E_e$  and  $E_c$  by computing  $E_c$  for each false positive match to a set of motifs, and computing  $E_e$  by counting the number of actual false positive matches that score with the same or lower RMSD. The motifs used were the high-scoring “surviving” motifs from each round of a single GASPS run for two randomly chosen structures from each of the five groups studied in Chapter 3. Figure 2 shows that while  $E_e$  and  $E_c$  tend to agree, there is a large degree of variation in both directions:  $E_c$  can both underestimate and overestimate  $E_e$  depending on the motif. Some of this discrepancy has been previously reported and is not unexpected: motifs that include physically favorable relationships such as salt bridges or disulfide bridges are

---

<sup>1</sup> Equation 5, which is used to compute a P value (the probability of any random match) from an expectation value (the average number of random matches) is based on a Poisson distribution, which assumes independence of counted objects. Repeat matches in the same structure are usually not independent and it would be best to not count them. Regardless, the difference is minor at low expectation values.



expected to occur with greater frequency than a model based solely on geometry would predict. There is an association with motif sizes—matches to the larger motifs are predicted to be more numerous than are actually observed ( $E_c > E_e$ ). This is due in part to the fact that  $E_c$  counts similar matches within the same protein, but also because GASPS restricts single atom-pair distance deviations between match and motif to reasonable distances (sidechain-sidechain deviation  $< 3.8 \text{ \AA}$ ;  $\alpha$ -carbon -  $\alpha$ -carbon deviation  $< 5.0 \text{ \AA}$ ). Relaxing both the sidechain and  $\alpha$ -carbon distance deviations to  $15 \text{ \AA}$  allows for many more matches at the same RMSD giving a closer correspondence between  $E_e$  and  $E_c$ , but we feel that deviations up to  $15 \text{ \AA}$  cannot represent true correspondences. It appears from this that GASPS with computed scores would favor smaller motifs as compared to GASPS with empirical scores because Equation 4 tends to over-predict the numbers of false positives for the larger motifs.

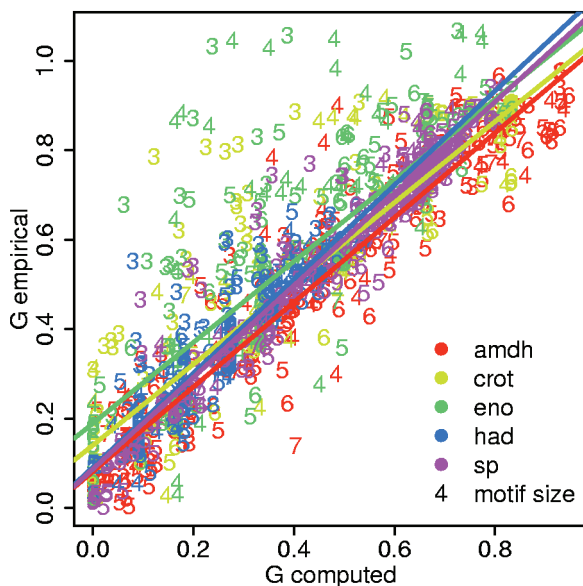


**Figure 2. Empirical counts of false positives versus computed expectation values.** Each point represents a single false positive match to a motif sampled by GASPS.  $E_{\text{empirical}}$  is the count of false positives that match with equal or better RMSD, and  $E_{\text{computed}}$  is calculated based on Equation 4. The solid line is the equivalence point where  $E_e = E_c$ . Points are shaded by the number of residues in the motif.

### **GASPS scores compared.**

The above discussion of expectation values focuses on the expectation values at values of RMSD from false positives alone. The computation of G scores, both by empirical ( $G_e$ ) and computational ( $G_c$ ) methods, depends on the expectation values of true positive or group member matches. While the correspondence should hold between expectations of true positives as it does for false positives, the relevant expectation values are at values that are too low to compare accurately by an analysis like Figure 2. We directly compared the  $G_e$  and  $G_c$  values for each of the motifs used in the above  $E_e$  and  $E_c$  comparisons. While we see the same overall correspondence, the main difference being that  $G_c$  has a range of 0-1.0 and  $G_e$  a range of 0-1.1, there are two noteworthy trends: First we see the expected trend for  $G_c$  to be an underestimate of  $G_e$  for the larger motifs. Second, there

appears to be more distant outliers above the line than there are below. These motifs are predicted to occur very often by  $G_c$  but in fact occur much less often than expected. These may represent motifs that are usually physically unfavorable, but are necessary for a unique characteristic of function. For the enolase and amidohydrolyase superfamilies, these outliers ( $G_e > 0.9$  AND  $G_c < 0.6$ ) are small motifs dominated by very close negatively charged residues that act as metal ligands.

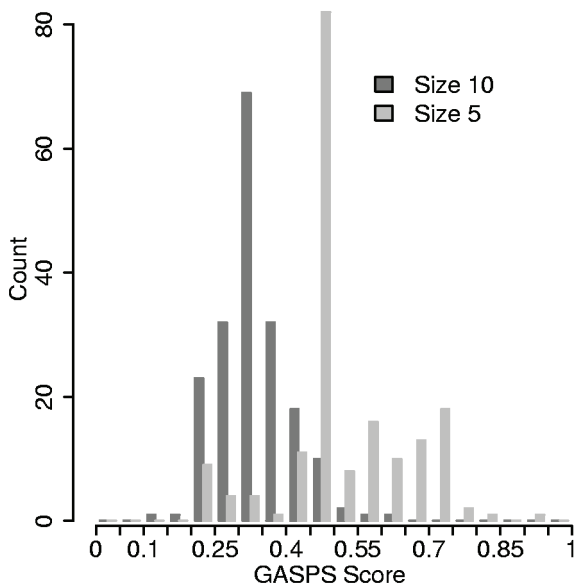


**Figure 3. GASPS scores (G) compared between empirical and computed methods.** Points are drawn as numbers which describe the number of residues in a motif. The color indicates the group: amidohydrolyase superfamily (amdh), crotonase superfamily (crot), enolase superfamily (eno), halo-acid dehalogenase (had) or serine proteases (sp). Colored lines are lines fit by linear least-squares regression on the data split by groups.

### GASPS with $G_c$ on Random Groups

The significance of a motif found by GASPS is measured by comparing it against the distribution of motifs generated on randomly constructed groups. For an accurate comparison I used the same randomly constructed groups as used previously for GASPS that used  $G_c$  (Chapter 2, Supplementary Materials). Here we see the first major difference

in results given by using  $G_c$  instead of  $G_e$ . GASPS with  $G_c$  generates motifs with higher scores for unrelated proteins (Figure 4). The biggest peak for the distributions of groups with 5 and 10 structures corresponds to matching two or three other structures in addition to the query structure. It appears that three unrelated structures chosen at random often share a motif by chance alone. For groups of 10 structures we can set a significance cutoff at approximately  $G_c=0.5$  and for groups of 5 structures the same threshold is as high as  $G_c=0.75$ .

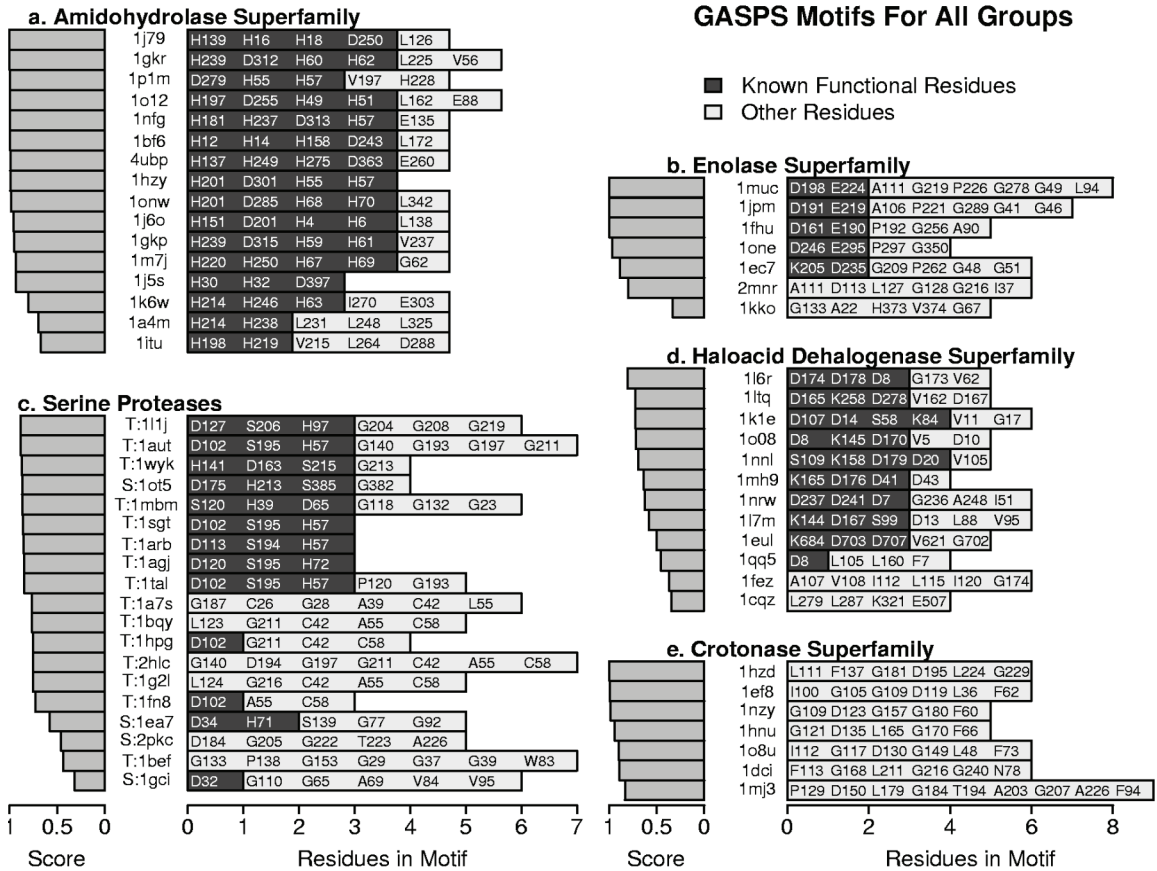


**Figure 4. Distributions of motifs by GASPS with  $G_c$  on random groups.**

### **Composition of motifs**

I next made use of  $G_c$  in an updated version of GASPS on the same set of superfamilies and structures as used in Chapter 2. While the new GASPS was much faster, both the distributions of scores and the makeup of motifs were similar (see Figure 5; compare with Chapter 2, Figure 3). I show only the result of a single run of GASPS for each query structure so differences at the level of individual motifs are not significant. Surprisingly

we do not see the expected trend for smaller motifs to be favored by  $G_c$ . Instead, we actually see larger motifs. This is likely the result of trends at very low expectation values, which cannot be adequately represented by the data in Figure 2.



**Figure 5. Composition of motifs generated by GASPS with computed G scores.** The results of a single GASPS run are presented for each named query structure. Residues in the motif that correspond to previously identified functional residues or known active-site motif residues are darkly shaded. All other residues are lightly shaded regardless of subsequent determination of their functional significance. For the serine proteases, query structures are labeled “T:” to denote trypsin-like folds or “S:” for subtilisin-like folds.

On 92 groups of enzymes defined by the Enzyme Commission (EC) enzyme naming scheme (International Union of Biochemistry and Molecular Biology. Nomenclature Committee. et al. 1992), I asked whether the degree of overlap with catalytic residues was different depending on the method of computing G scores. I used the Catalytic Site

Atlas (CSA) as the source of catalytic residues (Porter et al. 2004). Only the motifs that were considered significant by the G thresholds discussed above were counted. There was no significant difference in the number of overlaps with residues called catalytic by the CSA (see Table 1).

**Table 1. Overlap of significant motifs with catalytic sites in CSA.**

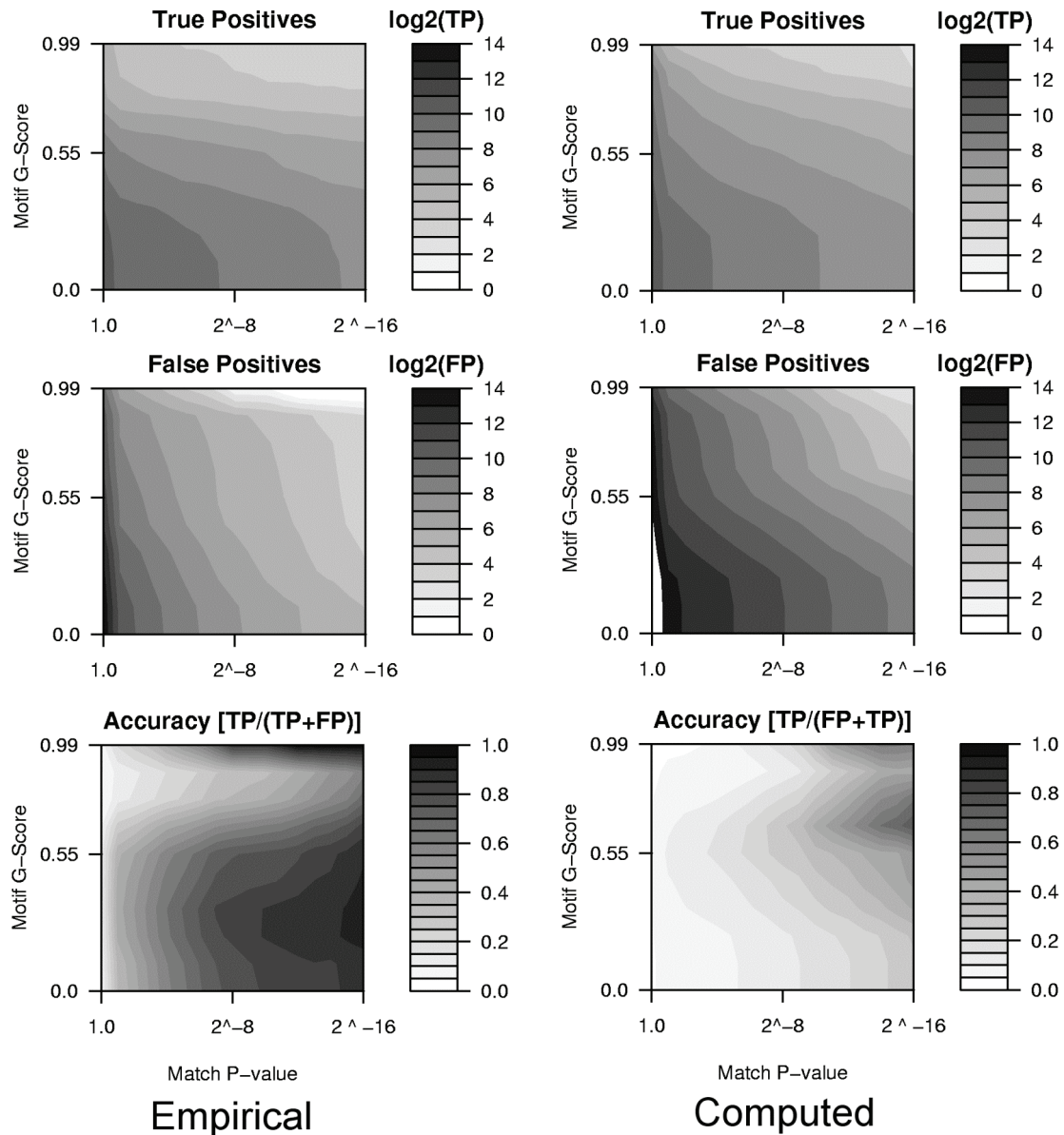
	CSA overlap	No overlap	% with CSA overlap
Empirical	88	59	60%
Computed	93	46	67%

Chi-squared = 1.237, df=1, p-value = 0.2660

### **Accuracy of motifs at identifying homologous groups**

The final test of using GASPS with  $G_c$  scores was to test the effectiveness of the generated motifs at identifying new protein structures to groups. I chose to use homologous groups, the families and superfamilies in the Structural Classification of Proteins (SCOP) version 1.65 (Murzin et al. 1995), because these groups produced higher scoring motifs than available classifications based on function. These motifs are tied to function and are useful at identifying function (see Chapter 4). Only groups with at least 7 structures, after removing redundancy at a threshold of 25% sequence identity, were included. One motif was generated for each structure in each group for both GASPS with  $G_c$  and GASPS with  $G_e$ . The domains newly added in SCOP version 1.67, compared to version 1.65, were searched against libraries of generated motifs with the program RIGOR (Kleywegt 1999). Only the first match to all motifs by each new domain were counted. Each match was given a P value using Equation 4 and 5, and both sets of motifs

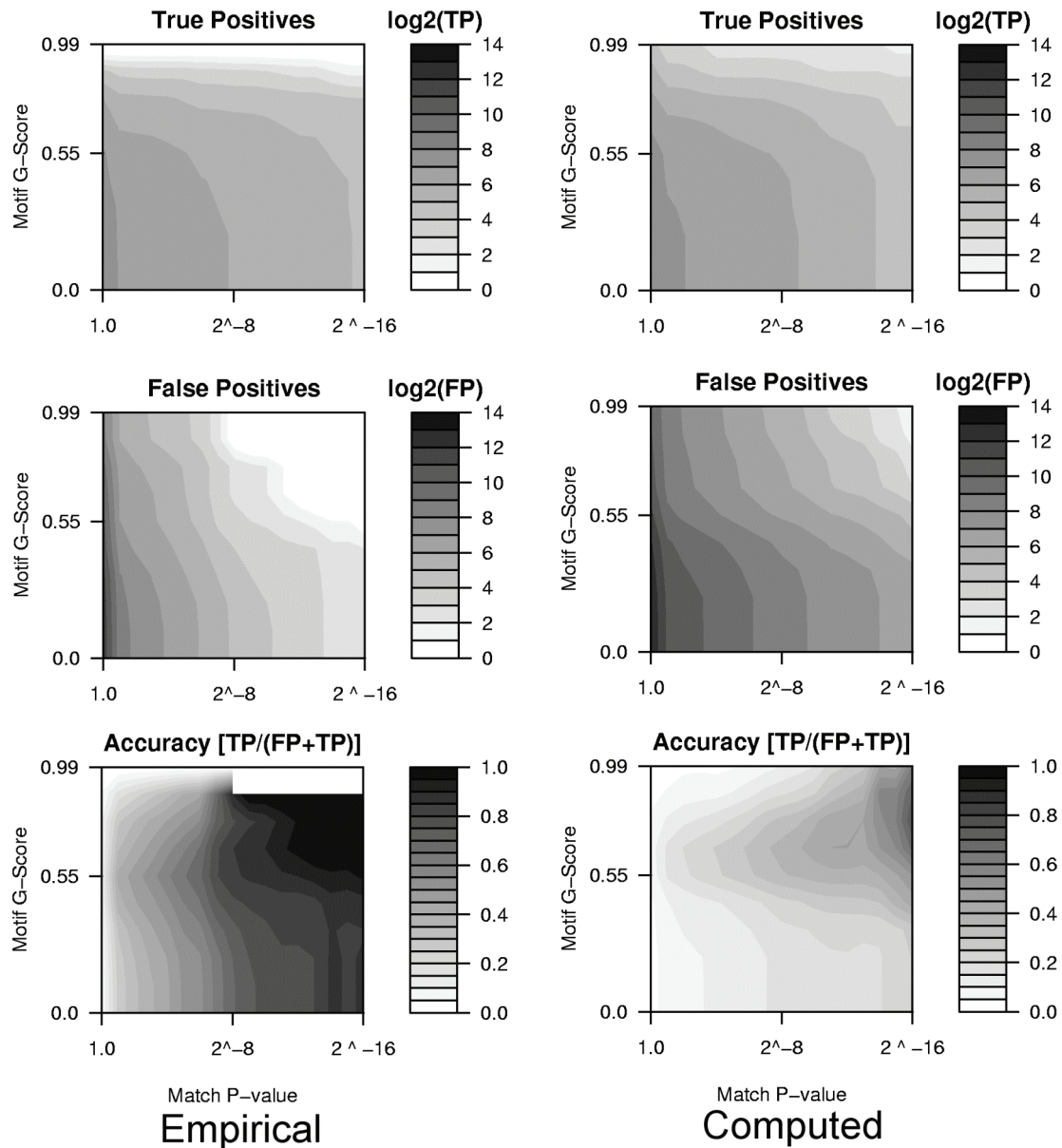
matched family or superfamily members at the same rates (true positives), but the motifs generated by computed G scores matched non-group members (false positives) at higher rates, significantly decreasing accuracy. Clearly, the motifs generated using computed G scores are not as specific to the group as those generated by empirical G scores. Several factors may account for this loss of specificity. Computed G scores allow for motifs that may identify a broader homologous class: fold instead of superfamily, or superfamily instead of fold, or they may simply consist of a very general protein motif such as disulfides or hydrophobic clusters. Most of the GASPS with  $G_c$  runs that resulted in a much higher scoring motif than the corresponding run with  $G_e$  included disulfides or hydrophobic clusters, especially leucines.



**Figure 6. SCOP superfamilies identified by motifs generated by empirical G scores compared to computed G scores.**

In all plots, the P value and G score axes are treated as thresholds. All matches with higher G scores or lower P values are counted and plotted at that location. The true positives and false positives are reported as the base 2 logarithms of the count. The x-axis is a logarithmic scale that decreases from left to right.





**Figure 7. SCOP families identified by motifs generated by empirical G scores compared to computed G scores.**  
See legend for Figure 6.

## ***Discussion***

I tested a new scoring function that could potentially make the running of GASPS much faster. It uses a geometrical model to predict the expected number of random matches to a motif to identify how significant are the true positive matches. Motifs generated with this scoring function appeared to be very similar to those generated with the empirical scoring function. Functional residues were identified with similar rates. As a tool to find a conserved structural pattern within a group, this faster scoring function appears adequate, and future studies should consider using the faster approach. I have discovered two significant shortcomings of this new scoring. First, motifs discovered by the new scoring function have to achieve a much higher score to stand out as significant. This problem is made greater by the fact that  $G_e$  scores tend to be greater than  $G_c$  scores overall (see Figure 3). Second, the specificity of motifs made using  $G_c$  have lower specificity. To make these motifs as useful as possible for annotation of new protein structures I need to eliminate as many false positives as possible. These problems with  $G_c$  were serious enough that I chose to not use it for further analysis here.

## ***References***

- International Union of Biochemistry and Molecular Biology. Nomenclature Committee. and E. C. Webb (1992). *Enzyme nomenclature 1992: recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology on the nomenclature and classification of enzymes*. San Diego, Published for the International Union of Biochemistry and Molecular Biology by Academic Press.
- Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." *J Mol Biol* **285**(4): 1887-97.
- Meng, E. C., B. J. Polacco, et al. (2004). "Superfamily active site templates." *Proteins* **55**(4): 962-76.

- Murzin, A. G., S. E. Brenner, et al. (1995). "SCOP: a structural classification of proteins database for the investigation of sequences and structures." *J Mol Biol* **247**(4): 536-40.
- Polacco, B. J. and P. C. Babbitt (2006). "Automated discovery of 3D motifs for protein function annotation." *Bioinformatics* **22**(6): 723-30.
- Porter, C. T., G. J. Bartlett, et al. (2004). "The Catalytic Site Atlas: a resource of catalytic sites and residues identified in enzymes using structural data." *Nucleic Acids Res* **32 Database issue**: D129-33.
- Stark, A., S. Sunyaev, et al. (2003). "A model for statistical significance of local similarities in structure." *J Mol Biol* **326**(5): 1307-16.

## Introduction to Chapter 4

All the work described in previous chapters sets the stage for a broad application of GASPS across as many proteins as possible. To understand the benefit of these generated motifs it is important to remember the two main scientific functions that motifs serve. First, motifs provide a way to classify proteins to groups, so such a broad study seeks to generate the maximum impact that GASPS can provide in this area. Beyond this, it also can provide an assessment of how broadly and on what groups the technique of 3D motifs can be effective. Second, motifs represent the most conserved elements in protein structures. The distribution of these motifs among protein groups can answer questions about how local protein structure has evolved. Without a shift in overall function, how often does the evolution of proteins maintain the same set of critical residues? As functions change, how often do proteins make use of existing functional components? Furthermore, the composition of the motifs provides information about what features of proteins tends to be the most conserved and therefore the most critical for maintaining the function of proteins as they evolve. The distribution of these motifs and features among different types of proteins and different types of classifications indicate trends in evolution.

This chapter forms the basis of a manuscript that will be submitted for publication to a peer-reviewed journal.

## Chapter 4: An exhaustive survey of 3D motifs.

### ***Abstract***

Knowledge of local protein structure, such as individual residues or clusters of interacting residues, is essential for understanding how protein structure delivers function, and especially how structure and function evolve together.

Here we examine the evolution of fine scale protein structure by determining the distribution of conserved 3D motifs, and what structural features tend to be conserved. We apply GASPS, which identifies the most conserved and unique motif from an input group of protein structures, to SCOP superfamilies and families as well as isofunctional groups defined by the Gene Ontology.

We find that homologous relationships are more important than functional relationships for the presence of a highly conserved motif. Non-homologous but isofunctional groups do not commonly share a motif. This suggests that protein functions, as they are commonly described, are usually accomplished by different means in unrelated proteins. About one third of all SCOP groups show a strongly conserved motif. The lack of a conserved motif in the remaining two thirds of groups reveals that evolution of new functions is commonly not constrained to maintain the positions of a critical set of residues.

We describe the patterns of structural elements and residue types among motifs to reveal trends in conservation of local structure. As expected, the motifs from all groups show a strong link to function, frequently overlapping with known catalytic, metal and other

ligand sites. Additionally, disulfides as well as stabilized charged residue pairs are overrepresented among the most conserved motifs. Residue distribution among the motifs is mostly as expected based on these common elements: cysteine, histidine, aspartate and glutamate are among the most frequent. More surprisingly, glycine is the most common motif residue and glycine in motifs has the greatest rate of non-metal ligand interaction among all other motif residues.

The motifs generated in this study are available via a web resource named GASPSdb, which is effective for annotating protein structures as well as highlighting important residues in new structures. Using these motifs, we show that 3D motifs offer promise for annotating low quality homology models built on distantly related templates.

## ***Introduction***

With the recognition that diverse and varied proteins can make use of the same overall fold, recent efforts in computational protein structural biology have shifted from examining large-scale features such as an entire sequence or fold, to a more focused examination of fine-scale features such as the orientations of a small number of sidechains. While the large-scale approach is useful for identifying homology, the fine scale approach allows for the identification of shared and distinct functional differences not apparent from the wider, large-scale view (Watson et al. 2007). One fine-scale approach, the use of 3D or structural motifs (sometimes called templates), relies on the cataloging of functionally related structural components comprised of the types and orientations of a small number of residues or their functional atoms (Fetrow et al. 1998; Laskowski et al. 2005; Torrance et al. 2005). Because these motifs normally contain the

elements that actually deliver function, finding these motifs in newly solved or modeled structures can imply a likely function for the protein and a hypothesis at where and how the function is performed.

The development and investigation of this fine-scaled 3D motif approach requires two main components. First, we need tools that can search for matches among protein structures. There are already many well developed available motif search methods (for example, Artymiuk et al. 1994; Fetrow et al. 1998; Barker et al. 2003). Second, we need to know what motifs are indicators of what function. In other words, we need knowledge about how the requirements of function constrain fine scale protein structure, and whether the constraints are strong enough so that the same fine-scale structure will be conserved over great evolutionary distances. While there are observed cases of 3D motifs being tied to a specific function over great evolutionary distances (Meng et al. 2004), no study has yet shown how common such motifs are among a broad cross section of the protein universe.

Here we use our previously described technique named GASPS (Polacco et al. 2006), to describe the patterns of motifs from all groups of proteins, defined by homology and/or function, restricted only by available structures. Such a broad collection of these fine scale motifs provides a resource to aid annotation of protein structures, answers how universally we might be able to apply these focused 3D motif methods, and provides us with a source of more basic biological knowledge. Each motif is a hypothesis about which residues are so important for the function of a protein as to be irreplaceable. These critical residues can provide, for example, a specific step in a catalytic mechanism (Gerlt et al. 2001), a very specific orientation of a binding partner, an important stabilization of

an active site, or a specific geometry of the peptide backbone (Dym et al. 2001), to name just a few. Having a broad collection will also allow us to answer what types of residues and structural features tend to be most rigidly conserved. Knowing this can inform future studies of the evolution of protein function within families.

By considering groups defined by homology and function, we recognize that the mechanism of protein function and other conserved details of protein structure can be products of both natural selection and evolutionary history. By examining groups defined by function alone, we are testing the hypothesis that natural selection to perform the same function on unrelated proteins can shape similar functional sites. For groups defined by homology alone, presence of a motif, especially when it is clearly related to functional sites, supports the hypothesis that an existing functional component can be recruited to perform new overall functions. In these cases, the motif cannot indicate a protein's overall function, but instead can indicate a conserved functional step. For example, this might be a single step in an enzyme's reaction pathway, or simply the binding of a metal ligand (Gerlt et al. 2001). Finally when we examine the most specific groups, those that are both homologous and isofunctional, we allow for the detection of motifs in cases where a new function evolved with a new set of critically important residues with no need for ancestral residues.

Ours is not the first publicly available collection of motifs (Stark et al. 2003; Torrance et al. 2005), but it does provide new coverage and a new emphasis on classification ability. Moreover, our technique is a protein-group driven approach that lets us find motifs, if they exist, for all protein groups. Other motifs have typically been chosen based on available prior knowledge of functionally important residues such as the catalytic triad of



the serine proteases. These motifs have been shown successful at identifying specific enzymatic activities (Torrance et al. 2005), binding relationships (Artymiuk et al. 1994), and superfamily membership (Meng et al. 2004). The catalytic site atlas (CSA) is a database of high quality catalytic residue designations gleaned from the literature (Porter et al. 2004). It currently provides 147 non-redundant active site motifs for enzymes (Torrance et al. 2005). In lieu of a literature search, a faster method, though more error-prone, is to use the information imbedded directly in protein coordinate files, such as SITE records, or proximity to ligands. This technique has been used by PINTS (Stark et al. 2003) and others (Artymiuk et al. 1994; Kleywegt 1999; Laskowski et al. 2005). Our approach, GASPS, is automated, unbiased and applicable to any group of proteins with sufficient structures (Polacco et al. 2006). It seeks to choose motifs with a high degree of classification ability, measured by a motif's tendency to match group members with high sensitivity and specificity. While it is not biased by accepted trends of functionally important residues, the motifs it finds have been shown to overlap with known functional sites.

We describe here the motifs generated by applying GASPS to isofunctional groups defined by Gene Ontology (GO) molecular function terms (Ashburner et al. 2000), homologous groups defined by the Structural Classification of Proteins (SCOP) superfamilies and families (Murzin et al. 1995), and by homologous isofunctional groups defined by both GO terms and SCOP superfamily. For purposes of protein annotation, we find that GASPS motifs can provide coverage of proteins unavailable in existing motif libraries. Moreover, we find that while the motifs, for the most part, can be related to

known functional sites, homology is more important than function for determining the presence of a high-scoring motif.

## **Methods**

### **GASPS**

GASPS (Genetic Algorithm Search for Patterns in Structures) takes as input a group of proteins, the positive group, and a background set of other proteins. It seeks to choose the coordinates of a set of residues from a single positive group member that is well matched by all other members of the positive group, and not matched by members of the background set. Further details are described in an earlier publication (Polacco et al. 2006). GASPS was run once for each member (chain or domain) of each protein group, generating as many motifs for each group as there are members.

### **Protein groups**

For each classification, no groups of proteins were allowed with fewer than seven structures when reduced to a non-redundant set based on a 40% sequence identity cutoff. Where possible without going below seven structures, those groups with sufficient structures and diversity were further reduced based on a 25% sequence identity cutoff. This generated two sets of groups, those that could be reduced to a 25% sequence identity cutoff and those that could not. Homologous groups were created by gathering all domains in SCOP (version 1.65) families and superfamilies. Isofunctional groups were defined by gathering all protein chains that share a single GO molecular function term, including the terms implied by the “is a” hierarchy of GO. It is desirable to limit GO

terms to those that are not so general as to make highly improbable any motif. To generate motifs for all suitable terms, and eliminate the obvious artificial groupings (such as all structures sharing GO term 5488, “Binding”), groups were discarded if they had greater than 50 non-redundant structures. Isofunctional homologous groups were generated by gathering all protein chains that shared at least one homologous domain as defined by SCOP superfamilies, then generating unique, but not mutually exclusive, groups defined by GO molecular function terms. Table 1 gives the counts of groups and generated motifs.

### Searching motif libraries with proteins

We use the program RIGOR (Kleywegt 1999) to search the libraries of GASPS motifs. RIGOR returns all matches between a protein and each motif that satisfy a superposition RMSD threshold and a maximum distance deviation threshold. Because smaller motifs match randomly with much greater frequency than large motifs, the RMSD threshold was set per motif based on a computation of the number of expected random matches (E-value). GASPS uses relative RMSD between false positive and true positive matches to the same motif to determine the quality of the motif. When we compare matches involving different motifs, the RMSD becomes less meaningful. Random matches to a motif comprising three leucines are much more likely than to a motif with five tryptophans. We computed an E-value based on a slight modification of the method of Stark et al. (Stark et al. 2003) that accounts for residue background frequency, number of residues, and distance between atoms in each residue:

$$E = a_0 \Phi a_3^N R_M^{2.93N-5.88} \prod_{!gly} \frac{R_M^2}{d_r^2}.$$

Here,  $\Phi$  is the product of residue frequencies as percentages,  $N$  is the number of residues, and  $R_M$  is the rmsd. All other variables are experimentally determined constants. The constant  $a_0$  accounts for the size of search space; we used  $a_0=1.57 \times 10^{10}$  for all searches. While using the same value for  $a_0$  regardless of database size can lead to inaccurate estimates of true expectation values, doing so generates an accurate score that reflects the strength of a pairwise match, allowing for direct comparison of matches between different searches. The right-most product is over each non-glycine residue in the motif, and corrects for the two-atom nature of GASPS and RIGOR non-glycine residues: an  $\alpha$ -carbon and a sidechain centroid. Individual factors in this product are ignored when the ratio  $R_m/d_r$ , where  $d_r$  is the average distance between side chain centroid and  $\alpha$ -carbon for a residue type  $r$ , is greater than 1. The value for  $a_3$  was taken directly from Stark et al. (2003) at 0.00179.

**Table 1. Group and motif counts by classification.**

Classification	Groups	Redundancy Filter	# Groups	# Motifs	Avg # Motifs per Group
<b><u>Gene Ontology</u></b>	<b>Molecular Function</b> (7 < n < 50 )		<b><u>272</u></b>	<b><u>4385</u></b>	
		25 PID	177	2593	14.6
		40 PID	95	1792	18.9
<b><u>SCOP</u></b>	<b>Superfamilies</b>		<b><u>323</u></b>	<b><u>3599</u></b>	
			<b><u>186</u></b>	<b><u>2259</u></b>	
		25 PID	131	1801	13.7
	40 PID	55	458	8.3	
	<b>Families</b>		<b><u>137</u></b>	<b><u>1340</u></b>	
		25 PID	64	670	10.5
	40 PID	73	670	9.2	
<b><u>GO and SCOP</u></b>	<b>Superfamilies</b>		<b><u>376</u></b>	<b><u>4581</u></b>	
		25 PID	231	3318	14.4
		40 PID	145	1263	8.7

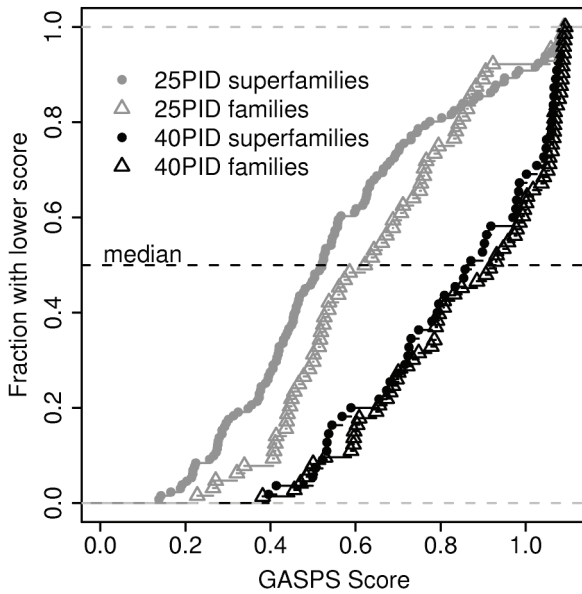
## **Results**

### **Quality of Motifs**

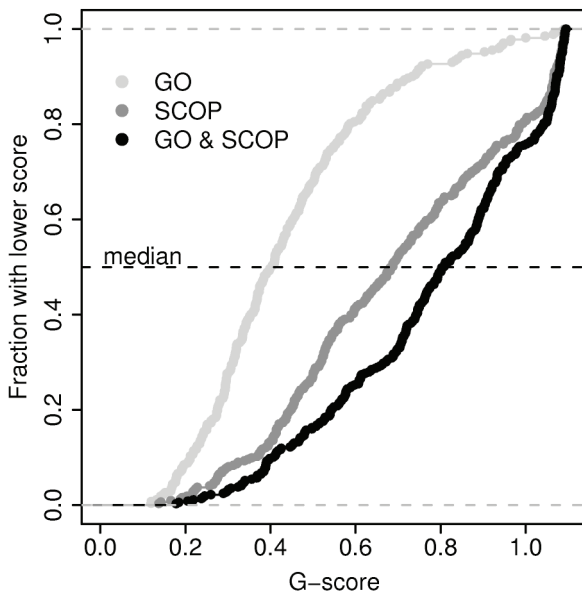
One goal of the current work is to study the broad applicability of 3D motifs across different types of proteins and different types of classifications. A motif's quality is described by its G score, short for GASPS score, which indicates the motif's ability to identify all other group members with high specificity. Ranging from 0 to 1.1, its main component is the area under a shortened receiver-operator characteristic (ROC) plot. The other component is the relative separation between true and false positive matches and accounts for only 0 to 0.1 of the total G-score. Therefore, most G-scores above 1.0 imply perfect separation in an ROC plot (ROC area = 1.0) though any score above 0.7 is highly significant and scores below about 0.4 are highly suspect. G-scores can give us a sense of how well groups of proteins can be identified by motifs. A high scoring motif for a group is evidence of a unique evolutionary constraint on that group. For homologous groups, the G-scores tell us whether there is a single structural pattern that evolution has not been able to alter, presumably without a loss of protein function. For non-homologous, functionally similar groups, high G-scores indicate convergent evolution where independent inventions of the same function required a common pattern of residues.

Previous work with GASPS showed it to be very effective on a small number of well-studied superfamilies. In this study we included all SCOP-defined superfamilies and families with sufficient structures (see Methods) as well as groups defined by common Gene Ontology (GO) annotations. Figure 1 and Figure 2 show the distribution of the highest scoring motif for each group in SCOP, GO and the GO/SCOP groupings. While a

large number of SCOP families and superfamilies have very high G-scores, the majority of protein groups produce motifs with G-scores lower than for the previously well-studied superfamilies (top G-score for haloacid dehalogenase superfamily is 0.8, (Polacco et al. 2006). However, the evolutionary distance between members in a group is important. Those groups composed of members that all share less than 25% sequence identity have significantly lower scores (median G-score = 0.54) than those where we permitted up to 40% sequence identity (median G-score = 0.90). Evolutionary distance measured by sequence identity is even more important than is the evolutionary distance implied by SCOP hierarchy depth: the superfamily and family distributions are much more similar than are the distributions grouped by sequence identity. Because group size is correlated with G-scores, and the groups where we allowed 40% sequence identity were smaller, this could also be a result of group size effects. However, a linear model fit to this data to predict G-scores with parameters for group size, SCOP hierarchy and percent identity, shows the greatest effect to come from percent identity, the second greatest from group size, and finally the effect from SCOP hierarchy depth is nearly insignificant.



**Figure 1. Distribution of motif G-scores on SCOP groups.**



**Figure 2. Distribution of motif G-scores on Gene Ontology and SCOP groups.**

For groups defined by GO molecular function annotations, the G-scores were even lower than for homologous groups. This is the result of GO defined groups containing unrelated proteins that perform the same molecular function, but doing so in a different way with a

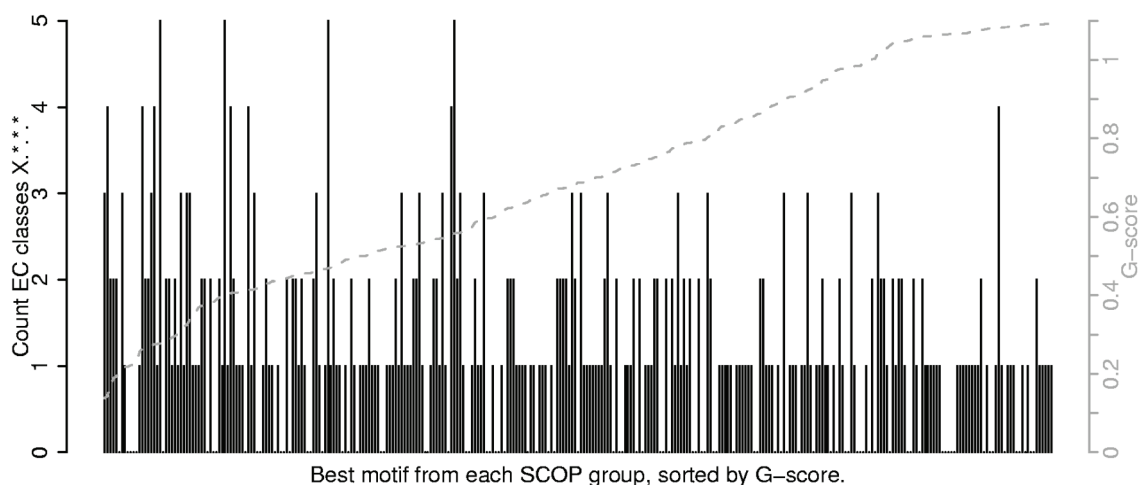
different set of important residues. Even when GO groups are matched by a high-scoring motif, the groups are made up of, or at least dominated by, a single homologous group that accounts for the high G-score (not shown). The frequency with which evolution has invented new ways of performing the same function indicates that most protein functions can be performed in many different ways.

If we further subdivide the SCOP superfamilies by GO molecular function annotations, we see an improvement in G-scores over both SCOP superfamilies and GO groups. This shift to higher scores is not simply the result of chance due to the use of dividing the superfamilies into smaller groups. In fact, the groups in this classification are larger on average owing to the groups not being mutually exclusive. For each superfamily, multiple overlapping groups can be defined depending on the precision of the GO annotations assigned to the individual structures. For example, SCOP superfamily “FAD/NAD(P)-binding domain” (c.3.1) has a group for GO term 16491 “oxidoreductase activity” as well as a group for GO term 15036 “disulfide oxidoreductase activity” that is a subset of the other. This actually results in a larger average group size compared to SCOP superfamilies (see Table 1). The larger groups are over-represented compared to the smaller groups because there are more ways to divide them.

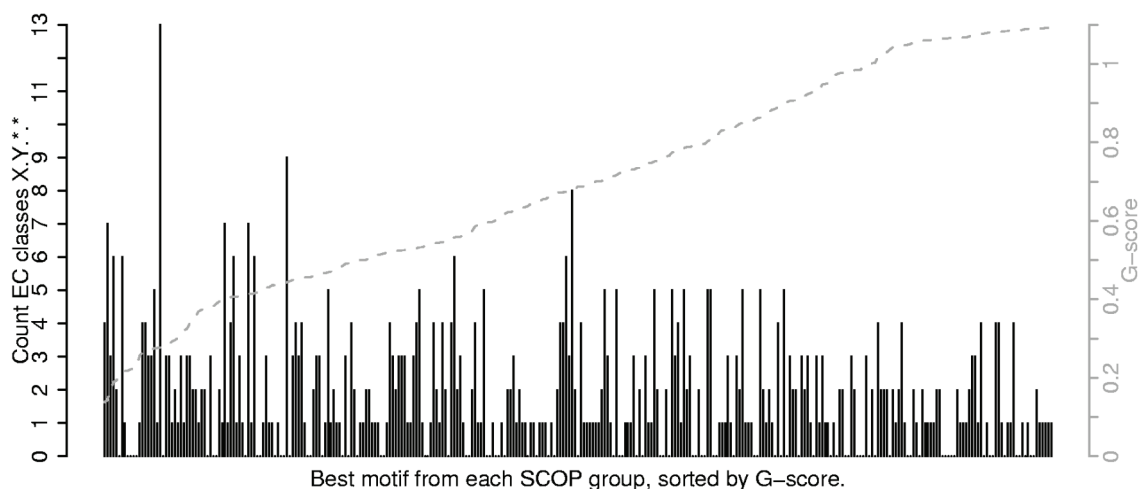
The results on the GO/SCOP groups indicate that groups with more functional diversity are less likely to have a conserved motif. To investigate this finding in more detail, we counted the number of distinct enzyme commission (EC) (International Union of Biochemistry and Molecular Biology. Nomenclature Committee. et al. 1992) classes at the first and second positions for each superfamily and family in SCOP. This measure of functional diversity is an underestimate because it only counts the enzyme functions and



only those given an EC number. Figure 3 and Figure 4 show the expected trend that the most functionally diverse are more likely to have a lower G-score, but there are still many groups with highly significant G-scores and significant functional diversity. Among these groups are the well known enolase (Babbitt et al. 1996), haloacid dehalogenase (Allen et al. 2004), and amidohydrolase superfamilies (Holm et al. 1997). The remainder are also good candidates for superfamilies that have evolved according to a similar evolutionary model.



**Figure 3. Number of distinct EC classes at first position in each SCOP group.** Each vertical bar shows the count of distinct EC codes, only counting the first position, in a SCOP group (superfamilies and families). The SCOP groups are sorted along the x axis by the G-score of their best motif. G-scores are indicated by the dashed gray line.



**Figure 4. Number of distinct EC classes at first two positions in each SCOP group.** Each vertical bar shows the count of distinct EC codes, only counting the first two positions, in a SCOP group (superfamilies and families). Remaining is as in Figure 3.

### Patterns of conservation in 3D

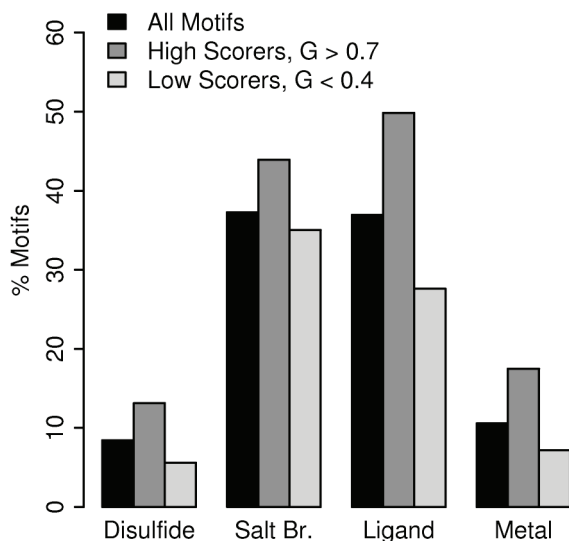
The motifs generated by GASPS are chosen for both their conservation in 3D space and their uniqueness, or lack of matches among unrelated protein structures. It appears from

this study, that repeated patterns between isofunctional but unrelated proteins are rare. Instead, most well conserved 3D structural patterns of more than two residues are the results of homology and are unlikely to be repeated in a non-homologous group (two residue motifs such as disulfide bridges are frequently conserved but not unique). While the results of chapter 3 indicate the opposite, that some sensitive motifs are not specific to a group, these motifs can be made more specific by the addition of one or two local residues without significantly changing the composition of the motif. Uniqueness therefore plays less of a role, so that GASPS motifs across broad protein groups describe primarily the patterns of conservation in 3D space.

It is well recognized that functionally significant residues are well conserved in both sequence and structure. It follows then that we can expect a large number of motifs generated by GASPS to contain residues that are known to be functional. Using the Catalytic Site Atlas (CSA) as an independent source of functional residue information, we do see significant overlap with GASPS motifs and CSA entries. In fact, 63% of protein groups in this study with representatives in the CSA have a motif with one or more residues directly involved in catalysis. As the name implies, the CSA limits its scope to enzymes and uses a strict definition of residues directly involved in catalysis, so that many residues involved in stabilizing or binding in a functional site are not included. Important binding sites from non-enzymes, such as the iron binding site in the ferritin superfamily, or the heme binding sites in globins, as well as binding sites for metals or other cofactors involved in catalysis (such as the metal binding sites in the enolase and amidohydrolase superfamilies) occur with high frequency among the GASPS motifs (Figure 5). We identified motif residues that interact with ligands by identifying residue

atoms that are within 4 Å of an atom described by a 'HETATM' record. Nearly half of the highest scoring motifs are associated with a ligand. The number is lower for the lower scoring motifs, confirming that G-scores correlate with functional significance, and that the association with ligands is not due to random sampling of protein residues.

Approximately 1/4 of all motif ligand interactions are to metal ligands. This number goes as high as 1/3 for the highest scoring motifs revealing that metal binding sites are among the best conserved structural features, and the most reliable to match by structural motifs. It is important to keep in mind that these computed frequencies of ligand interactions by GASPS motifs necessarily underestimate the actual number of motif residues that interact with a ligand biologically. The structures may not have been solved in the presence of a ligand (only two thirds of PDB files used in this study included any HETATM record), and if present, the ligands may have been unresolved. Furthermore, any static description of structure cannot fully represent the dynamic range of biologically relevant protein motions.



**Figure 5. Residue interactions captured by motifs.**

Another large group of features found in GASPS motifs are stabilizing residue interactions such as salt bridges or disulfide bridges. Salt bridges were identified by finding acidic (glutamate or aspartate) and basic (histidine, lysine, or arginine) sidechain atoms (O and N) within 4.0 Å. Disulfide bridges were identified by finding cysteine S atoms within 3 Å of each other. While a single such pair of residues is not unique to any group, when paired up with other neighboring conserved residues or other pairs of stabilizing residues they often become useful identifiers. Only 8% of all motifs include at least one cysteine involved in a disulfide bridge, however, the presence of a disulfide bridge is highly correlated with G-score. The motifs with the highest G-scores are twice as likely to contain at least one disulfide partner (13%) as those with the lowest G-scores (6%). Almost as common as the previously discussed ligand interactions are stabilized charged residues or salt bridges. These are also correlated with G-scores, though to a lesser degree.

## Residue types in motifs

Other studies with 3D motifs often focus on discovering functional sites, and therefore limit their analysis to those polar and charged residues assumed to be most likely functional. GASPS makes no such distinctions, so that any residue can be included in a motif provided it is conserved in 3D space in relation to other conserved residues. Our approach allows us to ask which residues dominate the motifs, tend to be the most conserved, and provide the most classification information. Indirectly this can tell us how critical to protein function is each residue's unique role. Residue prevalence in motifs was normalized by group size and motif size, so that larger groups or larger motifs would not bias the results. We also split the motifs by G-score to observe trends as classification ability increases. In detail, for each G-score range and residue, the normalized residue frequency for any residue type ( $f_t$ ) was calculated as:

$$f_t = \frac{1}{n_g} \sum_{groups} \frac{1}{n_m} \sum_{motifs} \frac{n_t}{n_r}$$

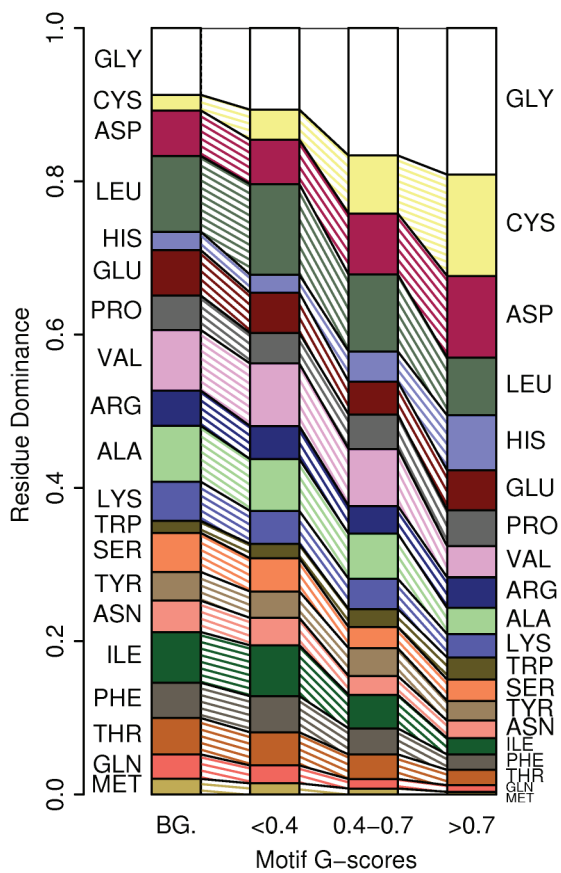
where  $n_t$  is the number of residues of type  $t$  in a motif,  $n_r$  is the number of total residues in a motif,  $n_m$  is the total number of motifs in a group, and  $n_g$  is the total number of groups.

The distributions of dominant amino acids in the highest scoring motifs shows some patterns that are expected from the previous discussion of common structural features, but also indicate that we have not yet described all important structural trends captured by motifs (Figure 6). The presence of cysteine, aspartate, histidine and glutamate among the top seven amino acids are expected from the previous discussion of catalytic sites, ligand—especially metal—sites, and salt bridges. Less expected is the dominance of glycine, and prominent role of leucine. The prevalence of cysteines is not surprising

given their unique role in disulfide bridges, catalytic sites and metal binding sites, all of which are well represented among motifs. Likewise, histidine and aspartate frequently play a role in metal binding sites and catalytic sites. The prevalence of leucine among motifs, on the other hand, may simply be a result of the high overall frequency of leucine in the entire proteins. The frequency of a residue type among the motifs is compared in Figure 6 to its frequency among the entire set of residues allowed by GASPS, the background frequency. The frequency of leucine among the high scoring motifs is actually reduced compared to this background frequency. Glycine also has a high background frequency, second only to leucine, but its prevalence among the high-scoring motifs is increased over this background. What accounts for glycine's high prevalence? While it has no sidechain to interact with ligands, a relatively high proportion of the glycines in motifs are within interaction distances of ligands. In fact, glycines in motifs rank fourth behind only histidine, cysteine and aspartate for their rate of interaction with ligands. Most ligand-interacting glycines interact with phosphate containing compounds, the most common being FAD, NAD and ADP. A smaller number of glycines interact with sulfur containing compounds, mostly sulfate. Most inter-atom (non hydrogen) distances for interactions with glycines are between 3 and 4 Å (median=3.25 Å), so that many are outside of the range of a typical hydrogen bond. Instead, the presence of glycine in binding pockets may provide for the tight bending of loops around ligands as well as space for a ligand to bind (Jornvall et al. 1984; Dym et al. 2001). Similarly, proline's unique backbone angles may account for its seventh highest frequency among high-scoring motifs. While the majority of prolines as well as many glycines are not within interaction distance of ligands, they are often near residues that do interact. Still,

not all of these conserved glycines and prolines are near known functional sites in motifs, and likely serve to stabilize an overall fold rather than the fine-scale geometries in a catalytic or binding site. Examples of these glycines can be found in motifs from the enolase superfamily (SCOP c.1.11; 2mnr, Gly291), and amylase families (SCOP c.1.8.1; 1esw, chain a, Gly40).





**Figure 6. Dominance of residue types, compared against background residue frequency, and at different G-scores.**

See text for computations of residue dominance that is shown here calculated for motifs pooled to three different groups by G-score. The column labeled BG. for background refers to the frequency of the residue type among all residues considered by GASPS for inclusion in any motif.

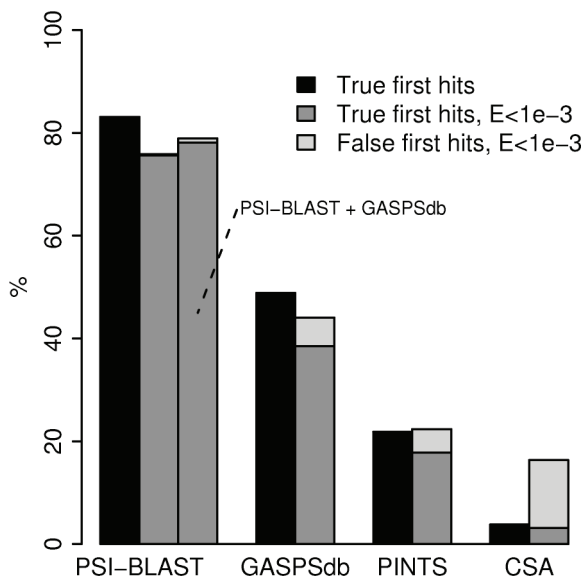
### Annotation of protein structures

Beyond a survey of structural conservation among diverse protein groups, the outcome of this study provides for a set of motifs that can be used to help annotate novel protein structures. We have packaged the motifs generated in this study together with structure matching and browsing tools as a web resource named GASPSdb (<http://gaspsdb.rbvi.ucsf.edu>). We demonstrate the benefit of GASPSdb by showing that its search results provide additional coverage with similar or better accuracy compared to

other available libraries of 3D motifs, CSA and PINTS described earlier. Between SCOP version 1.65 (used by GASPSdb) and SCOP version 1.69, an additional 1612 domains were added that were less than 40% identical to each other or earlier domains. GASPSdb contains superfamily motifs for about half of these new domains (790 domains, 49%). For the comparison, we used these 790 new domains as queries and for each search, we considered only the most significant match. Figure 7 shows how often these first matches identified a motif from the correct superfamily, and how often the first matches with significant scores ( $E < 0.001$ ) were true. No difference was found in the rate of matching structures labeled as “putative” in their “Structure Title” or “Structure Description” fields, suggesting that motifs perform as well on structures for which we have little prior knowledge.

By identifying functional sites and functional residues, the fine-scale information of 3D motifs from GASPSdb provides functional details that homology detection methods that use an entire protein such as PSI-BLAST(Altschul et al. 1997) cannot provide. On the other hand, PSI-BLAST is able to give accurate homology predictions for more proteins than 3D motifs generated by any method. Working together though, GASPSdb can extend the annotation power of PSI-BLAST by corroborating low-significance PSI-BLAST hits. On the same set of 790 domains from superfamilies with motifs in GASPSdb, 83% of PSI-BLAST searches return a sequence with the correct superfamily as the most significant match. However, the effectiveness of 3D motifs is mostly independent from sequence similarity, so that 3D motifs can complement homology searches. Considering only those proteins above where PSI-BLAST yields an ambiguous result (first match  $E > 1e-3$ ), only 57 of 184 PSI-BLAST first matches yield a true hit. A

GASPSdb RIGOR search can corroborate the true match 19 times (33%), not significantly different (chi-squared test,  $p = 0.46$ ) from its performance on the larger set (see Figure 7), and with only 4 false positives at  $E < 0.001$ .



**Figure 7. Coverage of GASPSdb compared to other 3D motif libraries and PSI-BLAST.**

The percentage on the y-axis is the number of structures giving true or false positives as their first hit when searched against the sequence or motif database on the x-axis.

### Homology models

A tool to annotate protein structures would be most useful if it worked on low quality structures, such as homology models generated from distant templates (e.g., less than 30% sequence identical), as well as structures solved empirically. We tested the performance of GASPSdb together with the homology-modeling tool MODELLER (Sali et al. 1993) used to model structures for the 790 new domains in SCOP 1.69 discussed previously. We generated models using only the new domain's sequence and an existing structure from the same superfamily in SCOP 1.65 as a template. To simulate the conditions of low-quality models, we required template structures to match the sequence

with a PSI-BLAST E-value worse than  $10^{-5}$ . The sequence-template pairs were automatically aligned and modeled using MODELLER with default settings and its built-in align2d program. Of the 618 resulting model-able pairs, 88 (14%) of the homology models were correctly annotated by GASPSdb motifs, while only 277 (45%) of the actual crystal structures were correctly annotated by any GASPSdb motif. These homology models are of sufficient quality to match a motif 32% (88/277) as often as their actual crystal structure. In another experiment, by purposely choosing a falsely homologous decoy template at similar PSI-BLAST E-values, we find that the rate of false positive motif matches to the decoy superfamily is only 0.005. Accuracy of homology models is known to be highly dependent on alignment accuracy (Martin et al. 1997; Venclovas et al. 2005), and this experiment included the simplest alignment protocol, with many expected alignment errors. Nevertheless, these low quality homology models are accurate enough to match the appropriate motif 32% (88 of 277) as often as their crystal structure. When applied to the large number of unknown sequences in available databases, even this relatively low proportion could prove useful.

The success of GASPS motifs on homology models is not due simply to the overall accuracy of homology modeling, but the tendency for GASPS to choose motifs that are accurately modeled. Among a set of homology models made for structures in GASPSdb, approximately 80% of the randomly generated motifs (those chosen before the first round of GASPS optimization) match with less significant E-values than the final GASPS-optimized motifs.

## ***Discussion***

Chosen only for their sensitivity and specificity, yet with frequent overlaps with functional sites, motifs presented in the GASPSdb make a useful tool to describe function and highlight likely functional residues of novel protein structures. These motifs are chosen for their ability to identify a group of proteins. Any protein that matches the motif is expected to share the same function, at least to the extent that function is shared among the original group that produced the motif. An alternative approach that does not require a mapping between motifs and specific functions is to find any significant similarity in local 3D structure between two proteins (Oldfield 2002; Laskowski et al. 2005). One challenge of this technique is the high number of false positives produced by these randomly chosen motifs. The method of Laskowski et al. (Laskowski et al. 2005), to further filter matches based on similarity of residues in the local structure of the motif source and its match, can effectively weed out the false positives and identify true homology between protein structures. While generally useful, these pairwise relationships do not identify to what extent the functions of two proteins is similar. The motifs generated here, like other motifs designed to identify specific groups of proteins, such as EC numbers (Torrance et al. 2005), can provide a useful complement to this pairwise motif technique.

In addition to this immediate practical application, taken together, the motifs generated here provide a view into the trends in evolutionary constraints on function. As discussed earlier, the high-scoring motifs provide evidence of evolutionary constraints. It is not surprising, then, that the high scoring motifs are mostly restricted to homologous groups. While cases of convergent evolution exist that can be described by a structural motif,

these cases are rare, and any common motif shared by two convergent groups may be washed out by additional isofunctional but independently evolved groups that do not share the motif. It is notable, though, that while the literature provides cases of diverse superfamilies and families matching a single motif, most groups with sufficient evolutionary distance do not share a single motif as defined here. This does not necessarily imply that the 3D motif approach to annotating function and identifying functional residues has limited application. While a single motif may not exist, subgroups within the larger group can share a common motif. Instead, it has more bearing on how function constrains protein structure. When function remains the same despite the lack of a shared motif, the diverse proteins must discover a new way of completing the function, or at least make do with a different set of residues. When these proteins evolve to perform new functions they do so without maintaining a superfamily-conserved set of residues. It is worth noting that the current study requires identical residue types. In many cases, while a recognized active site architecture is maintained, the roles such as acidic or basic side-chain, can be adequately performed by multiple residue types.

We have described common trends for evolutionary constraints on residue type and simple residue interactions. One interpretation is that the strongest constraints are on those residue types that perform unique roles: cysteines provide disulfides, histidines provide labile acid-base chemistry, glycines and prolines provide for unique backbone angles, and glycines' missing sidechains maximize available space. On the other hand, 3D motif methods require residues to be relatively unmoved across various proteins *and* the various experimental conditions used to solve the structures. Therefore, GASPS should favor features that confine sidechains to a specific position. As examples, metal

ions bind their ligands very tightly. Together with their functional importance, this explains why they were common features among motifs. This trend is observed at the residue level as well. The residues that frequently coordinate metal ions, cysteine, histidine, glutamate and aspartate, are among the most frequent. While this effect is not necessarily restricted to metal ligands, they seem to have the strongest effect. Both lysine and arginine in the motifs bind ligands (though not metals) with similar or greater rates as the negative charged and metal-binding aspartate and glutamate. However, the positive-charged residues rank several places below both negatively charged residues in their overall prevalence in motifs.

While this study provides one view on constraints on protein evolution, it is a view that is limited by the mechanics of GASPS. GASPS identifies just the strongest motif or constraint on each run. While repeated runs may reveal less well-conserved motifs, these secondary motifs can easily be missed, so that GASPS often cannot reveal all constraints on a protein group, but only the strongest. Similarly, GASPS looks for motifs that match all proteins in a group, and will give lower scores to those motifs that match only a fraction of the structures even at very high significance. Some of the moderately low G-scores seen here may simply represent a very well conserved motif, but only among members of a subgroup of the larger group. However, the trends for residue distributions of lower-scoring motifs to more closely match background distributions implies that many of the low scoring motifs are at least partially due to simple chance. As mentioned previously, GASPS requires identical residues at each position. It does not even consider conservative substitutions such as aspartate/glutamate. We previously examined modifying GASPS to allow for position-specific substitutions, but found little

improvement when applied to groups with already moderately high-scoring motifs.

Whether this substitution scheme allows for motifs with improved G-scores for the low scoring groups here is an open question.

GASPS requires sufficient diversity in available protein structures to weed out motifs that exist by chance or recent shared ancestry instead of by their importance to a shared function. As more structures are made available, the range of GASPS' effectiveness will extend to more groups. An automated method like GASPS can easily keep abreast of the latest developments.

## **References**

- Allen, K. N. and D. Dunaway-Mariano (2004). "Phosphoryl group transfer: evolution of a catalytic scaffold." Trends Biochem Sci **29**(9): 495-503.
- Altschul, S. F., T. L. Madden, et al. (1997). "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." Nucleic Acids Res **25**(17): 3389-402.
- Artymiuk, P. J., A. R. Poirrette, et al. (1994). "A graph-theoretic approach to the identification of three-dimensional patterns of amino acid side-chains in protein structures." J Mol Biol **243**(2): 327-44.
- Ashburner, M., C. A. Ball, et al. (2000). "Gene ontology: tool for the unification of biology. The Gene Ontology Consortium." Nat Genet **25**(1): 25-9.
- Babbitt, P. C., M. S. Hasson, et al. (1996). "The enolase superfamily: a general strategy for enzyme-catalyzed abstraction of the alpha-protons of carboxylic acids." Biochemistry **35**(51): 16489-501.
- Barker, J. A. and J. M. Thornton (2003). "An algorithm for constraint-based structural template matching: application to 3D templates with statistical analysis." Bioinformatics **19**(13): 1644-9.
- Dym, O. and D. Eisenberg (2001). "Sequence-structure analysis of FAD-containing proteins." Protein Sci **10**(9): 1712-28.
- Fetrow, J. S. and J. Skolnick (1998). "Method for prediction of protein function from sequence using the sequence-to-structure-to-function paradigm with application to glutaredoxins/thioredoxins and T1 ribonucleases." J Mol Biol **281**(5): 949-68.



- Gerlt, J. A. and P. C. Babbitt (2001). "Divergent evolution of enzymatic function: mechanistically diverse superfamilies and functionally distinct suprafamilies." Annu Rev Biochem **70**: 209-46.
- Holm, L. and C. Sander (1997). "An evolutionary treasure: unification of a broad set of amidohydrolases related to urease." Proteins **28**(1): 72-82.
- International Union of Biochemistry and Molecular Biology. Nomenclature Committee. and E. C. Webb (1992). Enzyme nomenclature 1992: recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology on the nomenclature and classification of enzymes. San Diego, Published for the International Union of Biochemistry and Molecular Biology by Academic Press.
- Jornvall, H., H. von Bahr-Lindstrom, et al. (1984). "Extensive variations and basic features in the alcohol dehydrogenase-sorbitol dehydrogenase family." Eur J Biochem **140**(1): 17-23.
- Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." J Mol Biol **285**(4): 1887-97.
- Laskowski, R. A., J. D. Watson, et al. (2005). "Protein function prediction using local 3D templates." J Mol Biol **351**(3): 614-26.
- Martin, A. C., M. W. MacArthur, et al. (1997). "Assessment of comparative modeling in CASP2." Proteins Suppl 1: 14-28.
- Meng, E. C., B. J. Polacco, et al. (2004). "Superfamily active site templates." Proteins **55**(4): 962-76.
- Murzin, A. G., S. E. Brenner, et al. (1995). "SCOP: a structural classification of proteins database for the investigation of sequences and structures." J Mol Biol **247**(4): 536-40.
- Oldfield, T. J. (2002). "Data mining the protein data bank: residue interactions." Proteins **49**(4): 510-28.
- Polacco, B. J. and P. C. Babbitt (2006). "Automated discovery of 3D motifs for protein function annotation." Bioinformatics **22**(6): 723-30.
- Porter, C. T., G. J. Bartlett, et al. (2004). "The Catalytic Site Atlas: a resource of catalytic sites and residues identified in enzymes using structural data." Nucleic Acids Res **32 Database issue**: D129-33.
- Sali, A. and T. L. Blundell (1993). "Comparative protein modelling by satisfaction of spatial restraints." J Mol Biol **234**(3): 779-815.
- Stark, A. and R. B. Russell (2003). "Annotation in three dimensions. PINTS: Patterns in Non-homologous Tertiary Structures." Nucleic Acids Res **31**(13): 3341-4.
- Stark, A., S. Sunyaev, et al. (2003). "A model for statistical significance of local similarities in structure." J Mol Biol **326**(5): 1307-16.

- Torrance, J. W., G. J. Bartlett, et al. (2005). "Using a Library of Structural Templates to Recognise Catalytic Sites and Explore their Evolution in Homologous Families." J Mol Biol **347**(3): 565-81.
- Venclovas, C. and M. Margelevicius (2005). "Comparative modeling in CASP6 using consensus approach to template selection, sequence-structure alignment, and structure assessment." Proteins **61 Suppl 7**: 99-105.
- Watson, J. D., S. Sanderson, et al. (2007). "Towards fully automated structure-based function prediction in structural genomics: a case study." J Mol Biol **367**(5): 1511-22.

## Conclusion

This work has focused on signature 3D motifs. These motifs are conserved within a group of proteins and can be used to identify the group. Through the development and application of GASPS, I have identified a large number of these signature 3D motifs that represent conserved functional components within protein structures. I have shown that these 3D motifs can be used to help annotate the functions of proteins but also that signature motifs are not a universally useful tool that can be used for all groups of proteins we may wish to classify. Many groups examined here provided no useful motif. In many cases, this can be attributed to the classification system used to define any particular group. While the group defined by the classification may not all share a motif, a natural sub-group may have. It may be that no perfect classification exists. The accepted classification of biological entities at any level is constantly changing. As long as there are new genomes to sequence, there will likely always be newly discovered protein structures that do not faithfully match any previously discovered group motif. The distribution of these 3D motifs and the patterns of residues within them also describe how local protein structure evolves. While we see a strong relationship between conserved elements and a protein's function, a protein's evolutionary history appears more important for determining local structure than just function alone. We see very few cases of convergent evolution here, where a single function has required the same set of residues in unrelated proteins. More often, we see that a 3D motif and the elements of a protein's function that the motif represents have been adopted from an ancestral function, even when the overall function has evolved to be different. This phenomenon is identified

here by homologous groups of proteins with diverse functions and well-conserved 3D motifs that are directly involved in a protein's function.

As more structures become available, I expect the automated method of GASPS can be used to generate motifs for protein groups that did not have enough structures to be included in this study. The motifs available in GASPSdb can therefore be a constantly growing resource available to the biological community. While I present some work here evaluating and discarding alternative techniques for GASPS, I expect that GASPS could be improved by other methods. GASPS was built to require no similarity in folds so that motifs could be detected even when overall folds could not be aligned. One of the conclusions of this work, however, is that cases of motifs across different folds are very rare. With this knowledge, a faster method could use a structural alignment to first identify conserved regions that can identify the group with high sensitivity, and, if necessary, adjustments can be made to ensure specificity. Such a system may be better able to deal with the added degrees of freedom provided by a method that allows substitutions at specific positions. Another alternative is to build motifs from single atoms, chemical groups, or physical and chemical descriptors instead of residues.

In many ways, the development of this work parallels both the evolution of protein structure and function, as well as the genetic algorithm that provided the majority of the results for my work. All three systems make use of fortuitous occurrences with trial and error, leverage existing resources, and culminate in a successful product. The main fortuitous occurrence (there are plenty of others that were less successful) that contributed to this dissertation was the opportunity to collaborate with the initial superfamily active site templates project, described in Chapter 1. Just as today's proteins

recruited functional components that originally may have provided a different overall function, GASPS makes use of many Python functions written for the work of Chapter 1. In fact, my using a genetic algorithm as opposed to other techniques for building motifs, is a result of its development history, when other techniques may be equally, if not more, appropriate (see above). Finally, while I suggest above that the parallel includes a successful end-product of my work, I leave the evaluation of the success of this dissertation to the scoring or fitness functions that my readers bring with them.

## Appendix 1: GASPS Package

This appendix contains the text files included in the GASPS software package. This package was distributed as a gzipped, tar-formatted archive to other researchers requesting the source code for GASPS.

### ***ReadMe***

```
1 |
2 |
3 |
4 |
5 |
6 |
7 |
8 |
9 | This package should contain all the python code that is necessary to
   | run GASPS. In addition you will need the motif searching tool SPASM
   | and MKSPAZ to generate your own libraries. Together with their
   | manuals, these can be downloaded from:
10 |
11 | http://alpha2.bmc.uu.se/usf/spasm.html
12 |
13 | Just in case we need to say it: We are not responsible, nor hold any
   | ownership for any of the SPASM and MKSPAZ software available from
   | the above site.
14 |
15 | GASPS.py can construct a multiple sequence alignment by running PSI-
   | BLAST if it is so instructed. To use this you will need a copy of
   | the blastpgp program and sequence database available from NCBI:
16 |
17 | http://www.ncbi.nlm.nih.gov/BLAST/download.shtml
18 |
19 | Again: We are not responsible, nor hold any ownership for any of the
   | NCBI software.
20 |
21 | Minimally, once SPASM is installed you are ready to go. GASPS can be
   | controlled via many different command line arguments. The following
   | is a typical command and makes a reasonable test of your
   | installation. It should run from the GASPS_package directory with
   | the files located in test/. Any errors should result in a failed
   | execution within the first few minutes. A sequence of happy SPASM
   | messages to stderr "..Toodle Pip.." indicates that things are
   | probably running properly. Running to completion may take a few
   | hours or more depending on your system's speed. As configured (--
   | writeTables=1), it will probably use about 30MB or more to store the
   | output of all its SPASM runs, turn this off (--writeTables=0) if
   | you're tight on disk space.
22 |
23 |
```

```

24 | [~/GASPS_package] % python GASPS.py --pdbFile=test/d2mnr_1.pdb --
    | chain=' ' --filesPath=2mnr.test --tpFile=test/enolase.list --
    | trueLibrary=test/enolase.lib --lengthTrueLibrary=7 --
    | doNotCountQuery=d2mnr_1 --falseLibrary=test/astral_1.65_SF.lib --
    | alignFile=test/d2mnr_1.fasta.psiblast.xml.faln --refRowName=d2mnr_1
    | --writeTables=1 --useFileNames
25 |
26 |
27 | For more information on the purpose of the specific command line
    | arguments, read the GASPS.py file. Most useful output will appear
    | in 2mnr.test_log.txt, and the files used to run SPASM will all be in
    | 2mnr.test/*. The log.txt file will contain at least one line for
    | each motif attempted and its GASPS score. The winning motif can be
    | located near the end of the file. Additionally, the residues that
    | appear frequently among the top scoring motifs are given a score
    | that is simply the number of top scoring motifs they appear in. When
    | these residues are described, the potential list of substitutions
    | are listed at each of these positions regardless of whether
    | substitutions were turned on (--noSubs=0).
28 |
29 | Running GASPS from any other directory may require configuring your
    | python environment to find my modules located in
    | GASPS_package/polacco/. If you need help and don't know where to
    | look, try a web search for PYTHONPATH. For example:
30 |
31 | http://www.google.com/search?q=pythonpath.
32 |
33 | Additionally, if the spasm and blastpgp binaries (or databases) are
    | not in your shell's search paths, you may have to modify the
    | following lines in GASPS.py to point to the absolute paths of these
    | files:
34 |
35 | __spasmBinaryPath = "spasm"
36 | __blastpgpPath = "blastpgp"
37 | __blastDB = "nrdb90"
38 |
39 |
40 |
41 | The files in test/ are typical files that GASPS depends on for a
    | typical run:
42 | files in test/
43 |
44 |
45 |         d2mnr_1.fasta Sequence file corresponding to
46 |                               2mnr.pdb. use:
47 |                               --generateAlign=d2mnr_1.fasta
48 |
49 |         d2mnr_1.fasta.psiblast.xml.faln Multiple sequence alignment
    | already
50 |                               generated by GASPS. use:
51 |                               --
    | align=d2mnr_1.fasta.psiblast.xml.faln
52 |
53 |         d2mnr_1.pdb Structure file as a source of
    | motif
54 |                               coordinates. use: --
    | pdbFile=d2mnr_1.pdb
55 |

```







```

41 | #
42 | # May have to move polacco/*.py to a directory listed in PYTHONPATH
43 | # or modify PYTHONPATH to include parent directory of polacco/*.py
44 | import polacco.Spasm, polacco.MultiAlign, polacco.utils
45 | #
46 |
47 | # Trouble locating these, you may have to use absolute path here.
48 | __spasmBinaryPath = "spasm"
49 | __blastpgpPath = "blastpgp"
50 | __blastDB = "nrdb90"
51 |
52 |
53 | def DescribeMembers (info, openFile):
54 |     kees = info.__dict__.keys()
55 |     kees.sort()
56 |     for key in kees:
57 |         openFile.write( "%30s\t:\t%s\n" % (key, info.__dict__[key]) )
58 |     openFile.flush()
59 |
60 | def FileExists (filePath):
61 |     try:
62 |         fp = open (filePath)
63 |         fp.close()
64 |         return 1
65 |     except IOError:
66 |         return 0
67 |
68 | def PatternSampled (info, subset):
69 |     directory = os.path.join (info.filesPath, string.join (subset,
70 |     "_"))
71 |     doneFile = os.path.join (directory, "spasm.table")
72 |     return FileExists (doneFile) or FileExists (doneFile + ".gz")
73 |
74 | def WriteLog (logFile, string, newLine = 1):
75 |     fp = open (logFile, 'a')
76 |     fp.write (string)
77 |     if newLine:
78 |         fp.write ("\n")
79 |     fp.close()
80 |
81 | def ChooseConservationCutoff (conservationScores, numWanted,
82 |     referenceRow=None, allowedResidues=None):
83 |     cons = []
84 |     for i in range (len (conservationScores)):
85 |         if referenceRow.chars[i] in allowedResidues:
86 |             cons.append (conservationScores[i])
87 |
88 |     if numWanted < 1.0 and numWanted > 0.0:
89 |         # User is asking for a fraction,
90 |         # convert it to a number of residues based on number of
91 |         allowed residues
92 |         numWanted = numWanted * len (cons)
93 |
94 |     numWanted = int (numWanted)
95 |     if numWanted > len (cons):
96 |         print "After removing gaps and unwanted residues:Only %d
97 |         residues to choose from (wanted %d)" % (len (cons), numWanted)
98 |     return 0.0

```

```

95 |
96 |     cons.sort()
97 |     cons.reverse()
98 |     return cons[numWanted-1]
99 |
100 |
101 | def GetAvailableResidues(pdbFile, chain=None, model=1):
102 |     chains = {}
103 |     fp = open (pdbFile)
104 |     lastRes = ''
105 |     lastChain = ''
106 |     index = 0
107 |     models = 0
108 |     curModel = None
109 |     while (1):
110 |         line = fp.readline()
111 |         if (line == ''):
112 |             break
113 |
114 |
115 |         if line[0:5] == 'MODEL':
116 |             models+=1
117 |             curModel = int(line[11:16])
118 |             continue
119 |         elif line[0:6] == 'ENDMDL':
120 |             if curModel == model:
121 |                 break
122 |             curModel = None
123 |             continue
124 |         elif line[0:4] != 'ATOM':
125 |             continue
126 |
127 |         if curModel and curModel != model:
128 |             continue
129 |
130 |         curChain = line[21]
131 |         if (chain and curChain != chain):
132 |             if not chain in '?*': #special cases, ? is first and * is
all
133 |                 continue
134 |         res = string.strip( line[22:27])
135 |         atom = line[13:16]
136 |
137 |         if (res == lastRes and curChain == lastChain):
138 |             if locatedCA:
139 |                 continue
140 |         else:
141 |             locatedCA = 0
142 |             recorded = 0
143 |
144 |         if not recorded:
145 |             try:
146 |                 c = chains[curChain]
147 |             except KeyError:
148 |                 index = 0
149 |                 c = chains[curChain] = {}
150 |             c[res] = ( )
151 |             recorded = 1

```

```

152 |
153 |     if not locatedCA and atom == "CA ":
154 |         locatedCA = 1
155 |         x = float (line[32:38])
156 |         y = float (line[40:46])
157 |         z = float (line[48:55])
158 |         type = line[17:20]
159 |
160 |         chains[curChain][res] = (x,y,z,index, type)
161 |         index += 1
162 |         lastChain = curChain
163 |         lastRes = res
164 |     fp.close()
165 |     if chain and not chain in '?*':
166 |         return chains[chain]
167 |     elif chain == '?':
168 |         assert len(chains) == 1
169 |         return chains[chains.keys()[0]]
170 |     else:
171 |         return chains
172 |
173 |
174 | #simple Needleman-Wunsch to map the alignment sequence on to the
175 | structure sequence
176 | def Needleman(s1, s2, scoreMatch=1.0, scoreMismatch=-3.0, scoreGap=-
177 | 1.0):
178 |     m = []
179 |     for i1 in range(len(s1) + 1):
180 |         m.append((len(s2) + 1) * [ 0.0])
181 |
182 |     for i1 in range(len(s1)):
183 |         for i2 in range(len(s2)):
184 |             if s1[i1] == s2[i2]:
185 |                 best = m[i1][i2] + scoreMatch
186 |             else:
187 |                 best = m[i1][i2] + scoreMismatch
188 |                 skip = m[i1][i2+1] + scoreGap
189 |                 if skip > best:
190 |                     best = skip
191 |                 skip = m[i1+1][i2] + scoreGap
192 |                 if skip > best:
193 |                     best = skip
194 |                 m[i1+1][i2+1] = best
195 |     i1 = len(s1)
196 |     i2 = len(s2)
197 |     matchList = []
198 |     while i1 > 0 and i2 > 0:
199 |         best = m[i1-1][i2-1]
200 |         action = 0      # match
201 |         if m[i1][i2-1] > best:
202 |             best = m[i1][i2-1]
203 |             action = 1  # skip i2
204 |         if m[i1-1][i2] > best:
205 |             best = m[i1-1][i2]
206 |             action = 2  # skip i1
207 |         if action == 0:
208 |             matchList.append((i1-1, i2-1))
209 |             i1 = i1 - 1

```

```

208 |         i2 = i2 - 1
209 |         elif action == 1:
210 |             i2 = i2 - 1
211 |         else:
212 |             i1 = i1 - 1
213 |     return matchList
214 |
215 |
216 | def GetAvailableConservedResidues (info, pdbFile, chain, multiAlign,
referenceRow,
217 |                                     minConservation = 0.0,
allowedResidues = "FILVPAGMCWYTSQNEHDHKR", numResiduesAllowed = -1):
218 |     print "Processing multiple sequence alignment..."
219 |     #first load all residues from the pdbFile
220 |     allResidues = GetAvailableResidues (pdbFile, chain)
221 |     residueNames = allResidues.keys()
222 |     #reconstruct their order
223 |     namesInOrder = []
224 |     typesInOrder = []
225 |     for i in range (len(residueNames)):
226 |         namesInOrder.append ('?')
227 |         typesInOrder.append ('???')
228 |
229 |     for name in residueNames:
230 |         #print name
231 |         if len (allResidues[name]) < 5:
232 |             print "Warning: Trouble reading information from pdb for
residue %s" % (name)
233 |             WriteLog (info.logFile, "Warning: Trouble reading
information from pdb for residue %s" % (name))
234 |             continue
235 |             (index, type) = allResidues[name][3:5]
236 |             try:
237 |                 namesInOrder[index] = name
238 |             except IndexError, data:
239 |                 print data
240 |                 print index
241 |                 print namesInOrder
242 |                 print typesInOrder
243 |                 print pdbFile
244 |                 print allResidues
245 |                 sys.exit(0)
246 |             typesInOrder[index] = type
247 |
248 |     #align pdbSequence with referenceRow
249 |     pdbSeq = (polacco.utils.SeqAA3to1(typesInOrder))
250 |     refChars, refIndexes = referenceRow.GetCharsAndIndexesNoGaps()
251 |     matches = Needleman (pdbSeq, refChars)
252 |
253 |     #compute conservation
254 |     vc = polacco.MultiAlign.ValdarConservation (multiAlign)
255 |     conservations = vc.Compute()
256 |     if numResiduesAllowed > 0:
257 |         minConservation = ChooseConservationCutoff (conservations,
numResiduesAllowed, referenceRow, allowedResidues)
258 |         print "Using conservation cutoff = %6.4f" % minConservation
259 |         #get possible substitutions per position
260 |         #subs = multiAlign.GetLettersPerColumn ()

```

```

261 |     subs = multiAlign.GetDominantLettersPerColumn(0.1)
262 |
263 |     #generate list of user requested residues from user requested
motifs.
264 |     #these are forced to be included (with all their substitutions!)
regardless of their conservation
265 |     userRequestedResidues = []
266 |     for motif in info.motifs:
267 |         for res in motif:
268 |             if not res in userRequestedResidues:
269 |                 userRequestedResidues.append (res)
270 |
271 |     #map conservation scores to pdbSeq and return result
272 |     conResidues = {}
273 |     for match in matches:
274 |         name = namesInOrder[match[0]]
275 |         conservation = conservations[refIndexes[match[1]]]
276 |
277 |         if name in userRequestedResidues:
278 |             pass
279 |         elif not pdbSeq[match[0]] in allowedResidues:
280 |             continue
281 |         elif conservation < minConservation:
282 |             continue
283 |
284 |         if not refChars[match[1]] in subs[refIndexes[match[1]]]:
285 |             subs[refIndexes[match[1]]].append (refChars[match[1]])
286 |
287 |         conResidues[name] = (allResidues[name] + (conservation,
polacco.utils.SeqAAto3(subs[refIndexes[match[1]]]))
288 |
289 |
290 |     return conResidues
291 |
292 | def EucDistance (a, b):
293 |     sumSquares = 0.0
294 |     for i in range (len (a)):
295 |         sumSquares += math.pow(a[i]-b[i], 2)
296 |     return math.sqrt(sumSquares)
297 |
298 | def GetDistanceMatrix (allResidueLocations):
299 |     mat = {}
300 |     for res in allResidueLocations.keys():
301 |         mat[res] = {}
302 |         for other in allResidueLocations.keys():
303 |             if res == other:
304 |                 continue
305 |                 mat[res][other] =
EucDistance(allResidueLocations[res][0:3],
allResidueLocations[other][0:3])
306 |     return mat
307 |
308 |
309 | def MatChooseSpatiallyCloseSubset (allResidueLocations,
distanceMatrix, numResidues, maxRadius, res = None):
310 |     numResidues = int (numResidues)
311 |     chosen = []
312 |     next = None

```

```

313 |     i = 0
314 |     resNames = allResidueLocations.keys()
315 |     center = None
316 |     while (i < numResidues):
317 |         while (1):
318 |             if not center:
319 |                 if not res:
320 |                     next = random.choice (resNames)
321 |                 else:
322 |                     next = res
323 |                 if maxRadius != 99.9:
324 |                     possibleOthers = [x for x in
distanceMatrix[next].keys() if distanceMatrix[next][x] < maxRadius]
325 |                 else:
326 |                     possibleOthers = distanceMatrix[next].keys()
327 |                 if len (possibleOthers) < numResidues-1:
328 |                     if res:
329 |                         return None
330 |                     else:
331 |                         continue
332 |                 center = allResidueLocations[next][0:3]
333 |             else:
334 |                 next = random.choice (possibleOthers)
335 |                 if EucDistance (center, allResidueLocations[next][0:3])
> maxRadius:
336 |                     continue
337 |                 if not next in chosen:
338 |                     break
339 |                 chosen.append (next)
340 |                 i += 1
341 |         return chosen
342 |
343 | #mostly for debugging:
344 | def DescribePossibilities (distanceMatrix, cutoff, number, info):
345 |     WriteLog (info.logFile, 'From %d residues at distance cutoff = %d
requiring %d neighbors' % (len (distanceMatrix.keys()), cutoff,
number))
346 |     num = int (number) - 1
347 |     for res in distanceMatrix.keys():
348 |         l = len ([x for x in distanceMatrix[res].values() if x <
cutoff])
349 |         if l >= num:
350 |             WriteLog (info.logFile, '%s %s %s' % (res, l, [x for x in
distanceMatrix[res].keys() if distanceMatrix[res][x] < cutoff] ))
351 |
352 | def GenerateMotifFile (pdbFile, chain, residues, directory,
allResidues = None):
353 |     motifPath = os.path.join (directory, "motif.pdb")
354 |     motfp = open (motifPath, "w")
355 |     pdbfp = open (pdbFile)
356 |     resTypes = []
357 |     lastres = ''
358 |     models = 0
359 |     for line in pdbfp:
360 |         if line[0:5] == 'MODEL':
361 |             models = 1
362 |             continue
363 |         if models and line[0:6] == 'ENDMDL':

```

```

364         break
365     if line[0:4] != 'ATOM':
366         continue
367     curChain = line[21]
368     if chain == '?':
369         chain = curChain
370     if (curChain != chain):
371         continue
372     res = string.strip( line[22:27])
373     if res in residues:
374         motfp.write (line)
375         if res != lastres:
376             try:
377                 if allResidues and len (allResidues[res]) > 6:
378                     restypes.append (allResidues[res][6])
379             else:
380                 restypes.append ((line[17:20],))
381         except KeyError:
382             #res not found in allResidues, must be a user
supplied motif
383             restypes.append ((line[17:20],))
384             lastres = res
385
386
387     # add remark indicating the allowed substitutions we expect:
SPASM and GASPS do not use this!
388     motfp.write ("REMARK * For note only, spasm does not use
this!\n")
389     motfp.write ("REMARK * restypes:")
390     for resType in restypes:
391         motfp.write (" %s" % string.join (resType, "/"))
392     motfp.write ("\n")
393     motfp.close()
394     pdbfp.close()
395     #print restypes
396     return (motifPath, restypes)
397
398
399
400 def GenerateSpasmRunFile (motifFile, restypes, directory, library =
"", runPath = 'spasm.com', outPath = 'spasm.out', info = None):
401     spasmBinaryPath = __spasmBinaryPath
402
403     maxHits = 100000 # set this arbitrarily high so that SPASM never
stops early
404     maxRMSD = 3.2
405     maxCADiff = 5.0
406     maxSCDiff = 3.8
407     scOnly = 0
408     if info:
409         maxRMSD = info.maxRMSD
410         maxCADiff = info.maxCADiff
411         maxSCDiff = info.maxSCDiff
412         scOnly = info.scOnly
413
414     maxResolution = 999.9
415     maxResidues = 9999
416

```



```

417 | subStrings = []
418 | for subList in restTypes:
419 |     subStrings.append (string.join (subList, " "))
420 |
421 | substituteString = string.join (subStrings, "\n")
422 |
423 | fp = open (runPath, "w")
424 | if scOnly:
425 |     spasmRunFileString = polacco.Spasm.runFileStringSTDOUT_scOnly
426 | else:
427 |     spasmRunFileString = polacco.Spasm.runFileStringSTDOUT
428 |
429 | fp.write (spasmRunFileString % (
430 |                                     spasmBinaryPath,
431 |                                     maxHits,
432 |                                     library,
433 |                                     motifFile,
434 |                                     'rand',
435 |                                     maxRMSD,
436 |                                     maxCADiff,
437 |                                     maxSCDiff,
438 |                                     maxResolution,
439 |                                     maxResidues,
440 |                                     5,
441 |                                     substituteString))
442 | fp.close()
443 | os.chmod (runPath, 0755)
444 | return runPath
445 |
446 |
447 | def LoadTrueHash (tpFile):
448 |     tpHash = {}
449 |     if tpFile:
450 |         fp = open (tpFile)
451 |         while (1):
452 |             item = string.strip (fp.readline())
453 |             if item == '':
454 |                 break
455 |             item = item.upper()
456 |             tpHash[item] = 1
457 |         fp.close()
458 |         print "Loaded %d unique identifiers from %s" %
459 |             (len(tpHash.keys()), tpFile)
460 |         return tpHash
461 |
462 | def SetupSpasmFiles (info, subset, allResidues):
463 |     directory = os.path.join (info.filesPath, string.join (subset,
464 |     "_"))
465 |     if info.scratchPath:
466 |         spasmTrueOutFile = os.path.join (info.scratchPath,
467 |         string.join (subset, "_") + "true_spasm.out")
468 |         spasmFalseOutFile = os.path.join (info.scratchPath,
469 |         string.join (subset, "_") + "false_spasm.out")
470 |     else:
471 |         spasmTrueOutFile = os.path.join (directory,
472 |         "true_spasm.out")
473 |         spasmFalseOutFile = os.path.join (directory,
474 |         "false_spasm.out")

```

```

469 |
470 |
471 |     spasmTableFile = os.path.join (directory, "spasm.table")
472 |     try:
473 |         os.makedirs (directory)
474 |     except OSError, data:
475 |         WriteLog (info.logFile, "Error (ignored) while generating
directory: %s" % directory)
476 |         WriteLog (info.logFile, data.strerror)
477 |         if info.noSubs:
478 |             (motif, resTypes) = GenerateMotifFile (info.pdbFile,
info.chain, subset, directory, 0)
479 |         else:
480 |             (motif, resTypes) = GenerateMotifFile (info.pdbFile,
info.chain, subset, directory, allResidues)
481 |             runFileTrue = os.path.join (directory, "true_spasm.com")
482 |             runFileFalse = os.path.join (directory, "false_spasm.com")
483 |             GenerateSpasmRunFile (motif, resTypes, directory,
info.falseLibrary, runFileFalse, spasmFalseOutFile, info)
484 |             GenerateSpasmRunFile (motif, resTypes, directory,
info.trueLibrary, runFileTrue, spasmTrueOutFile, info)
485 |             return info, runFileFalse, runFileTrue, spasmTableFile
486 |
487 |
488 | def DoSpasmRuns (info, runFileFalse, runFileTrue, spasmTableFile,
writeOutFile = 0):
489 |     fpFalseSpasm = os.popen ("csh %s" % runFileFalse)
490 |     falseSearch = polacco.Spasm.SpasmSearch (1)
491 |     falseSearch.titleFromFileName = info.useFileNames
492 |     falseSearch.ParseSpasmHits (fpFalseSpasm)
493 |     fpFalseSpasm.close()
494 |
495 |     fpTrueSpasm = os.popen ("csh %s" % runFileTrue)
496 |     trueSearch = polacco.Spasm.SpasmSearch (1)
497 |     trueSearch.titleFromFileName = info.useFileNames
498 |     trueSearch.ParseSpasmHits (fpTrueSpasm)
499 |     fpTrueSpasm.close()
500 |
501 |     returnAsString = 1
502 |     tableFileString =
polacco.Spasm.Convert2SpasmSearchesToSortedAndScoredTable
(trueSearch, falseSearch,
503 |                                     spasmTableFile,
info.trueHash, info.useDistanceRmsd, returnAsString,
504 |                                     writeOutFile,
info.trueSkipHash, info.falseSkipHash)
505 |     return tableFileString
506 |
507 | def ScoreMotifs (info, population, allResidues, scores = {}):
508 |     for subset in population:
509 |         info, runFileFalse, runFileTrue, spasmTableFile =
SetupSpasmFiles (info, subset, allResidues)
510 |
511 |         WriteLog (info.logFile, "guess %3d %40s" % (info.round,
string.join(subset, "_")), 0)
512 |
513 |         if not info.testing:

```

```

514 |         tableFileString = DoSpasmRuns(info, runFileFalse,
runFileTrue, spasmTableFile, info.writeTables)
515 |         if not tableFileString:
516 |             WriteLog (info.logFile, " WARNING!: No table file string
generated for guess %3d %40s; score set to 0.0" % (info.round,
string.join(subset, "_")))
517 |             print ( "WARNING!: No table file string generated for
guess %3d %40s; score set to 0.0" % (info.round, string.join(subset,
"_")))
518 |                 score = 0.0
519 |
520 |             elif info.rocArea:
521 |                 score = polacco.Spasm.ComputeAreaFromTableFile
(spasmTableFile, info.maxFalse, tableFileString,
info.useDistanceRmsd)
522 |             else:
523 |                 score =
polacco.Spasm.ComputeSeparationScoreFromTableFile3 (spasmTableFile,
info.maxFalse,
524 |                 info.maxRMSD, info.lengthTrueLibrary,
info.sepScoreImportance, info.useDistanceRmsd, tableFileString)
525 |             else:
526 |                 score = random.random ()
527 |                 scores[tuple(subset)] = score
528 |
529 |             WriteLog (info.logFile, " %8.4f" % ( score ))
530 |
531 |     return scores
532 |
533 | def GetResTypesFromMotifFile (motifFile):
534 |     fp = open (motifFile, "r")
535 |     restypes = []
536 |     while (1):
537 |         #"REMARK    1 restypes:"
538 |         line = fp.readline()
539 |         if line == '':
540 |             break
541 |         if line[0:6] == 'REMARK' and line[12:21] == "restypes:":
542 |             for res in line[21:].strip().split(' '):
543 |                 res = res.strip()
544 |                 if len(res) >= 3:
545 |                     restypes.append (res.split('/'))
546 |     fp.close()
547 |     return restypes
548 |
549 | def GetResTypesFromSpasmRunFile (runFile):
550 |     resList = polacco.utils.aa3to1.keys()
551 |     fp = open (runFile, "r")
552 |     restypes = []
553 |     while (1):
554 |         line = fp.readline()
555 |         if line == '':
556 |             break
557 |         words = line.strip().split (" ")
558 |         if words[0] in resList:
559 |             restypes.append (words)
560 |     fp.close()
561 |     return restypes

```

```

562 |
563 | def TestMotif (directory, trueLibrary, falseLibrary):
564 |     testTrueOutFile = os.path.join (directory,
565 |         "testTrue_spasm.out")
566 |     testFalseOutFile = os.path.join (directory,
567 |         "testFalse_spasm.out")
568 |
569 |     testTableFile = os.path.join (directory, "test.table")
570 |     motifFile = os.path.join (directory, "motif.pdb")
571 |
572 |     #load allowed residue types from motifFile
573 |     resTypes = GetResTypesFromMotifFile (motifFile)
574 |     if not resTypes:
575 |         print "WARNING! restypes not found in motif file!"
576 |         resTypes = GetResTypesFromSpasmRunFile (os.path.join
577 |             (directory, "false_spasm.com"))
578 |         if not resTypes:
579 |             print "Failed again reading from run file!"
580 |             print resTypes
581 |             runFileFalse = os.path.join (directory, "testFalse_spasm.com")
582 |             runFileTrue = os.path.join (directory, "testTrue_spasm.com")
583 |
584 |             GenerateSpasmRunFile (motifFile, resTypes, directory,
585 |                 falseLibrary, runFileFalse, testFalseOutFile)
586 |             GenerateSpasmRunFile (motifFile, resTypes, directory,
587 |                 trueLibrary, runFileTrue, testTrueOutFile)
588 |             print runFileFalse
589 |             os.spawnlp (os.P_WAIT, "sh", "sh", runFileFalse)
590 |             print runFileTrue
591 |             os.spawnlp (os.P_WAIT, "sh", "sh", runFileTrue)
592 |
593 |             polacco.Spasm.Convert2SpasmFilesToSortedAndScoredTable
594 |                 (testTrueOutFile, testFalseOutFile, testTableFile)
595 |             return testTableFile
596 |
597 |
598 |
599 | def IsSuperSet (super, other):
600 |     for item in other:
601 |         if not item in super:
602 |             return 0
603 |         else:
604 |             return 1
605 |
606 |
607 | def RemoveSameScoringSupersets (info, motifScores, motifHash,
608 |     upForRemoval):
609 |     toRemove = []
610 |     for mot in upForRemoval:
611 |         score = motifScores[mot]
612 |         otherMots = motifHash[score]
613 |         for otherMot in otherMots:
614 |             if otherMot == mot:
615 |                 continue
616 |             if IsSuperSet (mot, otherMot):
617 |                 WriteLog (info.logFile, "Removing the same scoring
618 |                     superset: %s (%s : %d)" % (mot, otherMot, score))
619 |                 toRemove.append (mot)
620 |                 break
621 |     for mot in toRemove:

```

```

612 |         score = motifScores[mot]
613 |         del (motifScores[mot])
614 |         motifHash[score].remove (mot)
615 |         if len (motifHash[score]) == 0:
616 |             del (motifHash[score])
617 |
618 | def GetTopScorers (info, number, motifScores, previousTop):
619 |     topScorers = {}
620 |     motHash = {}
621 |     #reverse the score hash:
622 |     for mot in motifScores.keys():
623 |         try:
624 |             motHash[motifScores[mot]].append(mot)
625 |         except KeyError:
626 |             motHash[motifScores[mot]] = [mot]
627 |
628 |     RemoveSameScoringSupersets (info, motifScores, motHash,
previousTop.keys())
629 |
630 |     scores = motHash.keys()
631 |     scores.sort()
632 |     scores.reverse()
633 |
634 |     for score in scores:
635 |         if score == 0.0:
636 |             break
637 |         motifs = motHash[score]
638 |         if len(topScorers) >= number:
639 |             break
640 |         if len(topScorers)==0 or
len(topScorers)+len(motifs)<=number*2:
641 |             for mot in motifs:
642 |                 topScorers[mot] = score
643 |         else:
644 |             break
645 |     return topScorers
646 |
647 | def MakeRandomGuessesWithCoverage (info,distanceMatrix, allResidues,
numGuesses, population = []):
648 |     if numGuesses==0:
649 |         return population
650 |     allRes = allResidues.keys()
651 |     lenAllRes = len(allRes)
652 |     if lenAllRes > numGuesses:
653 |         skip = lenAllRes/numGuesses
654 |     else:
655 |         skip = 1
656 |
657 |     for i in range (0, lenAllRes, skip):
658 |         res = allRes[i]
659 |         subset = MatChooseSpatiallyCloseSubset (allResidues,
distanceMatrix, info.numResidues, info.maxNeighborhood, res)
660 |         if not subset:
661 |             continue
662 |         subset.sort()
663 |         if subset in population:
664 |             continue
665 |

```

```

666         if PatternSampled (info, subset):
667             continue
668         population.append (subset)
669
670     return population
671
672 def MakeRandomGuesses (info, distanceMatrix, allResidues,
numGuesses, population = []):
673     totalSize = len (population) + numGuesses
674     while (len (population) < totalSize):
675         for i in range (10000): #TODO set up a better check to make
sure there are choices left before continuing the loop.
676             subset = MatChooseSpatiallyCloseSubset (allResidues,
distanceMatrix, info.numResidues, info.maxNeighborhood)
677             subset.sort()
678             if subset in population:
679                 continue
680             if PatternSampled (info, subset):
681                 continue
682             population.append (subset)
683             break
684         else:
685             WriteLog (info.logFile , "Could not find a random guess not
already tried after 10000 tries!")
686             break
687     return population
688
689
690 def MakeMutations (survivors, allResidues, info, numMutations,
population):
691     resNames = allResidues.keys()
692     totalSize = len (population) + numMutations
693
694     allMutations = []
695     for parent in survivors:
696         for pres in parent:
697             for mres in resNames:
698                 if mres in parent:
699                     continue
700                 newMotif = list (parent)
701                 newMotif.remove (pres)
702                 newMotif.append (mres)
703                 newMotif.sort()
704                 allMutations.append (newMotif)
705     while (allMutations and len (population) < totalSize):
706         newMotif = random.choice (allMutations)
707         allMutations.remove (newMotif)
708         if newMotif in population:
709             continue
710         if newMotif in survivors:
711             continue
712         if PatternSampled (info, newMotif):
713             continue
714         population.append (newMotif)
715     return population
716
717

```

```

718 | def MakeInsertions (survivors, allResidues, info, numInsertions,
      | population):
719 |     resNames = allResidues.keys()
720 |     totalSize = len (population) + numInsertions
721 |
722 |     allInsertions = []
723 |     for parent in survivors:
724 |         if len (parent) >= info.maxResidues:
725 |             continue
726 |         for res in resNames:
727 |             if res in parent:
728 |                 continue
729 |             newMotif = list (parent)
730 |             newMotif.append (res)
731 |             newMotif.sort()
732 |             allInsertions.append (newMotif)
733 |
734 |     while (allInsertions and len (population) < totalSize):
735 |         newMotif = random.choice (allInsertions)
736 |         allInsertions.remove (newMotif)
737 |         if newMotif in population:
738 |             continue
739 |         if newMotif in survivors:
740 |             continue
741 |         if PatternSampled (info, newMotif):
742 |             continue
743 |         population.append (newMotif)
744 |     return population
745 |
746 | def MakeDeletions (survivors, info, numDeletions, population):
747 |     totalSize = len (population) + numDeletions
748 |     allDeletions = []
749 |     for parent in survivors:
750 |         if len (parent) <= info.minResidues:
751 |             continue
752 |         for res in parent:
753 |             newMotif = list (parent)
754 |             newMotif.remove (res)
755 |             allDeletions.append (newMotif) # allDeletions may have
      | duplicates, but only a finite number so its okay
756 |
757 |     while (allDeletions and len (population) < totalSize):
758 |         newMotif = random.choice (allDeletions)
759 |         allDeletions.remove (newMotif)
760 |         if newMotif in population:
761 |             continue
762 |         if newMotif in survivors:
763 |             continue
764 |         if PatternSampled (info, newMotif):
765 |             continue
766 |         population.append (newMotif)
767 |     return population
768 |
769 | def MakeRecombinations (parents, info, numRecombinations,
      | population):
770 |     totalSize = len (population) + numRecombinations
771 |     attempts = 0
772 |     while (len (population) < totalSize and attempts < 10000):

```

```

773 |         attempts = attempts + 1
774 |         if not parents:
775 |             break
776 |         #first choose parents
777 |         mom = random.choice (parents)
778 |         tmp = parents[:]
779 |         if len (tmp) > 1:
780 |             tmp.remove(mom)
781 |         dad = random.choice (tmp)
782 | #         choose contributions
783 |         all = list(mom)
784 |         for res in dad:
785 |             if res in mom:
786 |                 all.remove(res)
787 |         all = all + list(dad)
788 |         newMotif = []
789 |         maxlen = min (info.maxResidues, len(all))
790 |         newLen = random.randrange (3, maxlen+1)
791 |         while (len (newMotif) < newLen):
792 |             newMotif.append (random.choice(all))
793 |             all.remove (newMotif[-1])
794 |
795 |         newMotif.sort()
796 |         if newMotif in population:
797 |             continue
798 |         if newMotif in parents:
799 |             continue
800 |         if PatternSampled (info, newMotif):
801 |             continue
802 |         population.append (newMotif)
803 |         if (attempts >= 10000): #lazy but sufficient:
804 |             WriteLog (info.logFile, "WARNING: Couldn't find a
recombination that hasn't been tried in 10000 tries!" )
805 |             return population
806 |
807 | def EvolveNextPopulation (survivors, info, allResidues,
distanceMatrix):
808 |     nextPop = []
809 |     l = len (nextPop)
810 |     if info.motifs:
811 |         nextPop.extend (info.motifs)
812 |         info.motifs = []
813 |         WriteLog (info.logFile, "***User specified starting motifs**")
814 |         PrintPopulation (info, nextPop[l:])
815 |         l = len (nextPop)
816 |
817 |     if len (survivors) < info.popFromPrevious/2: # no population
bottlenecks
818 |         nextPop = MakeRandomGuessesWithCoverage (info, distanceMatrix,
allResidues, info.populationSize, nextPop)
819 |         WriteLog (info.logFile, "***Random Guesses With Coverage")
820 |         PrintPopulation (info, nextPop[l:])
821 |
822 |
823 |         l = len (nextPop)
824 |         nextPop = MakeRandomGuesses (info, distanceMatrix,
allResidues, info.populationSize - len(nextPop), nextPop)
825 |         WriteLog (info.logFile, "***Random Guesses")

```



```

826 |         PrintPopulation (info, nextPop[l:])
827 |
828 |
829 |     else:
830 |         nextPop = MakeRandomGuesses (info, distanceMatrix,allResidues,
info.popFromRandom, nextPop)
831 |         WriteLog (info.logFile, "***Random Guesses")
832 |         PrintPopulation (info, nextPop[l:])
833 |
834 |         l = len (nextPop)
835 |         nextPop = MakeMutations (survivors, allResidues, info,
info.popMutations, nextPop)
836 |         WriteLog (info.logFile, "***Mutations")
837 |         PrintPopulation (info, nextPop[l:])
838 |
839 |         l = len (nextPop)
840 |         nextPop = MakeInsertions (survivors, allResidues, info,
info.popInsertions, nextPop)
841 |         WriteLog (info.logFile, "***Insertions")
842 |         PrintPopulation (info, nextPop[l:])
843 |
844 |         l = len (nextPop)
845 |         nextPop = MakeDeletions (survivors,info, info.popDeletions,
nextPop)
846 |         WriteLog (info.logFile, "***Deletions")
847 |         PrintPopulation (info, nextPop[l:])
848 |
849 |         l = len (nextPop)
850 |         nextPop = MakeRecombinations (survivors, info,
info.populationSize - len(nextPop), nextPop)
851 |         WriteLog (info.logFile, "***Recombinations")
852 |         PrintPopulation (info, nextPop[l:])
853 |     return nextPop
854 |
855 |
856 | def PrintPopulation (info, population):
857 |     for mot in population:
858 |         WriteLog (info.logFile, string.join (mot, '_' ) )
859 |
860 | def PrintMotifScores (info, motifScores, finalStats=0):
861 |     scorePairs = []
862 |     total = 0.0
863 |     maxScore = None
864 |     minScore = 2.0
865 |     meanScore = 0.0
866 |     count = 0
867 |
868 |     for mot in motifScores.keys():
869 |         count += 1
870 |         total = total + motifScores[mot]
871 |         maxScore = max (maxScore, motifScores[mot])
872 |         minScore = min (minScore, motifScores[mot])
873 |         scorePairs.append ( (motifScores[mot], string.join (mot, '_' )
) )
874 |
875 |     if not count:
876 |         WriteLog (info.logFile, "WARNING: No top scorers found to
describe!")

```

```

877 |         return
878 |
879 |     meanScore = total/count
880 |     scorePairs.sort()
881 |     #loop again to print, and calculate variance
882 |     total = 0.0
883 |     for pair in scorePairs:
884 |         total += math.pow (pair[0]-meanScore, 2)
885 |         if not finalStats:
886 |             WriteLog (info.logFile,"topScorer %d %40s %8.4f" %
(info.round, pair[1], pair[0]))
887 |             if count > 2:
888 |                 variance = total/(count - 1)
889 |             else:
890 |                 variance = 0.0
891 |             if finalStats:
892 |                 label = "finalTopScoreStats"
893 |             else:
894 |                 label = "topScoreStats"
895 |
896 |             WriteLog (info.logFile, "%s %d %6.4f %6.4f %6.4f %10.8f" %
(label, info.round, meanScore, maxScore, minScore,
math.sqrt(variance)))
897 |
898 | def GetTopScorer (motifScores):
899 |     scorePairs = []
900 |     for mot in motifScores.keys():
901 |         scorePairs.append ( (motifScores[mot], string.join (mot, '_')
) )
902 |     topPair = max (scorePairs)
903 |     return topPair
904 |
905 |
906 |
907 | def SummarizeResiduesInMotifScores (info, motifScores,
allResidues, finalSummary = 0):
908 |     res = {}
909 |     for mot in motifScores.keys():
910 |         for res in mot:
911 |             try:
912 |                 res[res] += 1
913 |             except KeyError:
914 |                 res[res] = 1
915 |     pairs = []
916 |     for res in res.keys():
917 |         pairs.append ( (res[res], res) )
918 |     pairs.sort()
919 |     pairs.reverse()
920 |     if not finalSummary:
921 |         for score,res in pairs:
922 |             if len (allResidues[res]) > 6:
923 |                 resTypes = (allResidues[res][6])
924 |             else:
925 |                 resTypes = []
926 |             subString = string.join (resTypes, ",")
927 |
928 |             WriteLog (info.logFile, "resScore %d %8s.%s %3d # %s" %
(info.round, res, allResidues[res][4], score, subString))

```

```

929 |     else:
930 |
931 |         summaryString = ""
932 |         for score,res in pairs:
933 |             summaryString += "%s.%s(%02d); "% (res,
allResidues[res][4], score)
934 |             WriteLog (info.logFile, "finalResidues " + summaryString)
935 |
936 |
937 | def PrintFinalSummary (info, motifScores, allResidues):
938 |     #first print top score stats:
939 |     PrintMotifScores (info, motifScores, 1)
940 |     #now single line for final residues
941 |     SummarizeResiduesInMotifScores (info, motifScores, allResidues,
1)
942 |
943 |
944 |
945 | # The meat of GASPS work is performed here.
946 |
947 | def DoGASpasm (info):
948 |     try:
949 |         os.makedirs (os.path.split (info.logFile)[0])
950 |     except OSError:
951 |         pass #most likely directory already exists.
952 |
953 |     WriteLog (info.logFile, "\n\nGASpasm started at %s" %
time.ctime(time.time()))
954 |     # First get information from multi align file if appropriate
955 |     if info.alignFile:
956 |         ma = polacco.MultiAlign.MultiAlign([])
957 |         WriteLog (info.logFile, "Opening multialign file at %s" %
info.alignFile)
958 |         fp = open (info.alignFile)
959 |         if (string.upper (info.alignFormat) == 'FASTA'):
960 |             ma.read_fasta (fp)
961 |         elif (string.upper (info.alignFormat) == 'CLUSTAL'):
962 |             ma.read_clustal (fp, 0) # 0 sets to not force strict
clustal: first line "CLUSTAL..." is optional
963 |         else:
964 |             print "Unrecognized align format: %s " % info.alignFormat
965 |             sys.exit(0)
966 |         fp.close()
967 |         WriteLog (info.logFile, "Finished align file at %s" %
time.ctime(time.time()))
968 |
969 |         if info.alignRange:
970 |             (start,stop) = info.alignRange.split(':')
971 |             if start:
972 |                 start = int (start)
973 |             else:
974 |                 start = ma.FIRST_POSITION
975 |             if stop:
976 |                 stop = int (stop)
977 |             else:
978 |                 stop = ma.FIRST_POSITION + ma.length()
979 |

```

```

980 |         WriteLog (info.logFile, "Restricting alignemnt to %s and
      | %s" % (start, stop))
981 |         ma = ma.getMultiAlignBlock (start, stop)
982 |
983 |         refRow = ma.get_row_by_name (info.refRowName)
984 |         ma.protectedRows.append (refRow)
985 |         if info.numSeqsInAlign > 0:
986 |             WriteLog (info.logFile, "Shrinking Multiple alignment")
987 |             ma.ShrinkByRemovingRedundancy_Efficient
      | (info.numSeqsInAlign)
988 |
989 |         allResidues = GetAvailableConservedResidues (info,
      | info.pdbFile, info.chain, ma, refRow, info.minConservation,
      | info.allowedResidues, info.numTopConservedResidues)
990 |     else:
991 |         allResidues = GetAvailableResidues (info.pdbFile, info.chain)
992 |
993 |         distanceMatrix = GetDistanceMatrix (allResidues)
994 |         #mostly for debugging purposes:
995 |         DescribePossibilities (distanceMatrix, info.maxNeighborhood,
      | info.numResidues, info)
996 |
997 |         topX = {}
998 |         motifScores = {}
999 |         for round in range (info.numRounds):
1000 |             info.round = round
1001 |             WriteLog (info.logFile, "Round %3d started at %s" % (round,
      | time.ctime(time.time())) )
1002 |             population = EvolveNextPopulation (topX.keys(), info,
      | allResidues, distanceMatrix)
1003 |             motifScores = ScoreMotifs (info, population, allResidues,
      | motifScores)
1004 |             topX = GetTopScorers (info, info.popFromPrevious, motifScores,
      | topX)
1005 |             WriteLog (info.logFile, "New Top Scorers:")
1006 |             PrintMotifScores (info, topX)
1007 |             WriteLog (info.logFile, "Residue scores:")
1008 |             SummarizeResiduesInMotifScores (info, topX, allResidues)
1009 |
1010 |             PrintFinalSummary (info, topX, allResidues)
1011 |             if info.xValidate:
1012 |                 score, motif_ = GetTopScorer (topX)
1013 |                 directory = os.path.join (info.filesPath, motif_)
1014 |                 ValidateDirectory (info, directory)
1015 |
1016 |             if info.writeTables:
1017 |                 score, motif_ = GetTopScorer (topX)
1018 |                 tableFile = os.path.join (info.filesPath, motif_,
      | "spasm.table.gz")
1019 |                 localName = info.filesPath + motif_ + "_table.gz"
1020 |                 cmd = 'cp %s %s' % (tableFile, localName)
1021 |                 os.system(cmd)
1022 |
1023 |             WriteLog (info.logFile, "Finished at %s" %
      | (time.ctime(time.time())))
1024 |
1025 |
1026 | def ValidateDirectory (info, directory):

```

```

1027 |     testTrueOutFile = os.path.join (directory,
      |     "testTrue_spasm.out")
1028 |     testFalseOutFile = os.path.join (directory,
      |     "testFalse_spasm.out")
1029 |
1030 |     testTableFile = os.path.join (directory, "test.table")
1031 |     motifFile = os.path.join (directory, "motif.pdb")
1032 |     #load allowed residue types from motifFile
1033 |     resTypes = GetResTypesFromMotifFile (motifFile)
1034 |     if not resTypes:
1035 |         print "WARNING!  restypes not found in motif file!"
1036 |         resTypes = GetResTypesFromSpasmRunFile (os.path.join
      | (directory, "false_spasm.com"))
1037 |         if not resTypes:
1038 |             print "Failed again reading from run file!"
1039 |         runFileFalse = os.path.join (directory, "testFalse_spasm.com")
1040 |         runFileTrue = os.path.join (directory, "testTrue_spasm.com")
1041 |
1042 |     GenerateSpasmRunFile (motifFile, resTypes, directory,
      | info.xFalseLibrary, runFileFalse, testFalseOutFile)
1043 |     GenerateSpasmRunFile (motifFile, resTypes, directory,
      | info.trueLibrary, runFileTrue, testTrueOutFile)
1044 |
1045 |     fpFalseSpasm = os.popen ("csh %s" % runFileFalse)
1046 |     falseSearch = polacco.Spasm.SpasmSearch (1)
1047 |     falseSearch.ParseSpasmHits (fpFalseSpasm)
1048 |     fpFalseSpasm.close()
1049 |
1050 |     fpTrueSpasm = os.popen ("csh %s" % runFileTrue)
1051 |     trueSearch = polacco.Spasm.SpasmSearch (1)
1052 |     trueSearch.ParseSpasmHits (fpTrueSpasm)
1053 |     fpTrueSpasm.close()
1054 |
1055 |
1056 |     trueSkipHash = LoadTrueHash (info.xTrueSkipFile)
1057 |
1058 |     returnAsString = 1
1059 |     writeOutFile = 1 #always write out this file, the most
      | interesting one!
1060 |     tableFileString =
      | polacco.Spasm.Convert2SpasmSearchesToSortedAndScoredTable
      | (trueSearch, falseSearch,
1061 |                                     testTableFile, info.trueHash,
      | info.useDistanceRmsd, returnAsString,
1062 |                                     writeOutFile, trueSkipHash)
1063 |     if info.rocArea:
1064 |         score = polacco.Spasm.ComputeAreaFromTableFile (testTableFile,
      | info.maxFalse, tableFileString)
1065 |     else:
1066 |         lengthTrueLibrary = 1
1067 |         score = polacco.Spasm.ComputeSeparationScoreFromTableFile3
      | (testTableFile, info.maxFalse,
1068 |         info.maxRMSD, lengthTrueLibrary, info.sepScoreImportance,
      | info.useDistanceRmsd, tableFileString)
1069 |
1070 |     WriteLog (info.logFile, "Cross-Validate Result %s  %8.4f" %
      | (directory , score ))
1071 |

```

```

1072 | def SetDefaults(info):
1073 |     info.chain = ' '
1074 |     info.numResidues = 5 # first guesses motif size, min and max
specify that allowed during optimization
1075 |     info.maxResidues = 10 # spasm won't display ca-ca or sc-sc
matrices if this is any higher than 10
1076 |     info.minResidues = 3
1077 |     info.maxNeighborhood = 12 # in angstroms, initial guess motifs
are built from residues within this distance
1078 |     info.alignFile = None #if alignment exists already, read it
here.
1079 |     info.alignFormat = 'FASTA' # format of above, alternative is
CLUSTAL
1080 |     info.minConservation = 0.6 # conservation necessary for inclc
1081 |     info.filesPath = None
1082 |     info.allowedResidues = 'FILVPAGMCWYTSQNEDHKR'
1083 |     info.logFile = None
1084 |     info.tpFile = None
1085 |     info.falseSkipFile = None #these specify which items in the true
or fales libraries should be skipped.
1086 |     info.trueSkipFile = None
1087 |
1088 |     info.populationSize = 36
1089 |     info.popFromPrevious = 16
1090 |     info.popFromRandom = 0
1091 |     info.popInsertions = 8
1092 |     info.popMutations = 12
1093 |     info.popDeletions = 8
1094 |     info.noSubs = 1 # true or false indicating wether to turn off
substitutions (beware double negative!)
1095 |     info.maxFalse = 5 # cut off for computing ROC scores
1096 |     info.numRounds = 50 # number of rounds to complete before
stopping
1097 |     info.testing = 0
1098 |     info.rocArea = 0 #binary indicating what kind of scoring to use
(rocArea, vs separation score)
1099 |     info.numSeqsInAlign = -1 #if greater than zero this determines
the size of the multi-align to use
1100 |     info.numTopConservedResidues = 100 #if greater than 0 this
determines the number of positions in the multialign to choose as
conserved residues
1101 |                                     #if less than 1 it specifies
fraction to accept, if greater than 1 it represents the number of
residues to accept
1102 |     info.alignRange = ':' #specifies a range of columns to limit the
multialign
1103 |
1104 |     info.trueLibrary = None
1105 |     info.falseLibrary = None
1106 |
1107 |     info.validateDirectory = None
1108 |
1109 |     info.lengthTrueLibrary = None
1110 |     info.sepScoreImportance = 0.1
1111 |
1112 |     info.useDistanceRmsd = 0
1113 |
1114 |     info.motifs = []

```

```

1115 |     info.scratchPath = ''
1116 |
1117 |     info.xValidate = 0 #perform cross validation at the end of a
      | completed run.
1118 |             #reuse tp library, but give different exclude list
1119 |             #must give new fp library.
1120 |     info.xTrueSkipFile = None
1121 |     info.xFalseLibrary = None
1122 |
1123 |     info.doNotCountQuery = None # used to pass the name of the query
      | in the the spasm library
1124 |             # so that matches to itself can be
      | ignored.
1125 |
1126 |     info.writeTables=0      #write out spasm output tables from each
      | run. Eats disk space.
1127 |     info.useFileNames=0    #use the file name given in the spasm
      | library to describe the matched structure.
1128 |             # the alternative is to use the four
      | character pdb-style code given in the spasm library.
1129 |             # if turned on, this will use "dlqcrd2"
      | from /pdbstyle-1.63/qc/dlqcrd2.ent
1130 |             # used by
      | polacco.Spasm.SpasmHit.ReadOpenFile()
1131 |
1132 |     info.generateAlignment = None #file to use as input to psiblast
      | for generation of alignment
1133 |
1134 |
1135 |     info.maxRMSD = 3.2 #thresholds to be passed to spasm
1136 |     info.maxCADiff = 5.0
1137 |     info.maxSCDiff = 3.8
1138 |
1139 |     info.scOnly = 0 #use spasm in sidechain-only mode; ignore alpha
      | carbons
1140 |
1141 |
1142 |
1143 |
1144 | def PrintUsage(short, long):
1145 |     for i in range (len(short)):
1146 |         print '%s      %s' % (short[i], long[i])
1147 |
1148 | def SetUpSkipHashes(info):
1149 |
1150 |     if info.trueSkipFile:
1151 |         info.trueSkipHash = LoadTrueHash (info.trueSkipFile)
1152 |         WriteLog (info.logFile, "loaded %d items to skip from %s" %
      | (len (info.trueSkipHash.keys()), info.trueSkipFile))
1153 |     else:
1154 |         info.trueSkipHash = {}
1155 |     if info.doNotCountQuery:
1156 |         info.trueSkipHash[string.upper(info.doNotCountQuery)] = 1
1157 |     if info.falseSkipFile:
1158 |         info.falseSkipHash = LoadTrueHash (info.falseSkipFile)
1159 |         WriteLog (info.logFile, "loaded %d items to skip from %s" %
      | (len (info.falseSkipHash.keys()), info.trueSkipFile))
1160 |     else:

```

```

1161 |         info.falseSkipHash = {}
1162 |
1163 | def GenerateAlignment (info):
1164 |     import polacco.BlastXML
1165 |
1166 |     psiBlastFile = info.generateAlignment + ".psiblast.xml"
1167 |     info.alignFile = psiBlastFile + ".faln"
1168 |     print "Looking for " + info.alignFile
1169 |     if FileExists (info.alignFile):
1170 |         print "Found align file, not repeating psiblast."
1171 |         return
1172 |
1173 |     cmd = "%s -d %s -i %s -o %s -m7 -j2" % (__blastpggpPath,
1174 |     __blastDB, info.generateAlignment, psiBlastFile)
1175 |     print cmd
1176 |     os.system (cmd)
1177 |
1178 |     polacco.BlastXML.GetAlignmentFromPsiBlastFile (psiBlastFile,
1179 |     info.alignFile)
1180 |
1181 |     os.remove (psiBlastFile)
1182 |
1183 | # dummy class mostly to allow me to easily store any number of
1184 | # configuration variables
1185 | class struct:
1186 |     pass
1187 |
1188 | def main():
1189 |     info= struct()
1190 |     #display the next line unwrapped for an easy mapping from short
1191 |     #option to long option, or simply do GASPS.py -h
1192 |     shortList = ['h', 'p:', 'H:', 'r:',
1193 |     'n:', 'a:', 'A:', 'F:',
1194 |     'c:', 'o:', 'R:', 'l:',
1195 |     't:', 'i:', 'P:', 'T', 'O:', 'm:'
1196 |     , 's:', 'C:', 'N:', 'G:'
1197 |     , 'U:', 'X:', 'v:', 'b:'
1198 |     , 'S:', 'D:', 'M:', 'z:'
1199 |     , 'k:', 'x'
1200 |     , '']
1201 |     longOptions = ['help', 'pdbFile=', 'chain=', 'numResidues=',
1202 |     'maxNeighborhood=', 'alignFile=', 'refRowName=', 'alignFormat=',
1203 |     'minConservation=', 'filesPath=', 'allowedResidues=',
1204 |     'logFile=', 'tpFile=', 'iterations=', 'popSize=', 'testing',
1205 |     'rocArea=', 'maxFalse=', 'numSeqsInMA=', 'numTopConservedResidues=', 'no
1206 |     tTpFile=', 'alignRange=', 'trueLibrary=', 'falseLibrary=',
1207 |     'validateDirectory=', 'lengthTrueLibrary=', 'sepScoreImportance=',
1208 |     'useDistanceRmsd=', 'motifs=', 'scratchPath=', 'trueSkipFile=',
1209 |     'falseSkipFile=', 'xValidate', 'xTrueSkipFile=', 'xFalseLibrary=',
1210 |     'maxResidues=', 'doNotCountQuery=', 'noSubs', 'writeTables=',
1211 |     'useFileNames', 'generateAlignment=', 'maxRMSD=', 'maxCADiff=',
1212 |     'maxSCDiff=', 'scOnly']
1213 |
1214 |     shortOptions = string.join (shortList, '')
1215 |     opts, args = getopt.getopt (sys.argv[1:], shortOptions,
1216 |     longOptions)

```



```

1194
1195 SetDefaults (info)
1196
1197 for o,a in opts:
1198     if o in ('-h', '--help'):
1199         PrintUsage (shortList, longOptions)
1200         sys.exit(0)
1201     elif o in ('-p', '--pdbFile'):
1202         info.pdbFile = a
1203     elif o in ('-H', '--chain'):
1204         if a == 'space':
1205             a = ' '
1206             info.chain = a
1207     elif o in ('-r', '--numResidues'):
1208         info.numResidues = int (a)
1209     elif o in ('-n', '--maxNeighborhood'):
1210         info.maxNeighborhood = float (a)
1211     elif o in ('-a', '--alignFile'):
1212         info.alignFile = a
1213     elif o in ('-A', '--refRowName'):
1214         info.refRowName = a
1215     elif o in ('-F', '--alignFormat'):
1216         info.alignFormat = a
1217     elif o in ('-c', '--minConservation'):
1218         info.minConservation = float (a)
1219     elif o in ('-o', '--filesPath'):
1220         info.filesPath = a
1221     elif o in ('-R', '--allowedResidues'):
1222         if a == "NOTBORING":
1223             info.allowedResidues = "GSTCMPDNEQKRHFYW"
1224         else:
1225             info.allowedResidues = string.upper (a)
1226     elif o in ('-l', '--logFile'):
1227         info.logFile = a
1228     elif o in ('-t', '--tpFile'):
1229         if a != 'none':
1230             info.tpFile = a
1231     elif o in ('-i', '--iterations'):
1232         info.numRounds = int(a)
1233     elif o in ('-P', '--popSize'):
1234         info.populationSize = int(a)
1235     elif o in ('-T', '--testing'):
1236         print "TESTING, no spasm runs will be attempted. Scores
1237         chosen randomly!!!!"
1238         info.testing = 1
1239     elif o in ('-L', '--library'):
1240         pass
1241     elif o in ('-O', '--rocArea'):
1242         info.rocArea = int (a)
1243     elif o in ('-m', '--maxFalse'):
1244         info.maxFalse = int (a)
1245     elif o in ('-s', '--numSeqsInMA'):
1246         info.numSeqsInAlign = int (a)
1247     elif o in ('-C', '--numTopConservedResidues'):
1248         info.numTopConservedResidues = float (a)
1249     elif o in ('-N', '--notTpFile'):
1250         info.notTpFile = a
1251     elif o in ('-G', '--alignRange'):

```

```

1251 |         info.alignRange = a
1252 |     elif o in ('-U', '--trueLibrary'):
1253 |         info.trueLibrary = a
1254 |     elif o in ('-X', '--falseLibrary'):
1255 |         info.falseLibrary = a
1256 |     elif o in ('-v', '--validateDirectory'):
1257 |         info.validateDirectory = a
1258 |     elif o in ('-b', '--lengthTrueLibrary'):
1259 |         info.lengthTrueLibrary = int(a)
1260 |     elif o in ('-S', '--sepScoreImportance'):
1261 |         info.sepScoreImportance = float (a)
1262 |     elif o in ('-D', '--useDistanceRmsd'):
1263 |         info.useDistanceRmsd = int (a)
1264 |     elif o in ('-k', '--trueSkipFile'):
1265 |         info.trueSkipFile = a
1266 |     elif o in ('--falseSkipFile',):
1267 |         info.falseSkipFile = a
1268 |     elif o in ('-x', '--xValidate'):
1269 |         info.xValidate = 1
1270 |     elif o == '--xTrueSkipFile':
1271 |         info.xTrueSkipFile = a
1272 |     elif o == '--xFalseLibrary':
1273 |         info.xFalseLibrary = a
1274 |     elif o == '--maxResidues' :
1275 |         info.maxResidues = int (a)
1276 |     elif o == '--doNotCountQuery':
1277 |         info.doNotCountQuery = a
1278 |     elif o in ('-M', '--motifs'):
1279 |         motifs = a.split(",")
1280 |         for mot in motifs:
1281 |             mot = mot.split("_")
1282 |             mot.sort()
1283 |             info.motifs.append (mot)
1284 |     elif o in ('-z', '--scratchPath'):
1285 |         info.scratchPath = a
1286 |         try:
1287 |             os.makedirs (a)
1288 |         except OSError, data:
1289 |             print ("Error (ignored) while generating scratchPath:
%s" % a)
1290 |             print (data.strerror)
1291 |
1292 |     elif o == '--noSubs':
1293 |         print "ATTENTION: No substitutions will be allowed."
1294 |         info.noSubs = 1
1295 |     elif o == '--writeTables':
1296 |         info.writeTables = int (a)
1297 |     elif o == '--useFileNames':
1298 |         info.useFileNames = 1
1299 |     elif o == '--generateAlignment':
1300 |         info.generateAlignment = a
1301 |
1302 |     elif o == '--maxRMSD':
1303 |         info.maxRMSD = float (a)
1304 |     elif o == '--maxCADiff':
1305 |         info.maxCADiff = float (a)
1306 |     elif o == '--maxSCDiff':
1307 |         info.maxSCDiff = float (a)

```

```

1308     elif o == '--scOnly':
1309         info.scOnly = 1
1310
1311     else:
1312         print "Unrecognized option: %s : %s, use '-h' for list of
possible arguments" % (o,a)
1313         sys.exit(0)
1314
1315     if info.motifs:
1316         print ( "Loaded %d motifs from input" % len (info.motifs))
1317
1318
1319     if not info.logFile:
1320         info.logFile = info.filesPath+"_log.txt"
1321
1322
1323     if info.tpFile:
1324         info.trueHash = LoadTrueHash (info.tpFile)
1325     else:
1326         info.trueHash = {}
1327     SetUpSkipHashes(info)
1328
1329     print " Settings from command line and defaults: "
1330     DescribeMembers (info, sys.stdout)
1331
1332     if info.generateAlignment:
1333         GenerateAlignment(info)
1334
1335     if info.doNotCountQuery:
1336         info.lengthTrueLibrary -= 1
1337
1338     if info.validateDirectory:
1339         print TestMotif (info.validateDirectory, info.trueLibrary,
info.falseLibrary)
1340     else:
1341
1342         DoGASpasm (info)
1343
1344 if __name__ == "__main__":
1345     main()
1346
1347

```

## ***polacco/BlastXML.py***

```

1 | import polacco.XML, polacco.MultiAlign
2 | import string, sys
3 |
4 |
5 | #
6 | #
7 | #
8 | #
9 | #
10 | #
11 |

```

```

12 |
13 | # Despite it's name, it works equally well with the XML output of
    | both blastall and blastpgp.
14 |
15 |
16 | class PsiBlastXMLFile (polacco.XML.XML_tree) :
17 |
18 |     def __init__ (self, openFileIn):
19 |         self.maxEValue = 0
20 |         self.minHitOverlapFraction = 0.0
21 |
22 |         polacco.XML.XML_tree.__init__(self, openFileIn)
23 |         self.queryLength = self.GetQueryLength()
24 |
25 |
26 |     def SetMaxEValue (self, maxEValue):
27 |         self.maxEValue = float(maxEValue)
28 |     def SetQueryLength (self, queryLength):
29 |         self.queryLength = int (queryLength)
30 |     def SetMinHitOverlapFraction (self, minHitOverlapFraction):
31 |         self.minHitOverlapFraction = float (minHitOverlapFraction)
32 |
33 |     def GetHits (self, iteration = -1):
34 |         hits = self.rootNode.subNodes['BlastOutput_iterations'][-
    | 1].subNodes['Iteration'][(iteration)].subNodes['Iteration_hits'][-
    | 1].subNodes['Hit']
35 |         return hits
36 |
37 |     def GetQueryLength (self):
38 |         return int (self.rootNode.subNodes['BlastOutput_query-len'][-
    | 1].value)
39 |
40 |     # returns tuples of id, accession, hitDef, evalue
41 |     def GetSimpleHits (self, iteration = -1):
42 |         simpleHits = []
43 |         hits = self.GetHits(iteration)
44 |
45 |         for hit in hits:
46 |             accession = hit.subNodes['Hit_accession'][-1].value
47 |             id = hit.subNodes['Hit_id'][-1].value
48 |             hitDef = hit.subNodes['Hit_def'][-1].value
49 |
50 |             #now pick best evalue from all hsps
51 |             hsps = hit.subNodes['Hit_hsps'][-1].subNodes['Hsp']
52 |             eValue = max ( [ hsp.subNodes['Hsp_evalue'][-1].value for
    | hsp in hsps] )
53 |             simpleHits.append ( (id, accession, hitDef, eValue) )
54 |
55 |         return simpleHits
56 |
57 |     def GetMultiAlignment (self, iteration = -1):
58 |         hits = self.GetHits (iteration)
59 |         multiAlign = None
60 |
61 |         queryName = self.rootNode.subNodes['BlastOutput_query-def'][-
    | 1].value
62 |
63 |         for hit in hits:

```

```

64 |         accession = hit.subNodes['Hit_accession'][-1].value
65 |         id = hit.subNodes['Hit_id'][-1].value
66 |         hitDef = hit.subNodes['Hit_def'][-1].value
67 |
68 |         hsps = hit.subNodes['Hit_hsps'][-1].subNodes['Hsp']
69 |         #I don't want to assume these are already sorted by evalute,
so sort them by evalute
70 |         tempToSort = [ (hsp.subNodes['Hsp_evalue'][-1].value, hsp)
for hsp in hsps]
71 |         tempToSort.sort()
72 |         tempToSort.reverse()
73 |         hsps = [row[1] for row in tempToSort]
74 |
75 |         #determine which hsps are worth keeping (IMHO)
76 |         #keep the most significant that do not overlap with any
others on either the query or match sequence
77 |         goodHsps = []
78 |         def _QueryOverLap (hsp1, hsp2):
79 |             if hsp1.subNodes['Hsp_query-to'][-1] <
hsp2.subNodes['Hsp_query-from'][-1]:
80 |                 return 0
81 |             elif hsp1.subNodes['Hsp_query-from'][-1] >
hsp2.subNodes['Hsp_query-to'][-1]:
82 |                 return 0
83 |             else:
84 |                 return 1
85 |         def _MatchOverLap (hsp1, hsp2):
86 |             if hsp1.subNodes['Hsp_hit-to'][-1] <
hsp2.subNodes['Hsp_hit-from'][-1]:
87 |                 return 0
88 |             elif hsp1.subNodes['Hsp_hit-from'][-1] >
hsp2.subNodes['Hsp_hit-to'][-1]:
89 |                 return 0
90 |             else:
91 |                 return 1
92 |
93 |         for hsp in hsps:
94 |             for goodHsp in goodHsps:
95 |                 if _QueryOverLap(hsp, goodHsp) and _MatchOverLap
(hsp, goodHsp):
96 |                     break
97 |                 else:
98 |                     goodHsps.append (hsp)
99 |
100 |         hsps = goodHsps
101 |
102 |         if len(hsps) > 1:
103 |             print "More than one hsp found and used for %s %s" %
(accession, id)
104 |             #print summaries of overlaps.
105 |             for hsp in hsps:
106 |                 sys.stdout.write ("query:")
107 |                 for i in range (0, hsp.subNodes['Hsp_query-from'][-
1], 5):
108 |                     sys.stdout.write (".")
109 |                     for i in range (hsp.subNodes['Hsp_query-from'][-1],
hsp.subNodes['Hsp_query-to'][-1], 5):
110 |                         sys.stdout.write ("Q")

```

```

111 |         sys.stdout.write ('\n')
112 |
113 |         sys.stdout.write ("hit  :")
114 |         for i in range (0, hsp.subNodes['Hsp_hit-from'][-1],
115 | 5):
116 |             sys.stdout.write (".")
117 |             for i in range (hsp.subNodes['Hsp_hit-from'][-1],
118 | hsp.subNodes['Hsp_hit-to'][-1], 5):
119 |                 sys.stdout.write ("H")
120 |                 sys.stdout.write ('\n')
121 |
122 |     #hsp = hsps[0]
123 |     i = 0
124 |     for hsp in hsps:
125 |         #make sure e value is significant
126 |         if float(hsp.subNodes['Hsp_evalue'][-1].value) >
127 | self.maxEValue:
128 |             continue
129 |         #make sure we are aligning to a significant fraction of
130 | the query
131 |         alignQueryLength = int(hsp.subNodes['Hsp_query-to'][-
132 | 1].value) - int(hsp.subNodes['Hsp_query-from'][-1].value)
133 |         if float(alignQueryLength)/self.queryLength <
134 | self.minHitOverlapFraction:
135 |             continue
136 |
137 |         hitName = accession
138 |         if i > 0:
139 |             hitName = hitName + ".%d" % i
140 |             i+=1
141 |
142 |         if not multiAlign:
143 |             parentRow = polacco.MultiAlign.AlignRow ( queryName ,
144 | int(hsp.subNodes['Hsp_query-from'][-1].value), list
145 | (hsp.subNodes['Hsp_qseq'][-1].value))
146 |             multiAlign = polacco.MultiAlign.MultiAlign
147 | ([parentRow])
148 |
149 |             master = polacco.MultiAlign.AlignRow ( queryName ,
150 | int(hsp.subNodes['Hsp_query-from'][-1].value),
151 | list(hsp.subNodes['Hsp_qseq'][-1].value))
152 |             #slave = polacco.MultiAlign.AlignRow (hitDef[0:50],
153 | int(hsp.subNodes['Hsp_hit-from'][-1].value),
154 | list(hsp.subNodes['Hsp_hseq'][-1].value))
155 |             slave = polacco.MultiAlign.AlignRow (hitName,
156 | int(hsp.subNodes['Hsp_hit-from'][-1].value),
157 | list(hsp.subNodes['Hsp_hseq'][-1].value))
158 |
159 |             try:
160 |                 multiAlign.addPair (master, slave)
161 |             except "AlignmentOutOfRange":
162 |                 raise "AlignmentOutOfRange"
163 |
164 |         print "Success reading alignment from (psi)blast file!"
165 |         return multiAlign, parentRow
166 |
167 |
168 |
169 |

```

```

154 |
155 |
156 |
157 |
158 | def GetAlignmentFromPsiBlastFile (fileName, outFile_name = None):
159 |     fp = open (fileName)
160 |
161 |     pbFile = PsiBlastXMLFile (fp)
162 |     fp.close()
163 |     if not outFile_name:
164 |         outFile_name = fileName + ".faln"
165 |
166 |     pbFile.SetMaxEValue (1.0e-10)
167 |     #pbFile.SetQueryLength (355)
168 |     pbFile.SetMinHitOverlapFraction (0.5)
169 |
170 |     ma,parentRow = pbFile.GetMultiAlignment()
171 |     del (pbFile)
172 |
173 |     ma.protectedRows.append(parentRow)
174 |     ma.ShrinkByRemovingRedundancy_Efficient(50)
175 |     ma.RemoveGappedColumns()
176 |     ma.DashifyGapCharacters()
177 |     fpout = open (outFile_name, "w")
178 |     ma.simple_print (fpout)
179 |     fpout.close()
180 |     return outFile_name
181 |
182 |
183 | if __name__ == '__main__':
184 |     GetAlignmentFromPsiBlastFile (sys.argv[1])
185 |
186 |
187 |
188 |
189 |
190 | def test():
191 |     fp = open ("longtest.xml")
192 |     pbtree = PsiBlastXMLFile (fp)
193 |
194 |     ma = pbtree.GetMultiAlignment()


```

## ***polacco/Data.py***

```

1 | #####
2 | #
3 | #
4 | #
5 | #
6 | #
7 | #
8 | #
9 | #
10 | #
11 | #

```



```

12 | #       Things were done to save typing time, not necessarily program
    |       running time
13 | #####
14 | import string
15 |
16 | global PET91_matrix
17 | PET91_matrix = None
18 |
19 | def GetPET91_matrix():
20 |     global PET91_matrix
21 |     if PET91_matrix:
22 |         return PET91_matrix
23 |
24 |     aaOrder = "ARNDCQEGHILKMFPSTWYV"
25 |     PET91_matrix = {}
26 |
27 |     # PET91 matrix for 120 PAM (Jones, Thornton and Taylor)
28 |     temp = {}
29 |     temp["A"] = string.split ("  6 -3 -1 -2 -3 -3 -2  0 -4
    | -1 -4 -4 -2 -6  0  2  2 -7 -7  1")
30 |     temp["R"] = string.split (" -3  8 -2 -4 -1  2 -3 -1  2
    | -6 -5  4 -4 -8 -2 -2 -3 -1 -5 -6")
31 |     temp["N"] = string.split (" -1 -2  8  3 -3 -1 -1 -1  2
    | -4 -6  1 -4 -6 -3  2  1 -8 -2 -5")
32 |     temp["D"] = string.split (" -2 -4  3  8 -7 -1  5  0 -1
    | -7 -8 -2 -7 -9 -5 -2 -3 -10 -4 -5")
33 |     temp["C"] = string.split (" -3 -2 -3 -7 14 -6 -8 -3 -2
    | -5 -5 -6 -5 -2 -5  0 -3  0  2 -3")
34 |     temp["Q"] = string.split (" -3  2 -1 -1 -6  9  2 -4  4
    | -6 -3  2 -4 -7  0 -3 -3 -6 -4 -6")
35 |     temp["E"] = string.split (" -2 -3 -1  5 -8  2  8 -1 -3
    | -7 -8  0 -6 -10 -5 -3 -4 -8 -7 -5")
36 |     temp["G"] = string.split ("  0 -1 -1  0 -3 -4 -1  8 -4
    | -6 -8 -4 -6 -9 -4  0 -3 -3 -8 -4")
37 |     temp["H"] = string.split (" -4  2  2 -1 -2  4 -3 -4 11
    | -6 -4 -1 -5 -2 -1 -2 -3 -6  4 -6")
38 |     temp["I"] = string.split (" -1 -6 -4 -7 -5 -6 -7 -6 -6
    |  7  2 -6  3 -1 -5 -3  0 -7 -5  5")
39 |     temp["L"] = string.split (" -4 -5 -6 -8 -5 -3 -8 -8 -4
    |  2  7 -6  3  2 -2 -3 -4 -4 -4  1")
40 |     temp["K"] = string.split (" -4  4  1 -2 -6  2  0 -4 -1
    | -6 -6  8 -4 -10 -4 -3 -2 -6 -7 -6")
41 |     temp["M"] = string.split (" -2 -4 -5 -7 -5 -4 -6 -6 -5
    |  3  3 -4 10 -2 -4 -3  0 -6 -6  2")
42 |     temp["F"] = string.split (" -6 -8 -6 -9 -2 -7 -10 -9 -2
    | -1  2 -9 -2 11 -5 -3 -6 -3  5 -2")
43 |     temp["P"] = string.split ("  0 -2 -3 -5 -5  0 -5 -4 -1
    | -5 -2 -4 -5 -5  9  1  0 -8 -6 -4")
44 |     temp["S"] = string.split ("  2 -2  2 -2  0 -3 -3  0 -2
    | -3 -3 -3 -3 -3  1  5  2 -5 -3 -3")
45 |     temp["T"] = string.split ("  2 -3  1 -3 -3 -3 -4 -2 -3
    |  0 -4 -2  0 -6  0  2  6 -7 -6  0")
46 |     temp["W"] = string.split (" -8 -1 -8 -9  0 -6 -9 -3 -6
    | -7 -4 -6 -6 -3 -8 -5 -7 17 -2 -6")
47 |     temp["Y"] = string.split (" -7 -5 -2 -4  2 -4 -7 -8  4
    | -5 -4 -7 -6  5 -7 -3 -6 -2 12 -6")
48 |     temp["V"] = string.split ("  1 -6 -5 -5 -3 -6 -5 -4 -6
    |  5  1 -6  2 -2 -4 -3  0 -6 -6  7")

```



```

49 | #temp["B"] = string.split (" 0 0 0 0 0 0 0 0 0 0 0 0 0
0 | 0 0 0 0 0 0 0 0 0 0 0 0")
50 |
51 | for aaRow in aaOrder:
52 |     PET91_matrix[aaRow] = {}
53 |     for i in range (len (aaOrder)):
54 |         PET91_matrix[aaRow][aaOrder[i]] = int(temp[aaRow][i])
55 |
56 | #now do specials
57 |
58 |
59 | return PET91_matrix
60 |
61 |

```

### ***polacco/MultiAlign.py***

```

1 | #! /sw/bin/python
2 |
3 |
4 | #
5 | #
6 | #
7 | #
8 | #
9 | #
10 |
11 |
12 |
13 | from string import *
14 | import sys, copy, math, string
15 |
16 |
17 |
18 | def is_gap(char):
19 |     if (char == '-' or char == '.' or char == '?'):
20 |         return 1
21 |     else:
22 |         return 0
23 |
24 | def is_ambiguous (char):
25 |     if (char in 'XBUxbu'):
26 |         return 1
27 |     else:
28 |         return 0
29 |
30 |
31 | class MultiAlign:
32 |     def __init__(self, rows = None):
33 |         if not rows:
34 |             self.rows = []
35 |         else:
36 |             self.rows = rows
37 |         self.FIRST_POSITION = 1 #how should the API call the first
position (0 or 1)
38 |         self.protectedRows = [] #these are rows that are of special
interest, see ShrinkByRemovingRedundancy_Efficient

```

```

39 |
40 | def LoadFromFastaFile (self, fastaFile):
41 |     fp = open (fastaFile, "r")
42 |     self.read_fasta (fp)
43 |     fp.close()
44 |
45 |
46 | def read_fasta(self, fasta_file):
47 |     if len (self.rows) > 0:
48 |         print "WARNING: Blindly adding new fasta alignment to
current alignment!"
49 |
50 |         name = ''
51 |         lines = []
52 |
53 |         for line in fasta_file.readlines():
54 |             if line[0] in (">"):
55 |                 #we're done with previous row
56 |                 if name:
57 |                     #first join lines, then split and join to remove
spaces, then repeat to remove \n, then make list
58 |                     chars = list (join(split(join (split(join (lines,
'', " "), ''), "\n"), ""))
59 |                     self.rows.append(AlignRow ( name, offset, chars))
60 |                     lines = []
61 |                     slash_split = split (line[1:-1], '/')
62 |                     if (len (slash_split) > 1):
63 |                         name = join(slash_split[:-1], '')
64 |                         try:
65 |                             offset = int (split (slash_split[-1], '-')[0])
66 |                         except ValueError:
67 |                             offset = 1
68 |
69 |                     else:
70 |                         offset = 1
71 |                         name = slash_split[0]
72 |                     else:
73 |                         lines.append(line)
74 |                 else:
75 |                     #if everything went well we have to add the last sequence.
76 |                     if lines:
77 |                         chars = list (join(split(join (split(join (lines, ''), "
"), ''), "\n"), ""))
78 |                         self.rows.append(AlignRow ( name, offset, chars))
79 |
80 |
81 | def LoadFromClustalFile (self, clustalFile, strictCLUSTAL = 1):
82 |     print "Loading alignment from %s" % clustalFile
83 |     fp = open (clustalFile, 'r')
84 |     self.read_clustal (fp, strictCLUSTAL)
85 |     fp.close()
86 |
87 |
88 |
89 | def read_clustal (self, clustal_file, strictCLUSTAL = 1):
90 |     if len(self.rows) > 0:
91 |         print "WARNING: Blindly adding new fasta alignment to
current alignment! Current row names:"

```

```

92 |         print [row.name for row in self.rows]
93 |     partsHash = {}
94 |     inorder = []
95 |     fileStarted = 0
96 |     while (1):
97 |         line = clustal_file.readline()
98 |         if line == '':
99 |             break
100 |         words = line.split()
101 |         if (len (words) == 0):
102 |             continue
103 |
104 |         elif (not fileStarted):
105 |             if words[0] == "CLUSTAL":
106 |                 fileStarted = 1
107 |                 continue
108 |             if strictCLUSTAL:
109 |                 continue
110 |             else:
111 |                 fileStarted = 1
112 |
113 |         elif (words[0][0] in ".*"):
114 |             if len (partsHash.keys()) == 0:
115 |                 print "WARNING: unexpected location of conservation
line, possibly illegal first character for sequence name?"
116 |                 print line
117 |                 #looks like a line indicating conservation, skip it
118 |                 continue
119 |
120 |         try:
121 |             partsHash[words[0]].append (string.join (words[1:], ''))
122 |         except KeyError:
123 |             inorder.append (words[0])
124 |             partsHash[words[0]] = [string.join (words[1:], '')]
125 |
126 |         length = -1
127 |         #error checking: make sure file was started.  and make sure we
have same number of parts for all
128 |         for key in inorder:
129 |             parts = partsHash[key]
130 |             if length < 0:
131 |                 length = len (parts)
132 |             elif len (parts) != length:
133 |                 print "Current: " + key
134 |                 print "Parts mismatch: %d with %d" % (len(parts),
length)
135 |             print inorder
136 |             print parts
137 |             raise "Bad Alignment Read!"
138 |         chars = list (string.join(parts, ""))
139 |         #see if we can get the offset from the name (key)
140 |         slash_split = split (key, '/')
141 |         if (len (slash_split) > 1):
142 |             name = join(slash_split[:-1], '')
143 |             try:
144 |                 offset = int (split (slash_split[-1], '-')[0])
145 |             except ValueError:
146 |                 offset = 1

```

```

147 |         else:
148 |             offset = 1
149 |             name = slash_split[0]
150 |             self.rows.append(AlignRow ( name, offset, chars))
151 |
152 |
153 | # very simple FASTA format
154 | def simple_print(self, outfile):
155 |     for row in self.rows:
156 |         outfile.write(">%s\n"%(row.description() ))
157 |         outfile.write("%s\n" % join(row.chars, ''))
158 |
159 | # roughly clustal format
160 | def PrettyPrint (self, outfile, columns = 60, nameWidth = 20):
161 |     minLength = self.min_length()
162 |     nameFormatString = '%%-%ds' % nameWidth
163 |
164 |     i = 0
165 |     while (i < minLength):
166 |         end = i+ columns
167 |         for row in self.rows:
168 |             outfile.write (nameFormatString %
169 | row.name.split()[0][0:nameWidth])
170 |             outfile.write (" %s\n" % string.join
171 | (row.chars[i:end], ''))
172 |             i = end
173 |             if i >= minLength:
174 |                 break
175 |             outfile.write ("\n\n")
176 |
177 | def get_row_by_name (self, name):
178 |     for row in self.rows:
179 |         if row.name == name:
180 |             return row
181 |     else:
182 |         print "row not found among: "
183 |         print map((lambda row: row.name), self.rows)
184 |         return 0
185 |
186 | def remove_row_by_name(self, name):
187 |     for row in self.rows:
188 |         if row.name == name:
189 |             self.rows.remove(row)
190 |             break
191 |     else:
192 |         raise ValueError
193 |
194 | def get_column (self, position):
195 |     column = {}
196 |     if position >= self.FIRST_POSITION:
197 |         for row in self.rows:
198 |             column[row] = (row.chars[position-self.FIRST_POSITION])
199 |     elif position < 0:
200 |         for row in self.rows:
201 |             column[row] = (row.chars[position])
202 |     else:

```

```

203         raise "Illegal MultiAlign Position"
204
205     return column
206
207 def get_column_aslist (self, position):
208     column = []
209     if position >= self.FIRST_POSITION:
210         for row in self.rows:
211             column.append(row.chars[position-self.FIRST_POSITION])
212     elif position < 0:
213         for row in self.rows:
214             column.append(row.chars[position])
215     else:
216         raise "Illegal MultiAlign Position"
217
218     return column
219
220
221 def insert_gap (self, position):
222     self.insertColumnIndex ( position)
223
224     if position == -1:
225         for row in self.rows:
226             if row.chars[-1] == '.':
227                 row.chars.append ('.')
228             else:
229                 row.chars.append('-')
230
231     else:
232         position = position - self.FIRST_POSITION
233         for row in self.rows:
234             if position == 0 or row.chars[position - 1] == '.':
235                 row.chars[position:position] = ['.']
236             else:
237                 row.chars[position:position] = ['-']
238
239 def max_length (self):
240     if len (self.rows) == 0:
241         return 0
242     return max (map (lambda r : len (r.chars), self.rows))
243
244 def min_length (self):
245     if len (self.rows) == 0:
246         return 0
247     return min (map (lambda r : len (r.chars), self.rows))
248
249 #only use if you KNOW that minlength and maxlength are the same
250 def length (self):
251     minim = self.min_length()
252     maxim = self.max_length()
253     assert (minim == maxim)
254     return maxim
255
256
257 def getMultiAlignBlock (self, start_column, stop_column):
258     subrows = []
259     for row in self.rows:

```

```

260 |         subrows.append (row.subRow (start_column -
self.FIRST_POSITION,
261 |                                     stop_column - self.FIRST_POSITION
+1)) #+1 includes the stop_column
262 |
263 |         return MultiAlign (subrows)
264 |
265 |     def MAPositionFromSeqPosition (self, refRow, rowPosition):
266 |         numChars = refRow.num_chars_at_offset (rowPosition )
267 |
268 |         return numChars + self.FIRST_POSITION
269 |
270 |
271 |     def ComputePIDTable (self):
272 |         pidTable = {}
273 |         for seqRow in self.rows:
274 |             pidTable[seqRow] = {}
275 |             for seqCol in self.rows:
276 |                 try:
277 |                     pidTable[seqRow][seqCol] = pidTable[seqCol][seqRow]
278 |                 except KeyError:
279 |                     pidTable[seqRow][seqCol] = seqRow.PercentIdentity
(seqCol)
280 |         return pidTable
281 |
282 |     def ComputePIDList (self, rows = None):
283 |         if not rows:
284 |             rows = self.rows
285 |         pidList = []
286 |         done = []
287 |         for seq1 in rows:
288 |             done.append(seq1)
289 |             for seq2 in rows:
290 |                 if seq2 in done:
291 |                     continue
292 |                 pidList.append ([seq1.PercentIdentity (seq2), seq1,
seq2])
293 |         return pidList
294 |
295 |     def GetHighestPIDPair (self, pidList):
296 |         max = 0.0
297 |         topPair = None
298 |
299 |         for pair in pidList:
300 |             if (pair[1] in self.protectedRows and pair[2] in
self.protectedRows):
301 |                 continue
302 |                 if pair[0] > max:
303 |                     max = pair[0]
304 |                     topPair = pair
305 |         if topPair == None:
306 |             print "Warning: All pairs are protected!"
307 |
308 |         return topPair
309 |
310 |     # At least it's more efficient than my last one which is now
deleted. No other
311 |     # claims are made to its efficiency!

```

```

312 |     def ShrinkByRemovingRedundancy_Efficient (self, maxSeqs):
313 |         if maxSeqs >= len (self.rows):
314 |             print "No shrinking needed (%d, %d)" % (maxSeqs, len
(self.rows))
315 |             return
316 |             #compute initial pidlist
317 |             #print "Computing initial PIDList"
318 |             keepRows = []
319 |             for row in self.protectedRows:
320 |                 if row in self.rows and not row in keepRows:
321 |                     keepRows.append (row)
322 |
323 |
324 |             i = 0
325 |             for row in self.rows:
326 |                 if len (keepRows) >= maxSeqs:
327 |                     break
328 |                 if not row in keepRows:
329 |                     keepRows.append (row)
330 |                 i = i + 1
331 |
332 |
333 |             pidList = self.ComputePIDList(keepRows)
334 |
335 |             #find highest identity, and rows that give it
336 |             topPidPair = self.GetHighestPIDPair (pidList)
337 |             #print topPidPair
338 |
339 |             #for each row not yet in keep, find if it is more similar than
highest identity in keep to any
340 |             #if it is, then skip it
341 |             #if it is not then add it, and its pidList
342 |             #then remove one or the other row of highest scoring pair
343 |             removed = [] # we have to keep a running tally of these to
remove when we're don, because we can't change self.rows while we're
iterating over it
344 |             for candidate in self.rows[i:]:
345 |                 if candidate in self.protectedRows:
346 |                     continue
347 |                 candPidList = []
348 |
349 |                 for row in keepRows:
350 |                     pid = row.PercentIdentity (candidate)
351 |                     candPidList.append ( [pid,candidate, row] )
352 |                     if pid >= topPidPair[0]:
353 |                         #print "Removing row: %s because it is %6.4f
identical to %s compared to running high of %6.4f (%s and %s) " %
(candidate.name, pid, row.name, topPidPair[0], topPidPair[1].name,
topPidPair[2].name)
354 |                         removed.append (candidate)
355 |                         break
356 |                 else: #it was not more similar than topPercentIdentity to
anything else, then we need to add it
357 |                     #first remove other
358 |                     length = [0,0,0]
359 |
360 |                     length[1] = topPidPair[1].numCharsNoGaps()
361 |                     length[2] = topPidPair[2].numCharsNoGaps()

```

```

362 |
363 |         #they can't both be in protectedRows because the
GetHighestPIDPair function prevents it
364 |         if topPidPair[1] in self.protectedRows:
365 |             keep = 1
366 |         elif topPidPair[2] in self.protectedRows:
367 |             keep = 2
368 |         elif length[1] >= length[2] :
369 |             keep = 1
370 |         else:
371 |             keep = 2
372 |
373 |         if keep == 1:
374 |             rem = 2
375 |         else:
376 |             rem = 1
377 |
378 |
379 |         removed.append (topPidPair[rem])
380 |         #print "Removing row: %s (%d chars) to keep %s with (%d
chars); similarity = %6.4f" % (topPidPair[rem].name, length[rem],
topPidPair[keep].name,length[keep], topPidPair[0])
381 |
382 |         #now combine pidList and remove old entries
383 |         pidList = pidList + candPidList
384 |         toRemove = []
385 |         for pair in pidList:
386 |             if pair[1] == topPidPair[rem] or pair[2] ==
topPidPair[rem]:
387 |                 toRemove.append(pair)
388 |             for r in toRemove:
389 |                 pidList.remove (r)
390 |
391 |         #now update keepRows
392 |         keepRows.remove (topPidPair[rem])
393 |         keepRows.append (candidate)
394 |
395 |         #find new highest identity, and rows that give it
396 |         topPidPair = self.GetHighestPIDPair (pidList)
397 |
398 |
399 |         print "Removing %d rows to keep %d rows that have at most
%6.4f sequence identity (%s)" % (len (removed), len(keepRows),
topPidPair[0], self.rows[0].name)
400 |         #now do actual removing from self
401 |         #print [r.name for r in removed]
402 |         for r in removed:
403 |             #print "Removing %s" % r.name
404 |             self.rows.remove (r)
405 |
406 |         #done?
407 |
408 |     def RemoveGappedColumns(self):
409 |         toRemove = []
410 |
411 |         for i in range (self.length()):
412 |             for j in range (len (self.rows)):
413 |                 if not is_gap (self.rows[j].chars[i]):

```



```

414         break
415     else:
416         toRemove.append(i)
417
418     toRemove.reverse()
419     for i in toRemove:
420         self.RemoveColumnBySimpleIndex (i)
421
422 def RemoveColumnBySimpleIndex (self, index):
423     for i in range (len (self.rows)):
424         self.rows[i].chars[index:index+1] = []
425
426 def DashifyGapCharacters (self):
427     for row in self.rows:
428         row.DashifyGapCharacters()
429
430 def GetDominantLettersPerColumn (self, minFraction = 0.1):
431     seqLetters = []
432     numRows = len (self.rows)
433     for i in range (self.length()):
434         columnLetters = {}
435         domLetters = []
436         for row in self.rows:
437             if is_gap (row.chars[i]):
438                 continue
439             try:
440                 columnLetters[row.chars[i]]+=1.0
441             except KeyError:
442                 columnLetters[row.chars[i]] = 1.0
443         for char in columnLetters.keys():
444             if columnLetters[char]/numRows > minFraction:
445                 domLetters.append (char)
446         seqLetters.append(domLetters)
447     return seqLetters
448 def addPair (self, master, slave):
449     return self.addMultiAlign (MultiAlign ([master, slave]))
450
451 #####
452 # The next four functions deal with adding or aligning two multiple
453 # alignments to each other. These alignments must share one row
454 # with the same name and with an identical overlapping region. No
455 # mismatches are allowed. anchor refers to that row in the self
456 # alignment
457 # tether refers to that row in the otherAlign.
458 #####
459 # Only generates a map of columns in one align to columns in the
460 # other
461 # No alignments (except temporary copies) should be modified
462 def MapColumns2OtherAlign (self, otherAlign):
463     (tether, anchor) = self.FindTetherAndAnchor(otherAlign)
464
465     #do all work on a copy of anchor and tether
466     anchor = AlignRow (anchor.name, anchor.offset, anchor.chars)
467     tether = AlignRow (tether.name, tether.offset, tether.chars)

```

```

468 |         #set up dummy clones of MultiAligns to keep track of
      | columnIndeces
469 |         otherAlign2 = MultiAlign ([tether])
470 |         otherAlign2.index = copy.deepcopy(otherAlign.columnIndex)
471 |         self2 = MultiAlign ([anchor])
472 |         self2.index = copy.deepcopy (self.columnIndex)
473 |
474 |         map = self.align2Aligns (self2, otherAlign2, tether, anchor)
475 |
476 |
477 |         return (map, self2)
478 |
479 | # This does the work of inserting gaps and extensions necessary so
      | that 2 alignments
480 | # can be directly compared or added by simply stacking
481 |
482 |     def align2Aligns (__self, firstAlign, otherAlign, tether,
      | anchor):
483 |         #print "Before dealing with offsets"
484 |         #pretty_print ( [join(anchor.chars,''), join(tether.chars,
      | '' ) ], 60)
485 |
486 |         #we're about to rely on offset values (a bad idea), so fix the
      | anchor's offset
487 |         #print anchor.offset, tether.offset
488 |         anchor.fixOffset (tether)
489 |         #print anchor.offset, tether.offset
490 |         #first deal with cases where the left end of tether begins
      | after anchor
491 |         if tether.offset > anchor.offset:
492 |             split = anchor.num_chars_at_offset (tether.offset)
493 |             if split > len (anchor.chars):
494 |                 sys.stderr.write ("alignment out of range, shares two
      | non overlapping regions of %s\n" % (anchor.name))
495 |                 raise "AlignmentOutOfRange"
496 |                 i = split
497 |                 while i:
498 |                     otherAlign.insert_column(otherAlign.FIRST_POSITION, '.')
499 |                     i = i - 1
500 |                 tether.chars[0:split] = anchor.chars[0:split]
501 |                 tether.offset = anchor.offset
502 |
503 |         #now deal with cases where the beginning of row is to the
      | right of beginning of master
504 |         if anchor.offset > tether.offset:
505 |             split = tether.num_chars_at_offset (anchor.offset)
506 |             if split > len (tether.chars):
507 |                 sys.stderr.write ("alignment out of range, shares two
      | non overlapping regions of %s\n" % (anchor.name))
508 |                 raise "AlignmentOutOfRange"
509 |                 i = split
510 |                 while i:
511 |                     firstAlign.insert_column (firstAlign.FIRST_POSITION,
      | '.')
512 |                     i = i - 1
513 |
514 |                 anchor.chars[0:split] = tether.chars[0:split]
515 |                 anchor.offset = tether.offset

```

```

516 |
517 |     #print "After dealing with offsets..."
518 |     #pretty_print ( [join(anchor.chars,''), join(tether.chars,
519 | '' ) ], 60)
520 |
521 |     i=0
522 |     minlength = min (len (anchor.chars), len (tether.chars))
523 |     while (i < minlength):
524 |         empty = ('.-')
525 |         #print "mas: %s row: %s" % (master.chars[i], row.chars[i])
526 |         if (anchor.chars[i] == tether.chars[i]):
527 |             # print "=",
528 |                 i = i + 1
529 |                 continue
530 |         elif (anchor.chars[i] in empty) and (tether.chars[i] in
531 | empty):
532 |             # print "g",
533 |                 tether.chars[i] = anchor.chars[i] = '-'
534 |                 i = i + 1
535 |                 continue
536 |         elif (tether.chars[i] in empty) and (anchor.chars[i] not in
537 | empty):
538 |             # print "^",
539 |                 firstAlign.insert_column (firstAlign.FIRST_POSITION + i,
540 | tether.chars[i])
541 |             elif (anchor.chars[i] in empty) and (tether.chars[i] not in
542 | empty):
543 |                 # print "v",
544 |                 otherAlign.insert_column (otherAlign.FIRST_POSITION + i,
545 | anchor.chars[i])
546 |             elif (tether.chars[i] == 'X'):
547 |                 # print "x",
548 |                 tether.chars[i] = anchor.chars[i]
549 |             elif (anchor.chars[i] == 'X'):
550 |                 # print "X",
551 |                 anchor.chars[i] = tether.chars[i]
552 |             else: # (row.chars[i] != master.chars[i]):
553 |                 pretty_print ( [join(anchor.chars,''),
554 | join(tether.chars, '' ) ], 60)
555 |                 print "WARNING!: tether sequence in foreign alignment
556 | does not match corresponding anchor sequence in firstAlign
557 | alignment! (%s %s)" % (anchor.chars[i], tether.chars[i])
558 |                 firstAlign.simple_print(sys.stdout)
559 |                 print "\notherAlign:\n"
560 |                 otherAlign.simple_print(sys.stdout)
561 |                 raise "exception"
562 |                 i = i + 1
563 |                 minlength = min (len (anchor.chars), len (tether.chars))
564 |
565 |     #pretty_print ( [join(anchor.chars,''), join(tether.chars,
566 | '' ) ], 60)
567 |
568 |     #now deal with the trailing end
569 |     diff = len(anchor.chars) - len (tether.chars)

```

```

564 |         if diff > 0:
565 |             for c in anchor.chars[-diff:]:
566 |                 otherAlign.insert_column (-1, '.')
567 |                 tether.chars[-1] = c
568 |         elif diff < 0:
569 |             for c in tether.chars[diff:]:
570 |                 firstAlign.insert_column (-1, '.')
571 |                 anchor.chars[-1] = c
572 |             #pretty_print ( [join(row.chars,''), join(master.chars,
573 | ''), join(slave.chars, '') ], 60)
574 |             #pretty_print ( [join(anchor.chars,''), join(tether.chars,
575 | ''), ], 60)
576 |             assert (len(anchor.chars) == len(tether.chars)),
577 | (len(anchor.chars) ,len (tether.chars))
578 |             assert (join(anchor.chars) == join (tether.chars)), "\n%s\n%s"
579 | % (join(anchor.chars, ''), join(tether.chars, ''))
580 |         #this does the work of finding the actual rows that are the
581 | tether and anchor. It chooses
582 | # the first it finds, not necessarily the best.
583 | def FindTetherAndAnchor (self, otherAlign, taName = None):
584 |     otherNames = [row.name for row in otherAlign.rows]
585 |     anchor = 0
586 |     for row in self.rows:
587 |         if taName and taName != row.name:
588 |             continue
589 |         if row.name in otherNames:
590 |             anchor = row
591 |             for r in otherAlign.rows:
592 |                 if r.name == anchor.name:
593 |                     tether = r
594 |                     break
595 |             break
596 |         else:
597 |             #corresponding master row not found, give warning
598 |             raise "Error finding an anchor and tether %s" %
599 | (otherNames), [row.name for row in self.rows]
600 |     return (tether, anchor)
601 |
602 | # This will add otherAlign to the self alignment.
603 | def insert_column(self, position, char):
604 |     if (len(char) > 1):
605 |         raise "insert too long", char
606 |     if position == -1:
607 |         for row in self.rows:
608 |             row.chars.append(char)
609 |     else:
610 |         position = position - self.FIRST_POSITION
611 |         for row in self.rows:
612 |             row.chars[position:position] = [char]
613 |
614 | def addMultiAlign (self, otherAlign, taName = None):
615 |     (tether, anchor) = self.FindTetherAndAnchor (otherAlign,
616 | taName)
617 |     #print "Using %s (%s) as tether and anchor sequences" %
618 | (tether.name, anchor.name)

```

```

614 |         self.align2Aligns (self, otherAlign, tether, anchor)
615 |         #otherAlign.rows.remove(tether)
616 |         self.rows.extend(otherAlign.rows)
617 |         self.rows.remove (tether)
618 |
619 |     #
620 |     #-----
621 |
622 |     #-----
623 |     #
624 |     # used by MultiAlign
625 |
626 |     class AlignRow:
627 |         def __init__(self, name, offset, chars):
628 |             self.name = name
629 |             self.offset = offset
630 |             self.chars = chars
631 |
632 |         def fixOffset (self, referenceRow):
633 |             minMatch = 20 # the reference row must minimally overlap by
634 |             this many characters.
635 |             # the larger it is the faster it is with obvious
636 |             problems if its too large
637 |             selfChars = self.GetCharsNoGaps()
638 |             reference = referenceRow.GetCharsNoGaps()
639 |
640 |             correction = minMatch - len (selfChars)
641 |             sub1Start = -correction
642 |             sub1End = len (selfChars)
643 |             sub2Start = 0
644 |             sub2End = minMatch
645 |
646 |             # find alignment by looking for matching substrings
647 |             # "slide" one sequence across the other and look for matching
648 |             substrings where substrings
649 |             # are the overlapping regions.
650 |             while (correction < len(reference) - minMatch):
651 |                 sub1 = selfChars[sub1Start:sub1End]
652 |                 sub2 = reference[sub2Start:sub2End]
653 |                 if sub1 == sub2:
654 |                     #print correction
655 |                     #print sub1
656 |                     #print sub2
657 |                     break
658 |
659 |                 #update for next iteration
660 |                 if sub1Start > 0:
661 |                     sub1Start = sub1Start - 1
662 |                 else:
663 |                     sub2Start = sub2Start + 1
664 |                 if sub1End + correction == len(reference):
665 |                     sub1End -= 1
666 |                 else:
667 |                     sub2End = sub2End + 1
668 |                     correction = correction + 1
669 |
670 |             self.offset = self.offset + correction - (self.offset -
671 |             referenceRow.offset)

```

```

668
669
670
671 def GetCharsNoGaps (self):
672     ch = []
673     for char in self.chars:
674         if not is_gap(char):
675             ch.append (char)
676     return ch
677
678 def GetCharsAndIndexesNoGaps (self):
679     ch = []
680     indexes = []
681     for i in range (len (self.chars)):
682         if not is_gap(self.chars[i]):
683             ch.append (self.chars[i])
684             indexes.append (i)
685     return ch, indexes
686
687 def num_chars_at_offset (self, offset):
688     num=0
689     if offset < self.offset and offset > 0:
690         return offset-self.offset
691     offset = offset - self.offset
692
693     for char in self.chars:
694         if not is_gap(char):
695             if offset == 0:
696                 return num
697             offset = offset - 1
698             num+=1
699     else:
700         return len (self.chars) + 1
701
702 def numCharsNoGaps(self):
703     count = 0
704     for char in self.chars:
705         if not is_gap(char):
706             count = count + 1
707     return count
708
709
710 def initial_gap_chars(self):
711     count = 0
712     for char in self.chars:
713         if is_gap(char):
714             count+=1
715         else:
716             break
717     return count
718
719 def terminal_gap_chars(self):
720     count = 0
721     pointer = -1
722     try:
723         while (1):
724             if is_gap(self.chars[pointer]):
725                 count = count + 1

```

```

726         pointer = pointer - 1
727     else:
728         break
729     except IndexError:
730         pass
731     return count
732
733     def description(self):
734         title = split (self.name, "\n")[0]
735         title = split (title, " ")[0]
736         desc = "%s/%03d" % (title, self.offset)
737         return desc
738
739     def subRow (self, start_index, stop_index):
740         return AlignRow (self.name, self.offset + start_index,
self.chars[start_index:stop_index])
741
742
743     def PercentIdentity (self, other):
744         seq1 = self.chars
745         seq2 = other.chars
746         numAligned = 0.0
747         ids = 0.0
748         acceptable = "ACDEFGHIKLMNPQRSTVWY"
749
750         for i in range (len (seq1)):
751             if seq1[i] not in acceptable and seq2[i] not in acceptable:
752                 continue
753             #if is_gap (seq1[i]) or is_gap (seq2[i]):
754             # continue
755             #if is_ambiguous (seq1[i]) or is_ambiguous (seq2[i]):
756             # continue
757             numAligned = numAligned + 1.0
758             if seq1[i] == seq2[i]:
759                 ids = ids + 1.0
760             if numAligned == 0:
761                 print "Warning:
polacco.Mulitalign.AlignRow.PercentIdentity, can't compare two
sequeunces because no meaningful characters align. pid set to 0."
762                 return 0
763             return ids/numAligned
764
765     def DashifyGapCharacters (self):
766         for i in range(len(self.chars)):
767             if is_gap (self.chars[i]):
768                 self.chars[i] = '-'
769
770
771
772 class ValdarConservation:
773     def __init__(self, multiAlign):
774         self.align = multiAlign
775         self.vMat = None
776
777     def Compute (self, multiAlign = None):
778         if not multiAlign:
779             multiAlign = self.align
780

```

```

781 |         if not self.vMat:
782 |             import polacco.Data
783 |             self.SetValdarMutMatrix (polacco.Data.GetPET91_matrix())
784 |         #first calculate sequence distance matrix
785 |         self.ComputeSequenceDistanceMatrix()
786 |         #second calculate sequence weights
787 |         self.ComputeSequenceWeights()
788 |         #calculate the normalizing denominator
789 |         self.ComputeTotalWeightProducts()
790 |         #finally calculate the conservation
791 |         conScores = []
792 |         numSeqs = len (self.align.rows)
793 |         for i in range (self.align.length()):
794 |             score = 0.0
795 |             for j in range (numSeqs):
796 |                 seqji = self.align.rows[j].chars[i].upper()
797 |                 for k in range (j+1, numSeqs):
798 |                     seqki = self.align.rows[k].chars[i].upper()
799 |                     if is_gap (seqji) or is_gap (seqki) or is_ambiguous
(800 | (seqji) or is_ambiguous(seqki):
801 |                         score = score + 0.0
802 |                     else:
803 |                         score += (self.seqWeights[self.align.rows[j]] *
self.seqWeights[self.align.rows[k]]) * self.vMat[seqji][seqki]
804 |                     conScores.append(score/self.sumProducts)
805 |                 self.conScores = conScores
806 |             return conScores
807 |
808 |     def SetValdarMutMatrix (self, mutMatrix):
809 |         self.vMat = {}
810 |         aas = mutMatrix.keys()
811 |         #first find min and max
812 |         minimum = maximum = None
813 |         for aaRow in aas:
814 |             minimum = min (minimum, min (mutMatrix[aaRow].values())) or
min (mutMatrix[aaRow].values())
815 |             maximum = max (maximum, max (mutMatrix[aaRow].values()))
816 |             range = float(maximum - minimum)
817 |             for aaRow in aas:
818 |                 self.vMat[aaRow] = {}
819 |                 for aaCol in aas:
820 |                     self.vMat[aaRow][aaCol] = (mutMatrix[aaRow][aaCol] -
minimum)/range
821 |             return self.vMat
822 |
823 |     def ComputeSequenceDistanceMatrix(self):
824 |         self.seqDistances = {}
825 |         for seqRow in self.align.rows:
826 |             self.seqDistances[seqRow] = {}
827 |             for seqCol in self.align.rows:
828 |                 try:
829 |                     self.seqDistances[seqRow][seqCol] =
self.seqDistances[seqCol][seqRow]
830 |                 except KeyError:
831 |                     self.seqDistances[seqRow][seqCol] = self.SeqDistance
(seqRow.chars, seqCol.chars)
832 |

```





```

14 |     self.hits = []
15 |     self.caMatrix = None
16 |     self.scMatrix = None
17 |     self.lowMemory = lowMemory
18 |     self.titleFromFileName = 0
19 |
20 |     def ParseSpasmHits (self, openfile, loadInMemory = 1, outFile =
None):
21 |         line = openfile.readline()
22 |         hit = None
23 |         while (1):
24 |             if line == '':
25 |                 break
26 |             #read until we find a " ==> HIT : (1CLX) " line
27 |             if len (line) < 8 or (line[:8] != " ==> HIT" and line[:8]
!= " ==> TRY"):
28 |                 if len (line) > 6 and line[1:6] == 'ERROR':
29 |                     print "Error found in current file"
30 |                     print line
31 |                     if (hit):
32 |                         print "File: " + hit.fileName
33 |                     line = openfile.readline()
34 |                     continue
35 |                 else:
36 |                     hit = SpasmHit(self)
37 |                     line = hit.ReadOpenFile (line, openfile)
38 |                     if loadInMemory:
39 |                         self.hits.append (hit)
40 |                     else:
41 |                         self.hits = (hit)
42 |                         if outFile:
43 |                             self.WriteOutTable (outFile)
44 |
45 |                 if loadInMemory and outFile:
46 |                     self.WriteOutTable (outFile)
47 |
48 |     def GetBestMatchesPerHit (self, useDistanceRmsd=0):
49 |         #get best match (residues) for each hit (pdb)
50 |         matches = []
51 |         for hit in self.hits:
52 |             matches.append (hit.GetBestRMSDMatch(useDistanceRmsd))
53 |             #save a reference to the hit
54 |             matches[-1].hit = hit
55 |
56 |         if useDistanceRmsd:
57 |             indexedMatches = [(m.distanceRmsd, m) for m in matches]
58 |         else:
59 |             indexedMatches = [(m.rmsd, m) for m in matches]
60 |         indexedMatches.sort()
61 |         matches = [m[1] for m in indexedMatches]
62 |
63 |
64 |     # def __matchSort (a,b):
65 |     #     if a.rmsd < b.rmsd:
66 |     #         return -1
67 |     #     elif a.rmsd > b.rmsd:
68 |     #         return 1
69 |     #     else:

```

```

70 | #         return 0
71 | #     matches.sort (__matchSort)
72 |     return matches
73 |
74 |
75 |     def WriteOutTable (self, openFile):
76 |         #every class writes out its prefix on the line.
77 |         prefix = ""
78 |         for hit in self.hits:
79 |             hit.WriteOutTable (openFile, prefix)
80 |
81 |
82 |     def WriteOutScoredSortedTable (self, openFile, trueHash,
useFileName = 0, notTrueHash = None):
83 |         matches = self.GetBestMatchesPerHit()
84 |
85 |         if not notTrueHash:
86 |             notTrueHash = {}
87 |
88 |         for match in matches:
89 |             if useFileName:
90 |                 name = os.path.splitext(os.path.basename
(match.hit.fileName))[0]
91 |             else:
92 |                 name = match.hit.title
93 |             try:
94 |                 trueHash[name]
95 |                 pre = '1' # a true positive
96 |             except KeyError:
97 |                 try:
98 |                     notTrueHash[name]
99 |                     pre = '*'
100 |                 except KeyError:
101 |                     pre = "0" # a false positive
102 |                 match.WriteOutTable (openFile,
match.hit.GetTablePrefix(pre))
103 |
104 |
105 |     #use to convert spasm to table with a small memory footprint
106 |     def ConvertToTable (self, openfile, outOpenFile):
107 |         doNotLoad = 0
108 |         self.ParseSpasmHits (openfile, doNotLoad, outOpenFile)
109 |
110 |
111 |
112 | class SpasmHit:
113 |     def __init__(self, search):
114 |         self.search = search
115 |         self.fileName = None
116 |         self.lowMemory = search.lowMemory
117 |         self.bestRMSD = 99999.9
118 |         self.titleFromFileName = self.search.titleFromFileName
119 |
120 |
121 |     def ReadOpenFile (self, titleLine, fp):
122 |         if self.titleFromFileName:
123 |             self.title = "noFileName!" #yet, will be set later
124 |         else:

```

```

125 |         self.title = titleLine.split(":")[1].strip()
126 |         self.title.strip
127 |         if self.title[0] == '(': # strip () if necessary
128 |             self.title = self.title[1:-1]
129 |
130 |     self.matches = []
131 |     line = fp.readline()
132 |     while (1):
133 |         if len (line) > 8 and line[:8] == " MATCH #":
134 |             match = SpasmMatch(self.search)
135 |             line = match.ReadOpenFile (line, fp)
136 |
137 |             if self.lowMemory: #only store the best scoring match
here
138 |                 if self.matches and match.rmsd > self.bestRMSD:
139 |                     match = None
140 |                     continue
141 |                 else:
142 |                     self.bestRMSD = match.rmsd
143 |                     self.matches = []
144 |                 self.matches.append (match)
145 |             elif line == '':
146 |                 break
147 |
148 |             elif len (line) > 5 and (line[:5] == " File"):
149 |                 self.fileName = line.split()[2]
150 |                 if self.fileName[0] == "(":
151 |                     self.fileName = self.fileName[1:-1]
152 |                 if self.titleFromFileName:
153 |                     #want dlqcrd2 from /pdbstyle-1.63/qc/dlqcrd2.ent
154 |                     self.title = string.upper
(os.path.basename(self.fileName))
155 |                     #remove extension:
156 |                     try:
157 |                         self.title = self.title
[:string.rindex(self.title, ".")]
158 |                     except ValueError:
159 |                         #no periods, leave title alone
160 |                         pass
161 |
162 |                 line = fp.readline()
163 |
164 |             elif len (line) > 8 and (line[:8] == " ==> HIT" or line[:8]
== " ==> TRY"):
165 |                 break
166 |             else:
167 |                 #blank lines and others that we will ignore?
168 |                 line = fp.readline()
169 |
170 |     return line
171 |
172 | def WriteOutTable (self, openFile, pre):
173 |     prefix = self.GetTablePrefix (pre)
174 |
175 |     for match in self.matches:
176 |         match.WriteOutTable (openFile, prefix)
177 |
178 | def GetTablePrefix (self, pre):

```

```

179         return pre + " %s %s" % (self.title, self.fileName)
180
181     def GetBestRMSDMatch(self, useDistanceRmsd = 0):
182         best = None
183         if useDistanceRmsd:
184             for match in self.matches:
185                 if not best or match.distanceRmsd < best.distanceRmsd:
186                     best = match
187         else:
188             for match in self.matches:
189                 if not best or match.rmsd < best.rmsd:
190                     best = match
191         return best
192
193
194     class SpasmMatch:
195         def __init__(self, search):
196             self.id = None
197             self.rmsd = None
198             self.scRmsd = None
199             self.caRmsd = None
200             self.distanceRmsd = None
201             self.numatoms = None
202             self.scMatrix = None
203             self.caMatrix = None
204             self.search = search
205             self.chain = ''
206
207     def ReadOpenFile (self, idLine, fp):
208         idwords = idLine.split()
209         self.id = idwords[2]
210
211         assert (idwords[4] == "RMSD")
212         self.rmsd = float(idwords[5])
213         self.numatoms = int (idwords[8])
214
215         #load the diff matrices
216         line = fp.readline()
217         while (1):
218             words = line.split()
219             if len(words) > 3 and words[3] == "matrix":
220                 matrix = SpasmMatrix()
221                 if words[1] == "SC":
222                     if words[0] == "Target":
223                         if self.search.scMatrix == None:
224                             self.search.scMatrix = matrix
225                     else:
226                         matrix.dummy = 1
227                 else:
228                     self.scMatrix = matrix
229             elif words[1] == "CA":
230                 if words[0] == "Target":
231                     if self.search.caMatrix == None:
232                         self.search.caMatrix = matrix
233                     else:
234                         matrix.dummy = 1
235             else:
236                 self.caMatrix = matrix

```

```

237 |         line = matrix.ReadOpenFile (fp)
238 |         elif line == '' or (len(words) > 1 and (words[0] == "MATCH"
or words[0] == "=>")):
239 |             break
240 |             else: #skip blank lines and other unexpected things
241 |                 line = fp.readline()
242 |                 if not self.search.scMatrix:
243 |                     print "WARNING! Did not read a scMatrix for : %s" % (
self.id)
244 |                 if not self.search.caMatrix:
245 |                     print "WARNING! Did not read a caMatrix for : %s" % (
self.id)
246 |
247 |             #calculate distance rmsd
248 |             if self.scMatrix and self.search.scMatrix:
249 |                 self.scRmsd = self.scMatrix.CalculateRMSD
(self.search.scMatrix)
250 |             if self.caMatrix and self.search.caMatrix:
251 |                 self.caRmsd = self.caMatrix.CalculateRMSD
(self.search.caMatrix)
252 |
253 |             if self.scRmsd != None and self.caRmsd != None:
254 |                 self.distanceRmsd = (self.scRmsd + self.caRmsd)/2
255 |             else:
256 |                 self.distanceRmsd = self.scRmsd
257 |                 if self.distanceRmsd == None:
258 |                     self.distanceRmsd = self.caRmsd
259 |                 if self.distanceRmsd == None:
260 |                     self.distanceRmsd = 9.99
261 |
262 |             return line
263 |
264 |         def CalcDiffs (self):
265 |             self.scDiffMatrix = SpasmMatrix()
266 |             self.caDiffMatrix = SpasmMatrix()
267 |
268 |             if self.scMatrix and self.search.scMatrix:
269 |                 self.scDiffMax =
self.scDiffMatrix.SetAsDifference(self.search.scMatrix,
self.scMatrix)
270 |             else:
271 |                 self.scDiffMax = 0.0
272 |
273 |             if self.caMatrix and self.search.caMatrix:
274 |                 self.caDiffMax =
self.caDiffMatrix.SetAsDifference(self.search.caMatrix,
self.caMatrix)
275 |             else:
276 |                 self.caDiffMax = 0.0
277 |
278 |         def WriteOutTable (self, outFile, pre):
279 |             self.CalcDiffs()
280 |             # if self.distanceRmsd != None:
281 |             #     calcRmsd = self.distanceRmsd
282 |             # else:
283 |             #     calcRmsd = -1.0
284 |             calcRmsd = self.distanceRmsd
285 |

```

```

286 |         outFile.write ("%6s %3s %3s %5.2f %5.2f %3d %4.1f %4.1f RES: "
%
287 |             (pre, self.id,
string.join(self.caMatrix.GetChains(), ''),
288 |             self.rmsd, calcRmsd, self.numatoms,
self.caDiffMax, self.scDiffMax))
289 |
290 |         outFile.write (self.caMatrix.OneLineResidues())
291 |
292 |         outFile.write ("\n")
293 |
294 | class SpasmMatrix:
295 |     def __init__(self):
296 |         self.dummy = 0
297 |
298 |     def ReadOpenFile (self, fp):
299 |         if not self.dummy:
300 |             self.rows = []
301 |             self.residues = []
302 |             numRows = -1
303 |             totalRows = 0
304 |             while (numRows < totalRows):
305 |                 line = fp.readline()
306 |                 if not self.dummy:
307 |                     res = SpasmResID (line[1:11])
308 |                     row = map (float, line[12:].split() )
309 |                     self.rows.append (row)
310 |                     self.residues.append (res)
311 |
312 |                 if not totalRows: #we're on the first row, so set number
of rows
313 |                     totalRows = len(line[12:].split())
314 |                     numRows = 1
315 |                 else:
316 |                     numRows = numRows + 1
317 |
318 |             return fp.readline()
319 |
320 |     def GetChains (self):
321 |         chains = []
322 |         for res in self.residues:
323 |             if res.chain in chains:
324 |                 continue
325 |             else:
326 |                 chains.append (res.chain)
327 |         if len (chains) == 0:
328 |             chains.append ('-')
329 |         return chains
330 |
331 |     def SetAsDifference (self, reference, other):
332 |         assert(len (reference.rows) == len (other.rows))
333 |         self.rows = []
334 |         self.residues = other.residues
335 |         maxDiff = -1
336 |         for i in range (len (other.rows)):
337 |             row = []
338 |             for j in range (len (other.rows)):
339 |                 diff = other.rows[i][j] - reference.rows[i][j]

```

```

340 |         maxDiff = max (maxDiff, abs(diff))
341 |         row.append (diff)
342 |         self.rows.append (row)
343 |     return maxDiff
344 |
345 |     def OneLineResidues (self):
346 |         return string.join (map (lambda a: "%s_%s_%s"%(a.type,
a.chain, a.position), self.residues), " ")
347 |
348 |     def CalculateRMSD (self, otherMatrix):
349 |         i = 0
350 |         sumSquareDevs=0
351 |         numSquareDevs=0
352 |         matSize = len (self.rows)
353 |         assert (matSize == len (otherMatrix.rows))
354 |         for i in range (0, matSize):
355 |             for j in range (i+1, matSize):
356 |                 dev = self.rows[i][j] - otherMatrix.rows[i][j]
357 |                 sumSquareDevs = sumSquareDevs + (dev * dev)
358 |                 numSquareDevs = numSquareDevs + 1
359 |         return math.sqrt (sumSquareDevs/numSquareDevs)
360 |
361 |
362 | class SpasmResID:
363 |     def __init__ (self, idString):
364 |         #string is "HIS A 339" or "HIS_A_339"
365 |         self.type = string.strip (idString[:3])
366 |         if idString[4] != " ":
367 |             self.chain = string.strip (idString[4])
368 |         else:
369 |             self.chain = string.strip ('-')
370 |         if idString[5] == '_':
371 |             self.position = string.strip (idString[6:])
372 |         else:
373 |             self.position = string.strip (idString[5:])
374 |
375 |
376 | def dummy__runFileStrings ():
377 |     pass
378 |
379 | #####
380 | # Formalized contents of a spasm *.com run file:
381 | #
382 | #
383 | #
384 | # use like:
385 | # spasmRunFileString % (
386 | #     spasmBinaryPath,
387 | #     maxhits,
388 | #     libpath,
389 | #     motifPath,
390 | #     fourLetterCode,
391 | #     maxRMSD,
392 | #     maxCADiff,
393 | #     MaxSCDiff,
394 | #     maxResolution,
395 | #     maxResidues,
396 | #     substitutionsAllowed, # 5 for user defined

```



```

397 | #           substitutionsString)    #multi line
398 | #
399 | #
400 | #####
401 | runFileStringSTDOUT = ""#automatically generated by Spasm.py
402 | %s maxhits %d<<EOI
403 | %s
404 | %s
405 | %s
406 | %f
407 | %f
408 | %f
409 | %f
410 | %d
411 | %d
412 | %s
413 | n
414 | n
415 | n
416 | y
417 | n
418 | n
419 | l
420 | n
421 | n
422 | n
423 | n
424 | n
425 | EOI
426 | ""
427 |
428 |
429 | runFileStringSTDOUT_scOnly = ""#automatically generated by Spasm.py
430 | %s maxhits %d<<EOI
431 | %s
432 | %s
433 | %s
434 | %f
435 | %f
436 | %f
437 | %f
438 | %d
439 | %d
440 | %s
441 | n
442 | n
443 | n
444 | y
445 | n
446 | n
447 | 2
448 | n
449 | n
450 | n
451 | n
452 | n
453 | EOI
454 | ""

```

```

455 | def Convert2SpasmFilesToSortedAndScoredTable (trueSpasm, falseSpasm,
      | tableFile, trueHash = None, useDistanceRmsd=0, returnString = 0):
456 |     fp = open (trueSpasm)
457 |     trueSearch = SpasmSearch (1)
458 |     trueSearch.ParseSpasmHits (fp)
459 |     fp.close()
460 |
461 |     fp = open(falseSpasm)
462 |     falseSearch = SpasmSearch (1)
463 |     falseSearch.ParseSpasmHits (fp)
464 |     fp.close()
465 |
466 |
467 |     return Convert2SpasmSearchesToSortedAndScoredTable (trueSearch,
      | falseSearch, tableFile, trueHash, useDistanceRmsd, returnString)
468 |
469 |
470 | def Convert2SpasmSearchesToSortedAndScoredTable (trueSearch,
      | falseSearch, tableFile, trueHash = None, useDistanceRmsd=0,
471 |                                                     returnString = 0,
      | writeFile=1, trueSkipHash=None, falseSkipHash=None):
472 |     if not trueHash:
473 |         trueHash = {}
474 |     if not trueSkipHash:
475 |         trueSkipHash = {}
476 |     if not falseSkipHash:
477 |         falseSkipHash = {}
478 |
479 |     if returnString:
480 |         stringFile = StringIO.StringIO()
481 |         fp = stringFile
482 |     else:
483 |         fp = open (tableFile, "w")
484 |
485 |     trueMatches = trueSearch.GetBestMatchesPerHit(useDistanceRmsd)
486 |     if falseSearch:
487 |         falseMatches =
      | falseSearch.GetBestMatchesPerHit(useDistanceRmsd)
488 |     else:
489 |         falseMatches = []
490 |     trueIndex = 0
491 |     for fm in falseMatches:
492 |         while (1):
493 |             if useDistanceRmsd:
494 |                 if not (trueIndex < len (trueMatches) and
      | trueMatches[trueIndex].distanceRmsd < fm.distanceRmsd):
495 |                     break
496 |             else:
497 |                 if not (trueIndex < len (trueMatches) and
      | trueMatches[trueIndex].rmsd < fm.rmsd):
498 |                     break
499 |
500 |             if trueSkipHash.has_key (trueMatches[trueIndex].hit.title):
501 |                 pre = "*"
502 |             else:
503 |                 pre = "1"
504 |

```

```

505 |         trueMatches[trueIndex].WriteOutTable (fp,
      | trueMatches[trueIndex].hit.GetTablePrefix(pre))
506 |         trueIndex += 1
507 |
508 |         name = fm.hit.title
509 |
510 |         #check that the false match isn't one that is marked as true
      | (lets me not have to remove these from false library)
511 |         if trueHash.has_key (name):
512 |             continue
513 | #
514 | #     try:
515 | #         trueHash[name]
516 | #         continue
517 | #     except KeyError:
518 | #         pass
519 |
520 |         if falseSkipHash.has_key (name):
521 |             pre = "#"
522 |         else:
523 |             pre = "0"
524 |         fm.WriteOutTable (fp, fm.hit.GetTablePrefix(pre))
525 |     #finally write any remaining true positives, important in cases
      | where no false positives were hit
526 |     while trueIndex < len (trueMatches):
527 |         if trueSkipHash.has_key (trueMatches[trueIndex].hit.title):
528 |             pre = "*"
529 |         else:
530 |             pre = "1"
531 |
532 |         trueMatches[trueIndex].WriteOutTable (fp,
      | trueMatches[trueIndex].hit.GetTablePrefix(pre))
533 |         trueIndex += 1
534 |
535 |     if returnString:
536 |         compressFileName = tableFile + ".gz"
537 |         fp = gzip.open (compressFileName, "w")
538 |         if writeFile:
539 |             fp.write (stringFile.getvalue())
540 |         else:
541 |             #still have to maintian file for how gasps checks which
      | motifs its tried already!
542 |             fp.write (":-)")
543 |             fp.close
544 |             return stringFile.getvalue()
545 |     else:
546 |         fp.close()
547 |
548 | def ComputeAreaFromTableFile (filePath, maxFalse, fileString = '',
      | useDistanceRMSD = 0):
549 |     if fileString:
550 |         fp = StringIO.StringIO (fileString)
551 |     else:
552 |         fp = open (filePath)
553 |     numFalse = 0
554 |     scoreHash = {}
555 |     #first load the lines we need in to scoreHash
556 |     while (numFalse < maxFalse):

```

```

557 |         line = fp.readline()
558 |         if line == '':
559 |             break
560 |         if line[0] == '#':
561 |             continue
562 |         elif line[0] == '*':
563 |             continue
564 |         elif line[0] == '0':
565 |             numFalse = numFalse + 1
566 |         if useDistanceRMSD:
567 |             score = float (line.split()[6])
568 |         else:
569 |             score = float (line.split()[5])
570 |         try:
571 |             scoreHash[score].append (line[0])
572 |         except KeyError:
573 |             scoreHash[score] = [line[0]]
574 |     fp.close()
575 |     numFalse = 0
576 |     numTrue = 0
577 |     area = 0
578 | # print scoreHash
579 | #next calculate area one false positive at a time.
580 | # if one score has several hits, false positives are counted
before true positives.
581 |     scores = scoreHash.keys()
582 |     scores.sort()
583 | # print scores
584 |     for score in scores:
585 |         hits = scoreHash[score]
586 |         newTrue = 0
587 |         for hit in hits:
588 |             if hit == '1':
589 |                 newTrue = newTrue + 1
590 |             elif hit == '0':
591 |                 numFalse = numFalse + 1
592 |                 area = area + numTrue
593 |         numTrue = numTrue + newTrue
594 |
595 |         if numFalse > maxFalse:
596 |             break
597 |
598 |     #if we don't have enough false positives assume all other hits
would be false.
599 |     while (numFalse < maxFalse):
600 |         print "Assuming a false positive"
601 |         area = area + numTrue
602 |         numFalse = numFalse + 1
603 |
604 |     return area
605 |
606 | def ComputeSeparationScoreFromTableFile3 (filePath, maxFalse,
maxRMSD, maxTrue, sepScoreImportance, useDistanceRmsd, fileString =
''):
607 |
608 |     if fileString:
609 |         fp = StringIO.StringIO (fileString)
610 |     else:

```

```

611 |         fp = open (filePath)
612 |         numFalse = 0
613 |         trueScores = []
614 |         falseScores = []
615 |         #read the relevant scores for true and false positives
616 |         while (numFalse < maxFalse):
617 |             line = fp.readline()
618 |             if line == '':
619 |                 break
620 |             if line[0] == '#':
621 |                 continue
622 |             elif line[0] == '*':
623 |                 continue
624 |             elif line[0] == '0':
625 |                 numFalse += 1
626 |                 if useDistanceRmsd:
627 |                     falseScores.append(float(line.split()[6]))
628 |                 else:
629 |                     falseScores.append(float(line.split()[5]))
630 |             elif line[0] == '1':
631 |                 if useDistanceRmsd:
632 |                     trueScores.append (float(line.split()[6]))
633 |                 else:
634 |                     trueScores.append (float(line.split()[5]))
635 |
636 |         #get score to assign to false positives that don't show up in
list.
637 |         if maxRMSD:
638 |             dummyFalseScore = maxRMSD
639 |         else:
640 |             dummyFalseScore = 0
641 |             if trueScores:
642 |                 dummyFalseScore = trueScores[-1]
643 |             if falseScores:
644 |                 dummyFalseScore = max (dummyFalseScore, falseScores[-1])
645 |             dummyFalseScore += 0.01
646 |
647 |         #fill out false positive list
648 |         while len(falseScores) < maxFalse:
649 |             falseScores.append (dummyFalseScore)
650 |
651 |         #remove trueScores that could just as easily have come after the
last false positive because they score equivalently
652 |         while (falseScores and trueScores and trueScores[-1] ==
falseScores[-1]):
653 |             del (trueScores[-1])
654 |
655 |
656 |         #calculate medianFP, relies on falseScores being sorted already
657 |         if len(falseScores) % 2:
658 |             #odd number just grab central scores
659 |             medianFP = falseScores[len(falseScores)/2]
660 |         else:
661 |             #even number, take center of two central scores
662 |             medianFP = falseScores[len(falseScores)/2] +
falseScores[len(falseScores)/2-1]
663 |             medianFP = medianFP/2
664 |

```

```

665 |
666 | #calculate meanFP while we also calculate rocArea
667 | meanFP = 0.0
668 | rocArea = 0.0
669 | for fs in falseScores:
670 |     meanFP += fs/maxFalse
671 |     for ts in trueScores:
672 |         if ts < fs:
673 |             rocArea += 1.0
674 | #normalize to 1
675 | rocArea = rocArea/(maxFalse*maxTrue)
676 |
677 |
678 | if trueScores:
679 |     #calculate medianTP
680 |     if len (trueScores ) % 2:
681 |         medianTP = trueScores[len(trueScores)/2]
682 |     else:
683 |         medianTP = trueScores[len(trueScores)/2] +
trueScores[len(trueScores)/2-1]
684 |         medianTP = medianTP/2
685 |
686 |     sepScore = (medianFP - medianTP)/medianFP # max 1, min 0
687 |     if sepScore < 0:
688 |         sepScore = 0
689 |
690 |
691 |     maxSepScore = (maxRMSD - 0.0) * len (trueScores) #potentially
undesirable; should use max number of trueScores instead?
692 |
693 |     #normalize:
694 |     sepScore = sepScore * sepScoreImportance
695 |
696 | else:
697 |     sepScore = 0
698 |
699 | return sepScore + rocArea
700 |

```

## ***polacco/XML.py***

```

1 |
2 | #
3 | #
4 | #
5 | #
6 | #
7 | #
8 | #
9 |
10 | # A VERY simple XML module. This was developed specifically to read
XML files as
11 | # output by ncbi blastall or blastpgp programs, and processed by
BlastXML.py
12 | # Many xml features may be poorly dealt with, if at all, if they do
not occur in

```

```

13 | # the output of the ncbi programs. Consider yourself warned if you
    | try to use this
14 | # as a full-featured XML parser.
15 |
16 |
17 |
18 | import string
19 | class XML_node:
20 |     def __init__ (self):
21 |         self.name = None
22 |         self.subNodes = {}
23 |         self.variables = {} #not common in blast output if any, these
    | are name value pairs within the first < >
24 |         self.value = None
25 |         self.parentNode = None
26 |
27 |     # builtins to get at the value :
28 |     def __int__(self):
29 |         return int (self.value)
30 |     def __str__(self):
31 |         return self.value
32 |     def __float__(self):
33 |         return float (self.value)
34 |
35 |
36 |     def keys(self):
37 |         return self.subNodes.keys()
38 |
39 |     def LoadFromOpenFile (self, openFileIn, startTag):
40 |         words = startTag[1:-1].split()
41 |
42 |         #get name from opening tag
43 |         self.name = words[0]
44 |
45 |         #get variableStrings -- this can be elaborated if anybody uses
    | them
46 |         self.variableStrings = words[1:]
47 |         if self.variableStrings:
48 |             for vs in self.variableStrings:
49 |                 name,value = vs.split("=")
50 |                 self.variables[name] = value[1:-1] #remove thequotes
51 |
52 |         #read value if any
53 |         valueBuffer = []
54 |         char = ''
55 |         while char != '<':
56 |             valueBuffer.append (char)
57 |             char = openFileIn.read(1)
58 |             if char == "":
59 |                 raise "Unexpected end of XML file %s" % (self.name)
60 |         self.value = string.join (valueBuffer, '')
61 |
62 |         #now start reading subnodes
63 |         tagType = ''
64 |         while tagType != '/':
65 |             tagType, tag = readXMLTag (openFileIn, char)
66 |             char = ''
67 |             if tagType == '':

```

```

68 |         nextNode = XML_node()
69 |         nextNode.LoadFromOpenFile (openFileIn, tag)
70 |         nextNode.parentNode = self
71 |         try:
72 |             self.subNodes[nextNode.name].append (nextNode)
73 |         except KeyError:
74 |             self.subNodes[nextNode.name] = [nextNode]
75 |     elif tagType != '/':
76 |         raise "Unexpected tag type %s in middle of file (%s)" %
      (tagType, self.name)
77 |
78 |     # tagType and tag should correspond to the closing tag for the
      current node.  verify it!
79 |     closingName = tag[2:-1]
80 |     if closingName != self.name:
81 |         raise "closingName doesn't match self.name: %s != %s" %
      (closingName, self.name)
82 |
83 |
84 | def readXMLTag (openFileIn, firstChar = ''):
85 |     tagBuffer= [firstChar]
86 |     type = ''
87 |
88 |     #first get opening carrot if we need to
89 |     if tagBuffer[-1] == '<':
90 |         started = 1
91 |     else:
92 |         char = ''
93 |         while (char != '<'):
94 |             char = openFileIn.read(1)
95 |             if char == "":
96 |                 raise "Unexpected end of XML file."
97 |             tagBuffer.append (char)
98 |
99 |     #get special first characters
100 |    char = openFileIn.read(1)
101 |    if char in '/?!':
102 |        type = char
103 |    else:
104 |        type = ''
105 |
106 |    tagBuffer.append (char)
107 |    done = 0
108 |    #get rest of tag
109 |    while (not done):
110 |        char = openFileIn.read(1)
111 |        if char == '>':
112 |            done = 1
113 |            if char in '\n\r':
114 |                continue
115 |            tagBuffer.append (char)
116 |
117 |    tag = string.join (tagBuffer, '')
118 |    #print tag
119 |
120 |    return type, tag
121 |
122 |

```



```

123 | class XML_tree:
124 |     def __init__ (self, openFileIn):
125 |         self.rootNode = None
126 |         self.xmlVersionString = None
127 |         self.doctype_string = None
128 |
129 |         if not openFileIn:
130 |             return
131 |
132 |         while(1):
133 |             char = openFileIn.read(1)
134 |             if char == '<':
135 |                 type,tag = readXMLTag (openFileIn, char)
136 |                 if type != '':
137 |                     if type == '?':
138 |                         self.xmlVersionString = tag[2:-2]
139 |                     elif type == '!':
140 |                         self.doctype_string = tag[2:-1]
141 |                     else:
142 |                         raise "Unexpected tag type at start of file: %s" %
type
143 |
144 |                 else: #found the start of the root node
145 |                     self.rootNode = XML_node()
146 |                     self.rootNode.LoadFromOpenFile (openFileIn, tag)
147 |                     break
148 |
149 |
150 |
151 |
152 | def test():
153 |     f = open ("longtest.xml")
154 |     xtree = XML_tree (f)

```

## ***polacco/utlis.py***

```

1 | #####3
2 | #
3 | #
4 | #
5 | #
6 | #
7 | #
8 | #####
9 |
10 | import string
11 |
12 |
13 |
14 | # This function will read a sequence in FASTA format up to and
including the '>'
15 | # character signifying the start of the next sequence. If I knew
how to peek
16 | # at this of put it back on the stream I would. If the first non-
white character

```

```

17 | # in the file is not '>' this function assumes that line is the name
    | of the sequence.
18 | # This never returns the first '>' on a line.  If for some reason
    | you have two, then you can expect one.
19 |
20 | def GetNextFASTASeq (seqFile):
    |     firstChar = " "
21 |     while firstChar in string.whitespace:
22 |         firstChar = seqFile.read (1)
23 |         if firstChar == ' ':
24 |             return ('', '')
25 |         if firstChar != '>':
26 |             name = string.rstrip (firstChar + seqFile.readline())
27 |         else:
28 |             name = string.strip (seqFile.readline())
29 |     parts = []
30 |     firstChar = seqFile.read(1)
31 |     while firstChar not in ( '>', '' ):
32 |         parts.append (string.strip (firstChar + seqFile.readline()))
33 |         firstChar = seqFile.read(1)
34 |     seq = string.join (parts, '')
35 |     return (name, seq)
36 |
37 |
38 |
39 | def WriteFastaSeq (fileName, seqName, sequence):
    |     fp = open (fileName, "w")
40 |     fp.write (">%s\n%s\n" % (seqName, sequence))
41 |     fp.close()
42 |
43 |
44 |
45 | aa3to1 = {
    |     "PHE" : 'F',
46 |     "ILE" : 'I',
47 |     "LEU" : 'L',
48 |     "VAL" : 'V',
49 |     "PRO" : 'P',
50 |     "ALA" : 'A',
51 |     "GLY" : 'G',
52 |     "MET" : 'M',
53 |     "CYS" : 'C',
54 |     "TRP" : 'W',
55 |     "TYR" : 'Y',
56 |     "THR" : 'T',
57 |     "SER" : 'S',
58 |     "GLN" : 'Q',
59 |     "ASN" : 'N',
60 |     "GLU" : 'E',
61 |     "ASP" : 'D',
62 |     "HIS" : 'H',
63 |     "LYS" : 'K',
64 |     "ARG" : 'R',
65 |     "GAP" : '-',
66 |     # non standard that we might run in to:
67 |     "MSE" : 'M',
68 |     "MME" : 'M',
69 |     "???" : 'X'
70 |
71 |     }
72 |

```

```

73 | aalto3 = {
74 |     'F' : "PHE",
75 |     'I' : "ILE",
76 |     'L' : "LEU",
77 |     'V' : "VAL",
78 |     'P' : "PRO",
79 |     'A' : "ALA",
80 |     'G' : "GLY",
81 |     'M' : "MET",
82 |     'C' : "CYS",
83 |     'W' : "TRP",
84 |     'Y' : "TYR",
85 |     'T' : "THR",
86 |     'S' : "SER",
87 |     'Q' : "GLN",
88 |     'N' : "ASN",
89 |     'E' : "GLU",
90 |     'D' : "ASP",
91 |     'H' : "HIS",
92 |     'K' : "LYS",
93 |     'R' : "ARG",
94 |     '-' : "GAP"
95 | }
96 |
97 | def AA3to1 (aaa):
98 |     return aa3tol[string.upper(aaa)]
99 |
100 | def AA1to3 (a):
101 |     return aalto3[string.upper (a)]
102 |
103 | def SeqAA3to1 (aaas):
104 |     as = []
105 |     for aaa in aaas:
106 |         as.append (AA3to1 (aaa))
107 |     return as
108 |
109 | def SeqAA1to3 (as):
110 |     aaas = []
111 |     for a in as:
112 |         try:
113 |             aaas.append (AA1to3 (a))
114 |         except KeyError:
115 |             aaa = a + a + a
116 |             aaas.append (aaa)
117 |     return aaas

```

### ***test/astral\_1.65\_SF.lib (partial)***

```

1 | ! Created by MKSPAZ V. 040618/2.3.3 at Thu Jan 27 10:13:05 2005 for
  | ben
2 | !
3 | PRO A6M_
4 | PDB /Users/ben/tmp_NoBackup/pdbstyle-1.65/a6/dla6m__.ent
5 | RES 1.00
6 | CMP
7 | VAL 1 -3.526 15.758 14.900 -4.746 16.634 16.149

```

```

8 | LEU      2  -0.689 14.190 16.862 1.731 14.740 15.905
9 | SER      3  -1.487 12.495 20.143 -1.027 10.685 20.278
10 | GLU     4   0.324 13.366 23.335 -0.274 15.439 25.046
11 | GLY     5   2.196 10.084 23.022 2.196 10.084 23.022
12 | GLU     6   3.317 10.981 19.508 1.056 10.123 17.224
13 | TRP     7   4.502 14.431 20.597 2.977 17.247 18.818
14 | GLN     8   6.475 12.812 23.418 7.581 10.572 24.982
15 | LEU     9   8.296 10.604 20.915 7.311 8.231 20.426
16 | VAL    10   9.019 13.628 18.670 7.839 14.702 17.503
17 | LEU    11  10.311 15.860 21.464 8.655 17.799 21.829
18 | HIS    12  12.315 13.068 23.090 13.323 10.588 24.655
19 | ...
20 | GLY A   74   1.747 -16.568 2.419 1.747 -16.568 2.419
21 | ARG A   75   4.859 -15.859 0.292 3.334 -18.781 -0.441
22 | VAL A   76   6.216 -12.816 -1.561 6.009 -11.548 -0.047
23 | GLU A   77   9.002 -11.725 -3.902 9.730 -10.534 -6.922
24 | ARG A   78  12.003 -10.462 -1.910 15.682 -10.880 -2.598
25 | SER A   79  11.298 -6.719 -1.500 9.563 -6.804 -2.269
26 | END
27 | !
28 | ! total residues 204251
29 | ! total proteins 1247

```

### test/d2mnr\_1.fasta

```

1 | >d2mnr_1 c.1.11.2 (133-359) Mandelate racemase {Pseudomonas putida}
  | RESOLUTION: 1.900000
2 | pvqaydshsldgvklateravtaaelgfravkktigypaldqdlavvrsirgavgdffgimvdynqsl
  | dvpaaikrsqalqqegvtwieeptlqhdyeghqriqsklnvpvqmgewlqgpeemfkalsigacrlam
  | pdamkiggvtgwirasalaqqfgipmsshlfqeisahllaatptahwlerldlagsvieptltfeggn
  | avipdlpgvggiwrekeigkylv

```

### test/d2mnr\_1.fasta.psiblast.xml.faln (partial)

```

1 | >d2mnr_1/001
2 | PVQAYD---S-H---S-LDGVKLATE---RA-VTA---A---EL-----GFRAV-KT-----
  | KI-----G-----YPA-----L-----DQ-----
  | -----DL-----AVVR-----SIRQAV-----GDDF--
  | G---I---MVD-Y-----NQS-L-----D-----V-P-----AA-
  | IKRSQAL-QQ---E-----G---VT-----W---IEE--PT-LQ---HD---YEGHQ-----
  | R-----I-QSKL---N-----V-P---VQM-GE-NW--L---GP-----E-
  | EMFKA-LS-IGAC---RL--AMPDAMKIGGVTGWIRASALAQ--Q--FG---I-P-M-----S-S-
  | ---H-----LF-----Q-----E-----IS-----AHL---LA-----
  | AT---PT---A-----H---W---LE---R-----L-----
  | -----DL---A-----G---SV---I---E---P-----T-----
  | LTFE-G-G--N--AV---I-P--D--LP--GVGIIWREKEIGKYLV
3 | >886661/143
4 | ---AWT---L-A---S-GDTARDIAE---AE-QML---E---AR-----RHRIF-KL-----
  | KI-----G-----ANP-----L-----EQ-----
  | -----DL-----KHVV-----AIKKAL-----GERA--
  | S---V---RVD-V-----NQY-W-----D-----E-S-----QA-
  | IRGCRVL-GD---N-----G---ID-----L---IEQ--PI-SR---VN---RSGQI-----
  | R-----L-NQRS---L-----A-P---IMA-DE-SI--E---SV-----E-
  | DAFSL-AA-DGAA---SV--FALKIAKNGGPRAVLRTAQIAE--A--AG---I-A-L-----Y-G-
  | ---G-----TM-----L-----E-----GS-----VGT---LA-----

```

```

SA-----HA-----FLTLRQLTWDTE----L----FG----P-----L-----
-----LL-----T-----E-----DI--V-----T--E-----R-----
POYR-D-F--H--LH----I-P--R--TP--GLGLTLDEERLARFR-
5 | ...

```

### **test/d2mnr\_1.pdb (partial)**

```

1 | HEADER      SCOP/ASTRAL domain d2mnr_1 [29245]          21-NOV-03    0000
2 | REMARK      99
3 | REMARK      99 ASTRAL ASTRAL-version: 1.65
4 | REMARK      99 ASTRAL SCOP-sid: d2mnr_1
5 | REMARK      99 ASTRAL SCOP-sun: 29245
6 | REMARK      99 ASTRAL SCOP-sccs: c.1.11.2
7 | REMARK      99 ASTRAL Source-PDB: 2mnr
8 | REMARK      99 ASTRAL Source-PDB-REVDAT: 31-JAN-94
9 | REMARK      99 ASTRAL Region: 133-359
10 | REMARK     99 ASTRAL ASTRAL-SPACI: 0.52
11 | REMARK     99 ASTRAL ASTRAL-AEROSPACI: 0.52
12 | REMARK     99 ASTRAL Data-updated-release: 1.61
13 | ATOM        954  N   PRO   133      26.117  28.195  19.354  1.00 16.42
    | 2MNR1142
14 | ATOM        955  CA  PRO   133      26.550  27.070  20.183  1.00 17.39
    | 2MNR1143
15 | ATOM        956  C   PRO   133      25.646  25.856  19.887  1.00 17.79
    | 2MNR1144
16 | ATOM        957  O   PRO   133      24.432  26.050  19.690  1.00 15.27
    | 2MNR1145
17 | ATOM        958  CB  PRO   133      26.423  27.575  21.607  1.00 16.88
    | 2MNR1146
18 | ATOM        959  CG  PRO   133      26.217  29.084  21.431  1.00 18.48
    | 2MNR1147
19 | ATOM        960  CD  PRO   133      25.391  29.164  20.163  1.00 16.15
    | 2MNR1148
20 | ATOM        961  N   VAL   134      26.203  24.632  19.793  1.00 14.37
    | 2MNR1149
21 | ATOM        962  CA  VAL   134      25.376  23.431  19.550  1.00 13.11
    | 2MNR1150
22 | ATOM        963  C   VAL   134      25.512  22.524  20.760  1.00 10.74
    | 2MNR1151
23 | ATOM        964  O   VAL   134      26.639  22.219  21.178  1.00 10.65
    | 2MNR1152
24 | ATOM        965  CB  VAL   134      25.822  22.626  18.293  1.00 14.98
    | 2MNR1153
25 | ATOM        966  CG1 VAL   134      24.813  21.481  18.039  1.00 15.66
    | 2MNR1154
26 | ATOM        967  CG2 VAL   134      25.879  23.546  17.071  1.00 18.17
    | 2MNR1155

```

### **test/enolase.lib (partial)**

```

1 | ! Created by MKSPAZ V. 040618/2.3.3 at Wed Jul 20 10:14:20 2005 for
    | ben
2 | !
3 | PRO ONEA
4 | PDB scopc.1.11/dlonea1.pdb
5 | RES 1.80

```

```
6 | CMP
7 | SER A 142 3.954 -21.471 10.866 3.132 -23.070 10.341
8 | PRO A 143 6.317 -20.493 9.738 6.437 -19.927 7.955
9 | TYR A 144 8.244 -18.231 11.988 8.330 -19.865 15.458
10 | VAL A 145 9.753 -15.108 10.378 8.597 -13.779 9.428
11 | LEU A 146 13.368 -14.050 10.947 14.866 -16.299 11.239
12 | PRO A 147 14.071 -10.342 10.694 14.559 -10.674 12.458
13 | VAL A 148 16.213 -8.369 8.330 15.891 -7.404 6.671
14 | PRO A 149 18.854 -6.741 10.645 19.911 -8.202 10.014
15 | PHE A 150 19.145 -2.899 10.152 16.286 -2.033 8.574
```

### ***test/enolase.list***

```
1 | d1onea1
2 | d1kkoa1
3 | d1fhua1
4 | d1muca1
5 | d1ec7a1
6 | d2mnr_1
7 | d1jpmal
```

## Appendix 2: GASPSdb CGI scripts

This appendix contains the scripts written in the python programming language that ran on the UCSF Resource for Biological Visualization and Informatics web server. The GASPSdb file contained the vast majority of the code to run the backend and interface with the database. The jsonMotif file's sole responsibility was to describe a requested motif so that the pages delivered by GASPSdb could show a motif in a popup window.

### ***GASPSdb***

```
1 | #!/usr/local/bin/python2.4
2 |
3 |
4 |
5 | import cgi, urllib
6 | import sys, tempfile, os, string, stat
7 | sys.path.insert(0, "/mol/sfld/gaspsdb/py.packages")
8 | import polacco.Rigor
9 |
10 | def LibraryLookup():
11 |     libs = {"scop3" : "/mol/sfld/gaspsdb/lib/scop3.emp.25and40.rig",
12 |            "scop4" : "/mol/sfld/gaspsdb/lib/scop4.emp.25and40.rig",
13 |            "GO"    : "/mol/sfld/gaspsdb/lib/go.emp.25and40.rig",
14 |            "goScop": "/mol/sfld/gaspsdb/lib/goScop.emp.25and40.rig"}
15 |     return libs
16 |
17 | def DrawMainForm():
18 |     print"""
19 |     <!-- DrawForm -->
20 |     <p> This search relies on <a href="/references.html">RIGOR</a>
21 |     which
22 |     is freely provided by its creator, Gerard Kleywegt, for use by
23 |     private
24 |     individuals, schools, academics, and not-for-profit institutions.
25 |     <strong>The use of this search by others is not allowed.</strong>
26 |     Contact us
27 |     if you wish to use our motifs, and contact Gerard Kleywegt if you
28 |     wish to
29 |     use <a href="/references.html">RIGOR</a>.
30 |     <table border=1>
31 |     <tr><th> Coordinates in PDB format:</th><th>Select a
32 |     library</th><th>Click to start search</th></tr>
33 |     <form action="./GASPSdb" enctype="multipart/form-data"
34 |     method="post">
35 |     <input type="hidden" name="do" value="rigor">
36 |     <tr>
```

```

32 |     <td>Enter a PDB code (e.g., '2mnr')<input type = "text"
    | name="pdb" size = "6">
33 |     <br>or select a file to upload:
34 |     <br><input type="file" name="fileUpload" size="20"></td>
35 |     <td>
36 |     <input name="lib" type="radio" value="scop3"> SCOP Superfamilies
37 |     <br> <input name="lib" type="radio" value="scop4"> SCOP Families
38 |     <br> <input name="lib" type="radio" value="GO"> Gene Ontology
    | (GO)
39 |     <br> <input name="lib" type="radio" value="goScop"> SCOP
    | Superfamilies / GO
40 |     </td>
41 |     <td>
42 |     <input type="submit" value="Send">
43 |     </td></tr>
44 | </form>
45 | </table>
46 | ""
47 |
48 |
49 | def DrawKeyWordSearchForm(instructions=True):
50 |     if instructions:
51 |         print""
52 |         <p> To look for a specific group, enter keywords below.
    | Separate words are automatically joined with 'AND'.
53 |         ""
54 |         print""
55 |         <p><form action="./GASPSdb" method="get">
56 |         <input type="hidden" name="do" value="keySearch">
57 |         <input name="keywords" type="text" size="40">
58 |         <button type="submit" id="Search">Search</button>
59 |         </form>
60 |         ""
61 | def DrawStructSearchForm (instructions = True):
62 |     if instructions:
63 |         print""
64 |         <p> To find groups that contain a structure, enter its PDB id
    | below. Multiple IDs are automatically joined with 'OR'.
65 |         ""
66 |         print""
67 |         <p><form action="./GASPSdb" method="get">
68 |         <input type="hidden" name="do" value="structSearch">
69 |         <input name="structs" type="text" size="40">
70 |         <input name="hideMotifs" type="hidden" value="TRUE">
71 |         <button type="submit" id="Search">Search</button>
72 |         </form>
73 |         ""
74 |
75 | def DrawSearchPage ():
76 |     print "<h2> Find motifs that match your structure.</h2>"
77 |     DrawMainForm()
78 |     print "<hr>"
79 |     print "<h2> Find groups by key word.</h2>"
80 |     DrawKeyWordSearchForm()
81 |     print "<hr>"
82 |     print "<h2> Find groups and motifs by PDB ID</h2>"
83 |     DrawStructSearchForm()
84 |

```



```

85 | def GetLocalPDBFile(pdbCode):
86 |     pdbCode = string.lower (pdbCode)
87 |     fp = None
88 |     if len (pdbCode) == 4:
89 |         path = "/databases/mol/pdb/%s/pdb%s.ent" % (pdbCode[1:3],
pdbCode)
90 |         try:
91 |             fp = open (path)
92 |         except IOError:
93 |             fp = None
94 |     else:
95 |         print "<p> PDB codes must be four characters long, '1one', for
example."
96 |         if not fp:
97 |             print "<p> There was an error retrieving the coordinates for
pdb %s. Please try to upload the coordinates." % pdbCode
98 |         return None
99 |     else:
100 |         return fp
101 |
102 | def GetFile(form):
103 |     pdbFile = None
104 |     pdbCode = form.getfirst ("pdb", "")
105 |     if pdbCode:
106 |         pdbFile = GetLocalPDBFile (pdbCode)
107 |         fileName = pdbCode
108 |     if not pdbFile:
109 |         try:
110 |             uploadedFile = form["fileUpload"]
111 |             fileName = uploadedFile.filename
112 |             pdbFile = uploadedFile.file
113 |         except KeyError:
114 |             print "<p> No file uploaded."
115 |             return None
116 |     lib = form.getfirst ("lib", "")
117 |     if pdbFile:
118 |         #set up a temporary file to write the uploaded file to:
119 |         fd, localFileName = tempfile.mkstemp (prefix= fileName +
"_" ,suffix = "_" +lib, dir="/var/tmp/gaspsdb/rigorRuns")
120 |         lineCount = 0
121 |         atomCount = 0
122 |         model = 0
123 |         endmdl = 0
124 |         for line in pdbFile:
125 |             if line[0:6] == "ATOM ":
126 |                 atomCount+=1
127 |                 lineCount +=1
128 |             if line[0:6] == "MODEL ":
129 |                 model += 1
130 |             if model and line[0:6] == "ENDMDL":
131 |                 endmdl = 1
132 |             if not (model and endmdl):
133 |                 os.write (fd, line)
134 |         os.close (fd)
135 |         os.chmod (localFileName, stat.S_IROTH | stat.S_IRUSR |
stat.S_IRGRP )
136 |         if atomCount:

```

```

137 |         print "<p>Received %d lines describing %d atoms from %s" %
      | (lineCount, atomCount, fileName)
138 |         if model and endmdl:
139 |             print "<p>Uploaded file contained %d models, only the first
      | model will be searched." % (model)
140 |             if atomCount < 100:
141 |                 print "<p><b>Warning:</b> The uploaded file (%s) seemed to
      | only contain %d atom records. This may not produce usable results."
      | % (uploadedFile.filename, atomCount)
142 |                 return localFileName
143 |             else:
144 |                 print "<p> fileUpload did not appear to be a file."
145 |                 return None
146 |
147 | def HTMLRedirectHead_DoRigor (code):
148 |     print ""
149 |     <html>
150 |     <meta HTTP-EQUIV="REFRESH" content="0;
      | url=http://babbittlab.ucsf.edu/cgi-
      | bin/GASPSdb?do=retrieveResults&code=%s">
151 |     </html>
152 |     "" % code
153 |
154 | def DrawFormatResultsForm (code = None, first=1, last=100):
155 |     print ""
156 |     <form action="./GASPSdb" enctype="application/x-www-form-
      | urlencoded" method="get">
157 |     <table width=600 border=1>
158 |     <tr><th> Click to format the results.</th><th> To view results
      | from a previous run, enter its code below.</th><tr>
159 |     <input type="hidden" name="do" value="retrieveResults">
160 |     <tr><td align="center"><input type="submit" value="Get
      | Results!"></td>
161 |     <td> <input type="text" name="code" size="40" value="%s"></td>
162 |     </tr>
163 |     <tr>
164 |         <td colspan="2" align="center">
165 |             <strong>Options:</strong>
166 |         </td>
167 |     </tr>
168 |     <tr>
169 |         <td> Which hits:</td>
170 |         <td><table>
171 |             <tr><td>First: <input type="text" name="first" size="8"
      | value="%d"></textare></td><td>Last: <input type="text" name="last"
      | size="8" value="%d"><small> Enter 'all' to show all.</small> </tr>
172 |             </table>
173 |         </td>
174 |     </tr>
175 |     </table>
176 |     </form>
177 |     "" % (code, first, last)
178 |
179 |
180 | def DoRigor(form):
181 |     localUploadFile = GetFile(form)
182 |     if not localUploadFile:
183 | #         HTMLHead()

```

```

184 | # DrawNavBar (form)
185 | print "<p> There was an error storing your uploaded file.
      | Can't run rigor."
186 |     return
187 |
188 |     lib = form.getfirst ("lib")
189 |     if not lib:
190 |         print "<p> You must select a motif library."
191 |         return
192 |     libraryPath = LibraryLookup()[form.getfirst("lib")]
193 |     rigorPath = "/mol/sfld/gaspsdb/al_rigor"
194 |     #files that will be written to by Rigor.py
195 |     rigorErr = localUploadFile + ".err"
196 |     rigorOut = localUploadFile + ".out"
197 |     rigorRun = localUploadFile + ".run"
198 |     pid = polacco.Rigor.RigorRun (localUploadFile, libraryPath,
      | rigorPath, rigorErr,
199 |                                     rigorOut, rigorRun,
      | useSubprocess=True)
200 |     #print "<p> rigor run started with pid %d." % pid
201 |     print "<p> Searching for matching motifs...<p> This run has been
      | given the code: <b>", os.path.split(localUploadFile)[1], "</b>"
202 |     DrawFormatResultsForm (os.path.split (localUploadFile)[1])
203 | # HTMLRedirectHead_DoRigor (os.path.split (localUploadFile)[1])
204 | # DrawNavBar(form)
205 | # print "Starting RIGOR run with code %s" % (os.path.split
      | (localUploadFile)[1])
206 |
207 | class TooManyHits (Exception):
208 |     pass
209 |
210 |
211 |
212 |
213 |
214 |
215 | def FixGroupName (group):
216 |     #two choices, the second char is a "." then this is a scop id
217 |     #(i.e., a.1.2.3) no need to fix
218 |     #otherwise the first is a number of a go id and needs a couple of
      | zeros prepended to it
219 |     if group[1] == ".":
220 |         return group
221 |     elif group[0] in "0123456789":
222 |         return "000" + group
223 |
224 |
225 |
226 | def DrawRigorHitsTable (hits, lib, code,first=0, last=None):
227 |     scores = hits.keys()
228 |     scores.sort()
229 |
230 |     if first > 1:
231 |         print '<p><strong>NOTE: Hits before hit #&#d are not shown
      | below.</strong>' % first
232 |
233 |     ToolTipScripts()
234 |     ToolTipDiv()

```

```

235 |
236 | print '<p><table class="smallTable">'
237 |
238 | if lib in ("scop3", "scop4", "goscop"):
239 |     domainHead = "Domain"
240 | elif lib == "go":
241 |     domainHead = "Chain"
242 | else:
243 |     domainHead = "Structure"
244 |
245 |     print
    '<th>#</th><th>E</th><th>RMSD</th><th>Group</th><th>%s</th><th>G-
    score</th><th>Residues</th><th>Matches</th><th></th>' % domainHead
246 |     rowSwitcher = TableRowGenerator()
247 |     count = 0
248 |     try:
249 |         for s in scores:
250 |             evText = "%4.1e" % s
251 |             if s<1e-4:
252 |                 evText = Green (evText)
253 |             elif s<5e-2:
254 |                 evText = Yellow(evText)
255 |             else:
256 |                 evText = Red(evText)
257 |             for h in hits[s]:
258 |                 motParts = h[0].split()
259 |                 if len (motParts) == 3:
260 |                     motGroup, motDomain, motScore = motParts
261 |                 elif len (motParts) == 2:
262 |                     motGroup, motDomain = motParts
263 |                     motScore = ""
264 |                 else:
265 |                     motGroup, motDomain, motScore = "* * *".split()
266 |
267 |                 if motScore and False: #currently broken so skip it for
now
268 |                     opString = h[5]
269 |                     imageLink = '<a
href="./MatchImage?code=%s&op=%s&group=%s&struct=%s"
target="_blank"> Image</a>' % (code, opString, motGroup, motDomain)
270 |                     else:
271 |                         imageLink = ''
272 |                         motGroup = FixGroupName(motGroup)
273 |                         groupName = motGroup
274 |                         structName = motDomain
275 |
276 |                         motGroup = '<a href
="./GASPSdb?do=describeGroup&group=%s" target="_blank">%s</a>' %
(motGroup, motGroup)
277 |                         if lib in ("scop3", "scop4"):
278 |                             #
                             motGroup = '<a
href="http://scop.berkeley.edu/search.cgi?ver=1.65&key=%s"
target="_blank">%s</a>' % (motGroup, motGroup)
279 |                             motDomain = '<a
href="http://scop.berkeley.edu/search.cgi?ver=1.65&key=%s"
target="_blank">%s</a>' % (motDomain, motDomain)
280 |                             elif lib == "GO":

```

```

281 | #             motGroup = '<a href="http://www.godatabase.org/cgi-
bin/amigo/go.cgi?action=replace_tree&search_constraint=terms&query=G
O:%07d" target="_blank">%s</a>' % (int(motGroup), motGroup)
282 |             pdbID = motDomain[0:4]
283 |             chain = motDomain[4]
284 |             motDomain = '<a
href="http://www.pdb.org/pdb/navbarsearch.do?inputQuickSearch=%s"
target="_blank">%s</a>%s' % (pdbID, pdbID, chain)
285 |
286 |             elif lib == "goScop":
287 |                 scop = string.join (motGroup.split(".")[0:3], ".")
288 |                 go = int (motGroup.split (".")[3])
289 |             #             motGroup = "" <a
href="http://scop.berkeley.edu/search.cgi?ver=1.65&key=%s"
target="_blank">%s </a>
290 |             #             <a
href="http://www.godatabase.org/cgi-
bin/amigo/go.cgi?action=replace_tree&search_constraint=terms&query=G
O:%07d" target="_blank"> %07d</a>
291 |             #             "" % (scop, scop, go, go)
292 |             pdbID = motDomain[0:4]
293 |             chain = motDomain[4]
294 |             motDomain = '<a
href="http://www.pdb.org/pdb/navbarsearch.do?inputQuickSearch=%s"
target="_blank">%s</a>%s' % (pdbID, pdbID, chain)
295 |
296 |             #             if altColor:
297 |             #                 print '<tr class="row2">'
298 |             #                 altColor = False
299 |             #             else:
300 |             #                 print '<tr class="row1">'
301 |             #                 altColor = True
302 |             #             residues=[]
303 |             #             resTypes = h[3].split(",")
304 |             #             resNames = h[4].split(",")
305 |             #             for i in range (len (resTypes)):
306 |             #                 residues.append (resNames[i]+resTypes[i])
307 |             #             resString = string.join (residues, "&res=")
308 |
309 |             count +=1
310 |             if count >= first:
311 |                 print TableRow (TableData(count, evText, h[1]) +
312 |                                 TooltipGroupName_td(motGroup,
groupName, structName=structName) +
313 |                                 TableData(motDomain, motScore, h[3],
h[4], imageLink),
314 |                                 rowSwitcher.next())
315 |             if last and count >= last:
316 |                 raise TooManyHits
317 |             print "</table>"
318 |
319 |             except TooManyHits:
320 |                 print "</table>"
321 |                 print "<strong> List of hits truncated after hit number %d
</strong><br>Modify 'first' and 'last' below to see other hits" %
last
322 |                 DrawFormatResultsForm(code, first, last)
323 |

```

```

324 |
325 | def RetrieveResults (form):
326 |     code = form.getfirst ("code")
327 |     if not code:
328 |         print "<p> Please enter a valid run code."
329 |         DrawFormatResultsForm()
330 |         return
331 |
332 |     last = form.getfirst ("last")
333 |     if not last:
334 |         last = "100"
335 |     try:
336 |         last = int (last)
337 |     except ValueError:
338 |         last = 0
339 |
340 |     first=form.getfirst ("first")
341 |     if not first:
342 |         first = "0"
343 |     try:
344 |         first = int (first)
345 |     except ValueError:
346 |         first = 0
347 |
348 |
349 |     rigorErr = os.path.join ("/var/tmp/gaspsdb/rigorRuns/", code +
".err")
350 |     try:
351 |         ferr = open (rigorErr)
352 |     except IOError:
353 |         print "<p> %s does not appear to be a valid code." % code
354 |         DrawFormatResultsForm (code)
355 |         return
356 |     for line in ferr:
357 |         #if line[0:4] == "STOP":# ... Toodle pip ... statement
executed"
358 |         if line[0:18] == "... Toodle pip ...":
359 |             ferr.close()
360 |             break
361 |     else:
362 |         ferr.close()
363 |         print "<p> Run %s not yet completed. Please try again in a few
seconds." % code
364 |         DrawFormatResultsForm (code)
365 |         return
366 |
367 |
368 |     rigorOut = os.path.join ("/var/tmp/gaspsdb/rigorRuns/", code +
".out")
369 |     fp = open (rigorOut)
370 |     hits,errorMessages = polacco.Rigor.SimpleRigorTranslate (fp,
silent=1)
371 |     fp.close()
372 |
373 |     lib = string.lower(code.split("_")[-1])
374 |     pdb = code.split ("_")[0]
375 |     if lib == "scop3":
376 |         library = "SCOP superfamilies"

```

```

377 | elif lib == "scop4":
378 |     library = "SCOP families"
379 | elif lib == "go":
380 |     library = "GO annotations"
381 | elif lib == "goscop":
382 |     library = "SCOP superfamilies subdivided by GO annotations"
383 | else:
384 |     library = "selected motifs"
385 | print ("<h2> Matches to %s among motifs from %s</h2>(code: %s)"
% (pdb, library, code) )
386 |
387 | if errorMessages:
388 |     print "<hr><p> <b> Warning:</b> The following error message(s)
was encountered while completing the rigor run."
389 |     for message in errorMessages:
390 |         print "<p>", message
391 |         print "<hr>"
392 |     DrawRigorHitsTable (hits, lib, code, first, last)
393 |
394 | def DrawDescribeForm ():
395 |     DrawBrowseChoices(title=False)
396 |     #print "<p> l a z y ."
397 |
398 | def GetCursorFromDatabase():
399 |     import MySQLdb
400 |     db = MySQLdb.connect ("""#### marked out for security #####""")
401 |     c = db.cursor()
402 |     return c
403 |
404 | def LibrarySort (toSort, index = None):
405 |     if index != None:
406 |         toSort = [ (string.split(l[index], "."), l) for l in toSort]
407 |     else:
408 |         toSort = [ (string.split(l, "."),l) for l in toSort]
409 |
410 |     for item in toSort:
411 |         for i in range (len (item[0]) ):
412 |             try:
413 |                 item[0][i] = int(item[0][i])
414 |             except ValueError:
415 |                 pass
416 |     toSort.sort()
417 |     return [item[1] for item in toSort]
418 |
419 |
420 | def DrawBrowseChoices (title=True):
421 |     if title:
422 |         print""
423 |         <div align="center" id="titleDiv">
424 |         <h1> Browse the Motifs </h1>
425 |         </div>
426 |         ""
427 |     print""
428 |     <p>
429 |     You can browse the motifs generated on
430 |     <ul>

```

```

431 |         <li><a href="/cgi-
bin/GASPSdb?do=browse&class=scop&depth=4">SCOP
families</a></li>
432 |         <li><a href="/cgi-
bin/GASPSdb?do=browse&class=scop&depth=3">SCOP
superfamilies</a></li>
433 |         <li><a href="/cgi-bin/GASPSdb?do=browse&class=go">GO
annotations</a></li>
434 |         <li><a href="/cgi-
bin/GASPSdb?do=browse&class=goscop">SCOP superfamilies
subdivided by GO annotations</a></li>
435 |     </ul>
436 |     ""
437 |     DrawKeywordSearchForm()
438 |
439 | def KeySearch (form):
440 |     c = GetCursorFromDatabase()
441 |     keyWords = form.getfirst ("keywords")
442 |     if not keyWords:
443 |         DrawKeywordSearchForm()
444 |         return
445 |     words = string.split (keyWords)
446 |     whereList = ["description like '%%s%%'" % word for word in words]
447 |     where = string.join (whereList , " AND ")
448 |
449 |     query = 'select g.name, g.description from groups g where
%s' % where
450 |     c.execute (query)
451 |     groups = c.fetchall()
452 |
453 |
454 |     print "<h2>Groups matching %s</h2>" % keyWords
455 |     if not groups:
456 |         print "No Groups Found!"
457 |     else:
458 |         print "Motif details can be viewed by clicking a group name."
459 |
460 |         print "<table>"
461 |         print TableHeader ("Group", "# Motifs", "Top G-Score",
"Description")
462 |         for (groupName, groupDesc) in groups:
463 |             c.execute ('select struct, gScore from motifs where
groupName = "%s" order by gScore' % groupName)
464 |             motifs = c.fetchall()
465 |             motCount = len (motifs)
466 |             if motifs:
467 |                 topScore = motifs[-1][1]
468 |             else:
469 |                 topScore = "NA"
470 |             groupName = '<a href = "./GASPSdb?do=describeGroup&group=%s"
>%s</a>' % (groupName, groupName)
471 |             print TableRow ( TableData (groupName, motCount, topScore,
groupDesc) )
472 |             print "</table>"
473 |             DrawKeywordSearchForm()
474 |
475 | def StructSearchGroups (form):
476 |     imagePath = "/images/"

```



```

477 | c = GetCursorFromDatabase()
478 | structs = form.getfirst ("structs")
479 | if not structs:
480 |     DrawStructSearchForm()
481 |     return
482 | hideMotifs = form.getfirst ("hideMotifs")
483 | if hideMotifs and string.upper (hideMotifs) == "TRUE":
484 |     hideMotifs = True
485 | else:
486 |     hideMotifs = False
487 |
488 | if hideMotifs:
489 |     toggle = "FALSE"
490 | else:
491 |     toggle = "TRUE"
492 | toggleMotifsURL = "./GASPSdb?" + urllib.urlencode
({ "do": "structSearch", "structs": structs, "hideMotifs": toggle })
493 | structs = string.split (structs)
494 | ignore = []
495 | for struct in structs:
496 |     if len (struct) != 4:
497 |         print "Structure %s excluded from search, pdb identifiers
should be at least four characters long."
498 |         ignore.append (struct)
499 | for struct in ignore:
500 |     structs.remove (struct)
501 |
502 | whereList = ["s.name like '%%s%%'" % struct for struct in structs]
503 | where = string.join (whereList, " OR ")
504 | query = 'select g.name, g.description, s.name from groups g inner
join group_struct gs on g.name = gs.groupName inner join structs s
on gs.structName = s.name where %s' % where
505 | #print query
506 | c.execute (query)
507 | groups = c.fetchall()
508 |
509 | print "<h2>Groups containing motifs from pdbs: %s</h2>" %
string.join (structs, " OR ")
510 | if not groups:
511 |     print "No Groups Found!"
512 | else:
513 |     print "Motif details can be viewed by clicking a group name."
514 |
515 |     print '<table class="smallTable">'
516 |     if hideMotifs:
517 |         print TableHeader ("Group", "Description", "Structure", '<a
href = "%s">Show Motifs</a>' % toggleMotifsURL)
518 |     else:
519 |         print TableHeader ("Group", "Description", "Structure",
"G-Score", "Motif", "Image")
520 |         tableRows=[]
521 |         rowSwitcher = TableRowGenerator()
522 |         for (groupName, groupDesc, matchStruct) in groups:
523 |             #get the motifs for each structure
524 |             c.execute ('select id, gScore from motifs where groupName =
"%s" and struct="%s"' % (groupName, matchStruct) )
525 |             motifs = c.fetchall()

```

```

526 |         groupNameLinked = '<a href
    | = "./GASPSdb?do=describeGroup&group=%s" >%s</a>' % (groupName,
    | groupName)
527 |         if motifs:
528 |             for motRow in motifs:
529 | #
530 | #                 #get the residues for each motif
531 | #                 c.execute ('select pdb, chain, name, resType from
    | motifs m inner join motif_residue mr on m.id = mr.motifID inner join
    | residues r on mr.resID = r.id where m.id = %d order by name' %
    | motRow[0] )
532 | #                 residues = c.fetchall()
533 | #                 try:
534 | #                     resToSort = [ (int(resRow[2]), resRow) for resRow
    | in residues]
535 | #                     resToSort.sort()
536 | #                     residues = [row[1] for row in resToSort]
537 | #                 except ValueError:
538 | #                     pass
539 | #                 resTable = '\n<table class="residue">' + string.join
    | ([TableRow(TableData (resRow[3], resRow[2], resRow[1]), resRow[3])
    | for resRow in residues], "\n") + "</table>"
540 |                 imageName = imagePath +
    | groupName+"_"+matchStruct+".r3d.png" #group_struct.r3d.png
541 |                 imageTag = ''
542 |                 try:
543 |                     gScore = float (motRow[1]) #gScore
544 |                 except TypeError:
545 |                     gScore = 0.0 #Acactually these should simply be
    | skipped
546 |                                     # but leave them in so I have a chance
    | of finding them!
547 |                                     #check http://babbittlab.ucsf.edu/cgi-
    | bin/GASPSdb?do=describeGroup&group=00008236
548 |                                     #first item in each item enables sorting on that item
    | (gScore)
549 |                 if hideMotifs:
550 |                     tableRows.append ( (gScore, TableData
    | (groupNameLinked, groupDesc, matchStruct, "")))
551 |                 else:
552 |                     resTable = GetResTableByMotifID (c, motRow[0])
553 |                     tableRows.append ( (gScore, TableData
    | (groupNameLinked, groupDesc, matchStruct, "%5.3f"%gScore, resTable,
    | imageTag)))
554 |                 else:
555 |                     if hideMotifs:
556 |                         tableRows.append ( (0.0, TableData (groupNameLinked,
    | groupDesc, matchStruct, "")))
557 |                     else:
558 |                         tableRows.append ( (0.0, TableData (groupNameLinked,
    | groupDesc, matchStruct, "", "", "No Motif Generated")))
559 |                 #if not motifs are found, still report a matching
    | group....
560 |                 for row in tableRows:
561 |                     print TableRow (row[1], rowSwitcher.next())
562 |                 print "</table>"
563 |

```

```

564
565
566
567
568
569
570 def Browse (form):
571     c = GetCursorFromDatabase()
572
573     minGScore = float (form.getfirst ("minGScore", "0.0"))
574
575     classification = form.getfirst ("class")
576
577     if not classification:
578         DrawBrowseChoices()
579         return
580
581     where = 'where r.classification = "%s"' % classification
582     depth = form.getfirst ("depth")
583
584     if depth:
585         try:
586             depth = int(depth)
587         except ValueError:
588             depth = 3
589         where = where + " and r.depth = %d" %depth
590
591     orderBy = ""
592     orderBy = form.getfirst ("orderBy", "")
593     if orderBy:
594         if orderBy == 'gScore':
595             orderBy = 'm.gScore'
596         else:
597             orderBy = ''
598     if orderBy:
599         orderBy = "order by " + orderBy
600
601     showMotifs= form.getfirst ("showMotifs", '')
602     c.execute ('select g.name, g.description, m.id, m.struct,
m.gScore, m.numMotifs from groups g inner join runs r on g.runID =
r.id inner join topMotifs m on g.name = m.groupName %s %s' % (where,
orderBy) )
603     groups = c.fetchall()
604
605     if not groups:
606         print "<p>Invalid classification</p><hr>"
607         DrawBrowseChoices()
608         return
609
610
611     #sort groups in an intuitive manner:
612     if not orderBy:
613         groups = LibrarySort (groups, 0)
614
615     classification = string.lower (classification)
616     if classification == "scop":
617         if depth == 3:
618             description = "SCOP Superfamilies"

```

```

619 |         elif depth == 4:
620 |             description = "SCOP Families"
621 |         else:
622 |             description = "SCOP groups"
623 |     elif classification == "go":
624 |         description = "GO Annotations"
625 |     elif classification == "goscop":
626 |         description = "SCOP Superfamilies subdivided by GO
annotations"
627 |     else:
628 |         description = "Selected Groups"
629 |
630 |     ToolTipScripts()
631 |     ToolTipDiv()
632 |     print "<h2>Browsing %s</h2>" % description
633 |     print "<strong>Hover</strong> over a group name to view a sample
motif."
634 |     print "<br><strong>Click</strong> a group name to view all motifs
for a group."
635 |
636 |     print "<table>"
637 |     print TableHeader ("Group", "# Motifs", "Top G-Score",
"Description")
638 |     for (groupName, groupDesc, motifID, struct, topScore, motCount)
in groups:
639 | #         c.execute ('select struct, gScore, id from motifs where
groupName = "%s" order by gScore' % groupName)
640 | #         motifs = c.fetchall()
641 | #         motCount = len (motifs)
642 | #         if motifs:
643 | #             topScore = motifs[-1][1]
644 | #             motifID = motifs[-1][2]
645 | #         else:
646 | #             topScore = "NA"
647 | #             motifID = ""
648 | #         if topScore == "NA" or topScore < minGScore:
649 | #             continue
650 |         groupNameHTML = '<a href
="./GASPSdb?do=describeGroup&group=%s" >%s</a>' % (groupName,
groupName)
651 |         if showMotifs:
652 |             groupDesc = '<table><th colspan=2>'+groupDesc+"</th>" +
TableRow(TableData(GetResTableByMotifID (c, motifID), GetImageTag
(groupName, struct))) + "</table>"
653 |             print TableRow ( ToolTipGroupName_td (groupNameHTML,
groupName, motifID) + TableData (motCount, topScore, groupDesc) )
654 |             print "</table>"
655 |
656 |
657 | def GetImageTag (groupName, struct):
658 |     imagePath = "/images/"
659 |     imageName = imagePath + groupName+"_"+struct+".r3d.png"
#group_struct.r3d.png
660 |     imageTag = ''
661 |     return imageTag
662 |
663 |

```

```

664 | def ToolTipGroupName_td (groupNameHTML, groupName, motifID=None,
      | structName=None):
665 |     if not motifID:
666 |         motifID=""
667 |     if not structName:
668 |         structName=""
669 |     td = ""
670 |     <td id = "%s" onmouseover="getGroupDataWithTO(this, '%s', '%s');"
      | onmouseout="cancelGroupData();">
671 |         %s
672 |     </td>
673 |     "" % (groupName, motifID, structName, groupNameHTML)
674 |     return td
675 |
676 | def ToolTipDiv ():
677 |     print ""
678 |     <div style="position:absolute;" id="popup" bgcolor="CFCFCF">
679 |         <table id="popupTable" bgcolor="CFCFCF" border="0"
      | cellspacing="0" cellpadding="0">
680 |             <thead>
681 |                 <tr><th id = "popupStatus" bgcolor = "CFCFCF"></th>
682 |                     <th id = "popupHeader" bgcolor = "CFCFCF" colspan="2"
      | align="left"></th></tr>
683 |             </thead>
684 |             <tbody id = "popupTableBody"></tbody>
685 |         </table>
686 |     </div>
687 |     ""
688 |
689 | def ToolTipScripts():
690 |     print ""
691 |     <script SRC="/js/MochiKit.js" TYPE="text/javascript"></script>
692 |     <script type = "text/javascript">
693 |         var lastDeferred;
694 |         var lastTO;
695 |
696 |         function getGroupDataWithTO (element, motifID, structName){
697 |             _doit = function(){
698 |                 getGroupData (element, motifID, structName);
699 |             }
700 |             lastTO = setTimeout ( "_doit();", 500);
701 |         }
702 |
703 |         function getGroupData(element, motifID, structName){
704 |             var url = "/cgi-bin/jsonMotif?group=" + escape (element.id) +
      | "&motifID=" + escape (motifID) + "&struct=" + (escape(structName));
705 |             var d = loadJSONDoc (url);
706 |             d.addCallback (partial (showMotifTable, element, motifID) );
707 |             lastDeferred = d;
708 |             showToolTip (element, status="loading...");
709 |         }
710 |
711 |         function cancelGroupData(){
712 |             document.getElementById ("popupHeader").innerHTML = "";
713 |             document.getElementById ("popupStatus").innerHTML = "";
714 |             if (lastDeferred) {
715 |                 lastDeferred.cancel();
716 |                 lastDeferred = null;

```

```

717 |     }
718 |     if (lastTO){
719 |         clearTimeout(lastTO);
720 |     }
721 |     clearData();
722 | }
723 |
724 | function showToolTip (element, status){
725 |     clearData();
726 |     setOffsets(element);
727 |     document.getElementById ("popupStatus").innerHTML = status;
728 |     document.getElementById ("popupTableBody").innerHTML= '<TR><td
width= "400" height="200" align="center"></td></TR>'
729 |     }
730 |
731 |
732 |     function showMotifTable (element, motifID, result){
733 |         document.getElementById ("popupStatus").innerHTML = "";
734 |         document.getElementById ("popupTableBody").innerHTML =
result["html"];
735 |         document.getElementById ("popupHeader").innerHTML =
result["group"]
736 |     }
737 |
738 |     function clearData(){
739 |         document.getElementById ("popupTableBody").innerHTML = "";
740 |         document.getElementById ("popup").style.border = "none";
741 |
742 |         document.getElementById
("popupTable").setAttribute('cellPadding',0);
743 |         document.getElementById
("popupTable").setAttribute('cellSpacing',0);
744 |     }
745 |
746 |
747 |     function calculateOffset (field, attr) {
748 |         var offset = 0;
749 |         while (field){
750 |             offset += field[attr];
751 |             field = field.offsetParent;
752 |         }
753 |         return offset;
754 |     }
755 |
756 |     function onTop (element) {
757 |         var height = 250;
758 |         var scrollTop = document.body.scrollTop;
759 |         var top = element.offsetHeight + calculateOffset (element,
"offsetTop");
760 |         return ( (top - height) > scrollTop);
761 |     }
762 |
763 |
764 |     function setOffsets(element){
765 |
766 |         var end = element.offsetWidth + calculateOffset (element,
"offsetLeft");

```

```

767 |         var top = element.offsetHeight + calculateOffset (element,
      | "offsetTop");
768 |         dataDiv = document.getElementById ("popup");
769 |         dataDiv.style.backgroundColor = "CFCFCF";
770 |         dataDiv.style.border = "black 1px solid";
771 |         dataDiv.style.left = end + 15 + "px";
772 |         var ot = onTop(element)
773 |         if (ot)
774 |             {dataDiv.style.top = top - 250 + "px";}
775 |         else
776 |             {dataDiv.style.top = top + 5 + "px";}
777 |
778 |         document.getElementById ("popupTable").setAttribute
      | ('cellPadding', 2);
779 |         document.getElementById ("popupTable").setAttribute
      | ('cellSpacing', 2);
780 |
781 |     }
782 | </script>
783 | ""
784 |
785 | def GetResTableByMotifID (c, motifID):
786 |     #get the residues for each motif
787 |     # c.execute ("""SELECT r.pdb, r.chain, r.name, r.resType,
      | l.shortLigandName, b.type, o.name
788 |     #         FROM motifs m
789 |     #         INNER JOIN motif_residue mr on m.id = mr.motifID
790 |     #         INNER JOIN residues r on mr.resID = r.id
791 |     #         LEFT JOIN ligandInts l on r.id = l.motResId
792 |     #         LEFT JOIN (SELECT * FROM bridges WHERE type="disulfide")
      | b on r.id = b.motResId
793 |     #         LEFT JOIN residues o on o.id = b.otherResId
794 |     #         WHERE m.id = %d ORDER BY r.name"" % motifID )
795 |     c.execute ("""SELECT r.pdb, r.chain, r.name, r.resType,
      | l.shortLigandName, b.type, o.name
796 |     FROM motifs m
797 |     INNER JOIN motif_residue mr on m.id = mr.motifID
798 |     INNER JOIN residues r on mr.resID = r.id
799 |     LEFT JOIN ligandInts l on r.id = l.motResId
800 |     LEFT JOIN bridges b on r.id = b.motResId
801 |     LEFT JOIN residues o on o.id = b.otherResId
802 |     WHERE m.id = %d ORDER BY r.name"" % motifID )
803 |     residues = c.fetchall()
804 |     try:
805 |         resToSort = [ (int(resRow[2]), resRow) for resRow in residues]
806 |         resToSort.sort()
807 |         residues = [row[1] for row in resToSort]
808 |     except ValueError:
809 |         pass
810 |
811 |     resTableRows = []
812 |     for resRow in residues:
813 |         if resRow[4]:
814 |             ligand = '~' + resRow[4]
815 |         else:
816 |             ligand = ''
817 |         if resRow[5] == 'disulfide':
818 |             bridge = "SS-"

```

```

819 |         if resRow[6]:
820 |             bridge = bridge + resRow[6]
821 |         elif resRow[5] == 'salt':
822 |             bridge = "+-"
823 |             if resRow[6]:
824 |                 bridge = bridge + resRow[6]
825 |         else:
826 |             bridge = ''
827 |         resTableRows.append ( TableRow(TableData (resRow[3],
resRow[2], resRow[1], ligand, bridge), resRow[3]) )
828 |
829 |     return '\n<table class="residue">' + string.join (resTableRows,
"\n")      +" </table>"
830 |
831 |
832 |
833 |
834 |
835 | def DescribeGroup (form):
836 |     imagePath = "/images/"
837 |     c = GetCursorFromDatabase()
838 |     # import MySQLdb
839 |     # c = db.cursor()
840 |
841 |     # get group id
842 |     # c.execute ('select id from groups where name = "%s";' %
groupName)
843 |     # groupIDs = c.fetchall()
844 |     # if not groupIDs:
845 |     #     print " <p> Group name %s not recognized. Please enter a valid
group name."
846 |     #     DrawDescribeFrom()
847 |     #     return
848 |     # #ignore the possibility of two or more groups with identical
names.
849 |     # groupID = int(groupIDs[0][0])
850 |
851 |     # don't be confused by the dummy for loop.  It's basically a
trick
852 |     # to make easier coding.  Basically a failure at any if
statements sends you to
853 |     # the same else statement.
854 |     for i in (1,):
855 |         groupName = form.getfirst ("group")
856 |         if groupName:
857 |             c.execute ('select description from groups where name =
"%s"' % groupName)
858 |             groupDesc = c.fetchone()
859 |             if groupDesc:
860 |                 break
861 |         else:
862 |             print "<p> Please enter a valid group name, or browse the
groups below."
863 |             DrawDescribeForm ()
864 |             return
865 |
866 |     #get all structures in a single group

```



```

867 |     c.execute ('select s.name, s.pdb, s.chain, s.description,
      | s.species, s.ec from structs s inner join group_struct gs on s.name
      | = gs.structName where gs.groupName="%s"' % groupName)
868 |     structures = c.fetchall()
869 |
870 |     print '<h2> Group: %s; %s </h2>' % (groupName, groupDesc[0])
871 |     print '<table class="smallTable">'
872 |     print TableHeader ("Structure", "G-score", "Motif", "Image")
873 |     tableRows = []
874 |     rowSwitcher = TableRowGenerator()
875 |     for structRow in structures:
876 |         # get the motifs if any for each structure
877 |         c.execute ('select id, gScore from motifs where groupName =
      | "%s" and struct="%s"' % (groupName, structRow[0]) )
878 |         motifs = c.fetchall()
879 |         for motRow in motifs:
880 |             resTable = GetResTableByMotifID (c, motRow[0])
881 |             imageName = imagePath +
      | groupName+"_"+structRow[0]+".r3d.png" #group_struct.r3d.png
882 |             imageTag = ''
883 |             try:
884 |                 gScore = float (motRow[1]) #gScore
885 |             except TypeError:
886 |                 gScore = 0.0 #Acdtually these should simply be skipped
887 |                 # but leave them in so I have a chance of
      | finding them!
888 |                 #check http://babbittlab.ucsf.edu/cgi-
      | bin/GASPSdb?do=describeGroup&group=00008236
889 |
890 |             tableRows.append ( (gScore, TableData (
      | DescribeStruct(structRow[0], structRow[3], structRow[4],
      | structRow[5]), "%5.3f"%gScore, resTable, imageTag)))
891 |             tableRows.sort()
892 |             tableRows.reverse()
893 |
894 |         for row in tableRows:
895 |             print TableRow (row[1], rowSwitcher.next())
896 |
897 |     print "</table>"
898 |
899 | def HotLinkStructName (structName):
900 |
901 |     #two choices 1a4ma and d1a4ma1
902 |     if len (structName) == 7:
903 |         return '<a
      | href="http://scop.berkeley.edu/search.cgi?ver=1.65&key=%s"
      | target="_blank">%s</a> <a href = "http://www.ebi.ac.uk/thornton-
      | srv/databases/cgi-
      | bin/pdbsum/GetPage.pl?pdbcode=%s&template=protein.html&l=1"
      | target="pdbSum">pdbSum</a>' % (structName, structName,
      | structName[1:5])
904 |     elif len (structName) == 5:
905 |         return '<a
      | href="http://www.pdb.org/pdb/explore/explore.do?structureId=%s"
      | target="_blank">%s</a>%s' % (structName[0:4], structName[0:4],
      | structName[4])
906 |

```

```

907 | def DescribeStruct (structName, description="", species="", ec=""):
908 |     if len (structName) ==7:
909 |         pdbID= structName[1:5]
910 |         chainID = structName[5]
911 |         astralName = structName
912 |
913 |     elif len (structName) == 5:
914 |         pdbID = structName[0:4]
915 |         chainID = structName[4]
916 |         astralName = None
917 |
918 |     if description:
919 |         descriptionLine = TableRow ("<td colspan=2>" + Bold(description)
+ "</td>")
920 |     else:
921 |         descriptionLine = ""
922 |     if species:
923 |         speciesLine = TableRow (TableData(Bold('Organism'), species))
924 |     else:
925 |         speciesLine=""
926 |
927 |     if ec:
928 |         ecLine = TableRow (TableData (Bold('EC'), ec))
929 |     else:
930 |         ecLine = ""
931 |
932 |     pdbLine = TableRow (TableData (Bold('pdb'), '<a
href="http://www.pdb.org/pdb/explore/explore.do?structureId=%s"
target="_blank">%s</a> chain %s' % (pdbID, pdbID, chainID))
933 |     otherLinks = TableRow (TableData (Bold('links'), "" "<a href =
"http://www.ebi.ac.uk/thornton-srv/databases/cgi-
bin/pdbsum/GetPage.pl?pdbcode=%s&template=protein.html&l=1"
target="pdbSum">pdbSum</a>
934 |                                     <br> <a
href="http://scop.berkeley.edu/search.cgi?key=%s"
target="scop">SCOP</a> "" % (pdbID, pdbID))
935 |     if astralName:
936 |         scopLine = TableRow (TableData (Bold('scop'), '<a
href="http://scop.berkeley.edu/search.cgi?ver=1.65&key=%s"
target="_blank">%s</a>' % (astralName, astralName))
937 |     else:
938 |         scopLine = ""
939 |
940 |     return '<table class="smallTable">' + descriptionLine + ecLine +
speciesLine + pdbLine + scopLine + otherLinks + "</table>"
941 |
942 |
943 | def TableRowGenerator():
944 |     switch = False
945 |     while True:
946 |         switch = not switch
947 |         if switch:
948 |             yield ("row1")
949 |         else:
950 |             yield ("row2")
951 |
952 |
953 | def TableRow (text, type = None):

```

```

954 |     if type:
955 |         first = "<tr class=\"%s\">" %type
956 |     else:
957 |         first = "<tr>"
958 |     return first + text + "</tr>"
959 |
960 |
961 | def TableData (*cells):
962 |     cells = ["%s" % s for s in cells]
963 |     return '<td>' + string.join (cells, '</td><td>') + '</td>'
964 | def TableHeader (*cells):
965 |     return '<th>' + string.join (cells, '</th><th>') + '</th>'
966 |
967 | def Bold (text):
968 |     return '<b>' + text+"</b>"
969 | def Green (text):
970 |     return '<font color="green">'+text+"</font>"
971 | def Yellow (text):
972 |     return '<font color="orange">'+text+"</font>"
973 | def Red (text):
974 |     return '<font color="red">'+text+"</font>"
975 |
976 |
977 |
978 | def HTMLHead():
979 |     print ""<html><HEAD>
980 |     <link rel=STYLESHEET type="text/css" href="/style.css">
981 |     <title>GASPS Motif Database</title>
982 |     </HEAD>
983 |     <body>
984 |     ""
985 |
986 | def DrawNavBar (form):
987 |     print""
988 |     <div id="header" class="header">
989 |     <map name="ucsfnosearch">
990 |
991 |         <area shape="rect" alt="UCSF home page" coords="38,3,84,27"
992 | href="http://www.ucsf.edu/">
993 |         <area shape="rect" alt="UCSF home page" coords="93,11,288,19"
994 | href="http://www.ucsf.edu/">
995 |         <area shape="rect" alt="About UCSF" coords="306,11,368,19"
996 | href="http://www.ucsf.edu/about_ucsf/">
997 |         <area shape="rect" alt="UCSF Medical Center"
998 | coords="387,11,498,19" href="http://www.ucsfhealth.org/" >
999 |
1000 |     </map>
1001 |
1002 |     <table border="0" width="100%" cellpadding="0" cellspacing="0"
summary="table used for layout purposes only">
1003 |     <tr bgcolor="#666666">
1004 |     <td></td>
1005 |     <td style="padding-top:5px; padding-right: 10px; vertical-
align:middle; text-align:right;">
1006 |     <a href="/"></a>

```

```

1003 |      &nbsp;&nbsp;&nbsp;<a href="http://www.rbvi.ucsf.edu"></a>
1004 |      </td>
1005 |    </tr>
1006 |  </table>
1007 |
1008 |    <!-- GASPSDB Nav Bar -->
1009 |    <table border="0" width="100%" cellpadding="0" cellspacing="0" >
1010 |      <tr bgcolor="#661166" style="text-align:center; color: white;
      font-family: arial,sans-serif; font-weight: bold;" >
1011 |        <td><a href="/" class="bigbutton">Main</a></td>
1012 |        <td><a href="/cgi-bin/GASPSdb?do=drawForm"
      class="bigbutton">Search</a></td>
1013 |        <td><a href="/cgi-bin/GASPSdb?do=retrieveResults"
      class="bigbutton">Get Results</a></td>
1014 |        <td><a href="/cgi-bin/GASPSdb?do=browse"
      class="bigbutton">Browse Motifs</a></td>
1015 |        <td><a href="/downloads.html" class =
      "bigbutton">Downloads</a></td>
1016 |        <td align="right">
1017 |          <table cellpadding="5px" >
1018 |            <tr>
1019 |              <td><a href="/help.html"
      class="smallbutton">Help</a></td>
1020 |              <td><a href="/references.html"
      class="smallbutton">References</a></td>
1021 |            </tr>
1022 |          </table></td>
1023 |        </tr>
1024 |      </table>
1025 |    </div>
1026 |
1027 |    <!-- Main body of page starts here -->
1028 |
1029 |    <div class = "body">
1030 |      ""
1031 |
1032 |    def HTMLTail():
1033 |      print "</div></body></html>"
1034 |      sys.exit(0)
1035 |
1036 |    def main():
1037 |
1038 |      print "Content-Type: text/html"      # HTML is following
1039 |      print                                # blank line, end of headers
1040 |
1041 |      form = cgi.FieldStorage()
1042 |
1043 |      debug = form.getfirst ("debug", "")
1044 |      if debug == "mince":
1045 |        import cgitb; cgitb.enable()
1046 |
1047 |
1048 |      do = form.getfirst ("do", "")
1049 |      if not do:
1050 |        do = "drawForm"
1051 |
1052 |    # if do=="rigor":

```

```

1053 | #     DoRigor(form)
1054 | # else:
1055 | #     #main switch function:
1056 |     HTMLHead()
1057 |     DrawNavBar (form)
1058 |     if do == "drawForm" or do == "search":
1059 |         DrawSearchPage()
1060 |     #     DrawMainForm()
1061 | # elif do == "search":
1062 | elif do == "rigor":
1063 |     DoRigor(form)
1064 | elif do == "retrieveResults":
1065 |     RetrieveResults(form)
1066 | elif do == "describeGroup":
1067 |     DescribeGroup (form)
1068 | elif do == "describeMotif":
1069 |     DescribeMotif (form)
1070 | elif do == "browse":
1071 |     Browse(form)
1072 | elif do == "keySearch":
1073 |     KeySearch(form)
1074 | elif do == "structSearch":
1075 |     StructSearchGroups (form)
1076 |
1077 | else:
1078 |     print "<p> Unrecognized do command", do, ". Try again."
1079 |     print "<hr>"
1080 |     DrawMainForm()
1081 | HTMLTail()
1082 |
1083 |
1084 |
1085 | if __name__ == "__main__":
1086 |     main()

```

## ***jsonMotif***

```

1 | #!/usr/local/bin/python2.4
2 |
3 | import cgi, string
4 |
5 | imagePath = "/images/"
6 |
7 | def GetCursorFromDatabase():
8 |     import MySQLdb
9 |     db = MySQLdb.connect ("""#### marked out for security ####""")
10 |     c = db.cursor()
11 |     return c
12 |
13 | def TableRow (text, type = None):
14 |     if type:
15 |         first = "<tr class=\"%s\">" %type
16 |     else:
17 |         first = "<tr>"
18 |     return first + text + "</tr>"
19 |

```

```

20 |
21 | def TableData (*cells):
22 |     cells = ["%s" % s for s in cells]
23 |     return '<td>' + string.join (cells, '</td><td>') + '</td>'
24 |
25 | def GetResTableByMotifID (c, motifID):
26 |     #get the residues for each motif
27 |     # c.execute ("""SELECT r.pdb, r.chain, r.name, r.resType,
28 |     l.shortLigandName, b.type, o.name
29 |     #         FROM motifs m
30 |     #         INNER JOIN motif_residue mr on m.id = mr.motifID
31 |     #         INNER JOIN residues r on mr.resID = r.id
32 |     #         LEFT JOIN ligandInts l on r.id = l.motResId
33 |     #         LEFT JOIN (SELECT * FROM bridges WHERE type="disulfide")
34 |     #         b on r.id = b.motResId
35 |     #         LEFT JOIN residues o on o.id = b.otherResId
36 |     #         WHERE m.id = %d ORDER BY r.name"" " % motifID )
37 |     c.execute ("""SELECT r.pdb, r.chain, r.name, r.resType,
38 |     l.shortLigandName, b.type, o.name
39 |     #         FROM motifs m
40 |     #         INNER JOIN motif_residue mr on m.id = mr.motifID
41 |     #         INNER JOIN residues r on mr.resID = r.id
42 |     #         LEFT JOIN ligandInts l on r.id= l.motResId
43 |     #         LEFT JOIN bridges b on r.id = b.motResId
44 |     #         LEFT JOIN residues o on o.id = b.otherResId
45 |     #         WHERE m.id = %d ORDER BY r.name"" " % motifID )
46 |     residues = c.fetchall()
47 |     try:
48 |         resToSort = [ (int(resRow[2]), resRow) for resRow in residues]
49 |         resToSort.sort()
50 |         residues = [row[1] for row in resToSort]
51 |     except ValueError:
52 |         pass
53 |
54 |     resTableRows = []
55 |     for resRow in residues:
56 |         if resRow[4]:
57 |             ligand = '~' + resRow[4]
58 |         else:
59 |             ligand = ''
60 |         if resRow[5] == 'disulfide':
61 |             bridge = "SS-"
62 |             if resRow[6]:
63 |                 bridge = bridge + resRow[6]
64 |         elif resRow[5] == 'salt':
65 |             bridge = "+-"
66 |             if resRow[6]:
67 |                 bridge = bridge + resRow[6]
68 |         else:
69 |             bridge = ''
70 |         resTableRows.append ( TableRow(TableData (resRow[3],
71 |         resRow[2], resRow[1], ligand, bridge), resRow[3]) )
72 |
73 |     return '\n<table class="residue">' + string.join (resTableRows,
74 |     "\n") + "</table>"

```

```

73 | def DescribeMotif (groupName, motifID=None, struct=None,
    | debug=False):
74 |     c = GetCursorFromDatabase();
75 |     if not motifID:
76 |         c.execute ('select id from motifs where groupName = "%s" and
    | struct = "%s"' % (groupName, struct))
77 |         motifIDS = c.fetchall()
78 |         assert len (motifIDS) == 1
79 |         motifID = motifIDS[0][0]
80 |
81 |     c.execute ("select groupName, struct, gScore from motifs where id
    | = %s" % motifID)
82 |     groupName, struct, gScore = c.fetchone()
83 |     resTable = GetResTableByMotifID (c, int(motifID))
84 | # c.execute ("""select resType, name, chain from residues r
    | #
    | #         inner join motif_residue mr on mr.resID = r.id
    | #         where mr.motifID = %s"""" % motifID)
87 | #
88 | # residues = c.fetchall()
89 | #
90 | # try:
91 | #     resToSort = [ (int(resRow[1]), resRow) for resRow in residues]
92 | #     resToSort.sort()
93 | #     resPairs = [row[1] for row in resToSort]
94 | # except ValueError:
95 | #     pass
96 | # resTable = '\n<table class="residue">' + string.join
    | ([TableRow(TableData (resRow[0], resRow[1], resRow[2]), resRow[0])
    | for resRow in residues], "\n") + "</table>"
97 |     imageName = imagePath + groupName+"_" + struct + ".r3d.png"
    | #group_struct.r3d.png
98 |     imageTag = ''
99 |     try:
100 |         gScore = float (gScore) #gScore
101 |     except TypeError:
102 |         gScore = 0.0 #Acdtually these should simply be skipped
103 |         # but leave them in so I have a chance of finding
    | them!
104 |         #check http://babbittlab.ucsf.edu/cgi-
    | bin/GASPSdb?do=describeGroup&group=00008236
105 |
106 |     textTable = """<table>
    |         %s
    |         %s
    |     </table>
    |     """%( TableRow (TableData
    | ("<strong>Structure:</strong>", struct)),
107 |         TableRow (TableData ("<strong>G-
    | Score:</strong>", "%5.3f" % gScore)) )
108 |
109 |     html = TableRow (TableData (textTable, resTable, imageTag))
110 |
111 |     c.execute ('select description from groups where name = "%s";' %
    | (groupName))
112 |     description = c.fetchall()
113 |     if len (description) == 1:
114 |         groupName = groupName + " : " + description[0][0]

```

```

119 |     if debug:
120 |         print "Content-Type: text/html"
121 |         print
122 |         print html
123 |         return
124 |
125 |     print "Content-Type: text/javascript"
126 |     print
127 |     print dict(html = html, group=groupName)
128 |
129 | def main():
130 |     form = cgi.FieldStorage()
131 |     group = form.getfirst ("group", "")
132 |     motifID = form.getfirst ("motifID", "")
133 |     struct = form.getfirst ("struct", "")
134 |     debug = form.getfirst ("debug", "")
135 |     if debug:
136 |         import cgitb; cgitb.enable()
137 |
138 |     DescribeMotif (group, motifID, struct, debug)
139 |
140 |
141 | if __name__ == "__main__":
142 |     main()
143 |
144 | """
145 | http://gaspsdb.rbvi.ucsf.edu/cgi-
    bin/jsonMotif?group=c.37.1&struct=dlin4a2&motifID=
146 | http://gaspsdb.rbvi.ucsf.edu/cgi-bin/jsonMotif?group=&struct=dlefval
147 | """

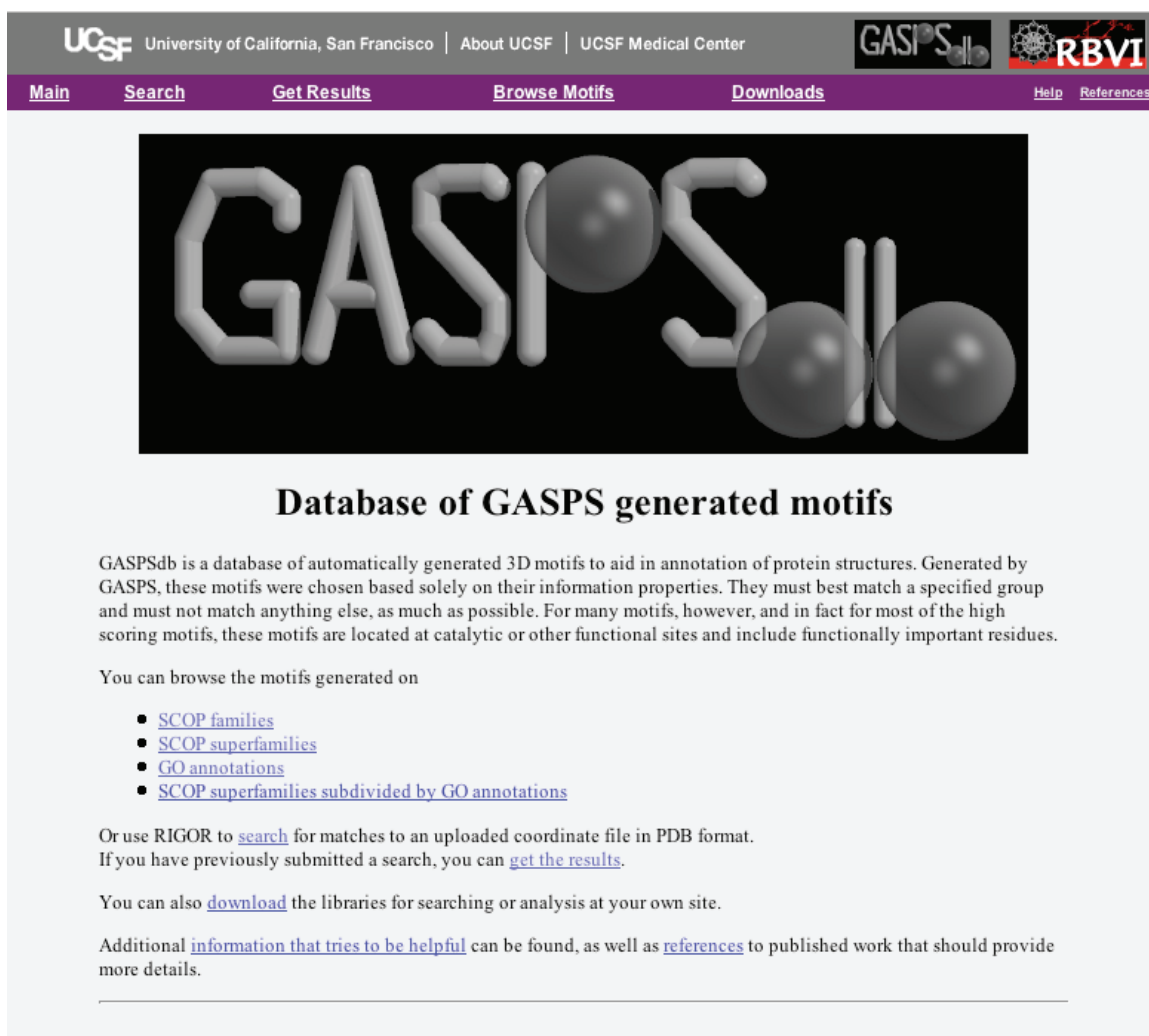
```



## Appendix 3: GASPSdb Web Interface

Shown here are screen shots and text from <http://gaspsdb.rbvi.ucsf.edu>.

### *GASPSdb Home Page*



UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

GASPS RBVI

[Main](#) [Search](#) [Get Results](#) [Browse Motifs](#) [Downloads](#) [Help](#) [References](#)

### Database of GASPS generated motifs

GASPSdb is a database of automatically generated 3D motifs to aid in annotation of protein structures. Generated by GASPS, these motifs were chosen based solely on their information properties. They must best match a specified group and must not match anything else, as much as possible. For many motifs, however, and in fact for most of the high scoring motifs, these motifs are located at catalytic or other functional sites and include functionally important residues.

You can browse the motifs generated on

- [SCOP families](#)
- [SCOP superfamilies](#)
- [GO annotations](#)
- [SCOP superfamilies subdivided by GO annotations](#)

Or use RIGOR to [search](#) for matches to an uploaded coordinate file in PDB format.  
If you have previously submitted a search, you can [get the results](#).

You can also [download](#) the libraries for searching or analysis at your own site.

Additional [information that tries to be helpful](#) can be found, as well as [references](#) to published work that should provide more details.

**Figure 1.** Home page of GASPSdb.

## GASPSdb Search Page

UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

GASPS RBVI

[Main](#) [Search](#) [Get Results](#) [Browse Motifs](#) [Downloads](#) [Help](#) [References](#)

### Find motifs that match your structure.

This search relies on [RIGOR](#) which is freely provided by its creator, Gerard Kleywegt, for use by private individuals, schools, academics, and not-for-profit institutions. **The use of this search by others is not allowed.** Contact us if you wish to use our motifs, and contact Gerard Kleywegt if you wish to use [RIGOR](#).

Coordinates in PDB format:	Select a library	Click to start search
Enter a PDB code (e.g., '2mnr') <input type="text"/> or select a file to upload: <input type="text"/> <input type="button" value="Browse..."/>	<input type="radio"/> SCOP Superfamilies <input type="radio"/> SCOP Families <input type="radio"/> Gene Ontology (GO) <input type="radio"/> SCOP Superfamilies / GO	<input type="button" value="Send"/>

---

### Find groups by key word.

To look for a specific group, enter keywords below. Separate words are automatically joined with 'AND'.

---

### Find groups and motifs by PDB ID

To find groups that contain a structure, enter its PDB id below. Multiple IDs are automatically joined with 'OR'.

Figure 2. GASPSdb Search Page.

## GASPSdb Browsing Page

UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

GASPS RBVI

Main Search Get Results Browse Motifs Downloads Help References

### Browsing SCOP Families

Hover over a group name to view a sample motif.  
Click a group name to view all motifs for a group.

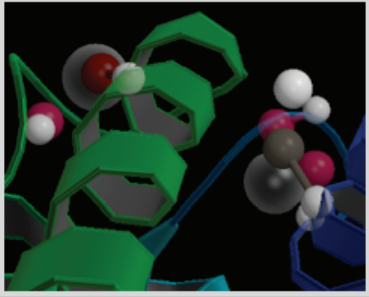
Group	# Motifs	Top G-Score	Description
<a href="#">a.1.1.2</a>	11	0.8622	Globins
<a href="#">a.3.1.1</a>	8	0.7661	
<a href="#">a.4.1.1</a>	9		
<a href="#">a.4.5.28</a>	7		
<a href="#">a.21.1.1</a>	7		
<a href="#">a.25.1.1</a>	8		
<a href="#">a.25.1.2</a>	7		
<a href="#">a.26.1.1</a>	7		
<a href="#">a.26.1.2</a>	6		
<a href="#">a.27.1.1</a>	8		
<a href="#">a.39.1.2</a>	7		
<a href="#">a.39.1.5</a>	7	0.7605	Calmodulin-like
<a href="#">a.45.1.1</a>	12	0.6219	Glutathione S-transferase (GST), C-terminal domain
<a href="#">a.74.1.1</a>	10	0.4135	Cyclin
<a href="#">a.104.1.1</a>	12	1.0732	Cytochrome P450
<a href="#">a.118.8.1</a>	6	0.7356	Tetratricopeptide repeat (TPR)
<a href="#">a.123.1.1</a>	7	1.09	Nuclear receptor ligand-binding domain
<a href="#">a.138.1.1</a>	8	1.0891	Cytochrome c3-like
<a href="#">a.138.1.3</a>	8	1.0906	Di-heme elbow motif
<a href="#">b.1.1.1</a>	16	0.6617	V set domains (antibody variable domain-like)
<a href="#">b.1.1.2</a>	8	0.8358	C1 set domains (antibody constant domain-like)
<a href="#">b.1.1.3</a>	7	0.7914	C2 set domains
<a href="#">b.1.1.4</a>	24	0.438	I set domains
<a href="#">b.1.2.1</a>	22	0.44	Fibronectin type III
<a href="#">b.1.18.2</a>	12	0.2649	E-set domains of sugar-utilizing enzymes

**a.39.1.5 : Calmodulin-like**

Structure: d1extra\_

G-Score: 0.760

PHE 16 A
ASP 20 A ~CA ++
ASP 24 A ~CA
GLY 25 A
ASP 56 A ~CA
GLU 67 A ~CA



tyl-tRNA

**Figure 3. GASPSdb Browse Page.**

Currently browsing scop superfamilies. The popup window is showing the top-scoring motif from the Calmodulin-like family in response to the pointer hovering over its group ID.

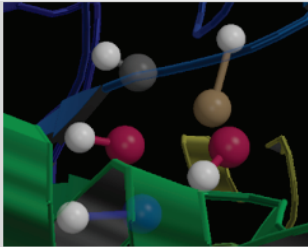
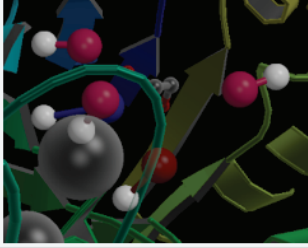
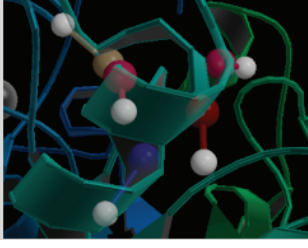
## GASPSdb Group Description Page

UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

GASPS RBVI

Main Search Get Results Browse Motifs Downloads Help References

Group: c.1.8.16160; SF: (Trans)glycosidases; GO: amylase activity

Structure	G-score	Motif	Image
<p><b>MALTOOLIGOSYL TREHALOSE SYNTHASE</b>            EC 5.4.99.15            Organism SULFOLOBUS ACIDOCALDARIUS (ARCHAEA)            pdb <a href="#">1iv8</a> chain a            links <a href="#">pdbSum</a>  <a href="#">SCOP</a></p>	1.088	<p>PRO 38 A            TYR 50 A            ASP 85 A +-            ARG 226 A +-228            ASP 228 A +-226</p>	
<p><b>ALPHA-AMYLASE ISOZYME 1</b>            EC 3.2.1.1            Organism HORDEUM VULGARE (BARLEY)            pdb <a href="#">1ht6</a> chain a            links <a href="#">pdbSum</a>  <a href="#">SCOP</a></p>	1.087	<p>ASP 88 A +-178            ARG 178 A +-88            ASP 180 A ~EDO +-178            GLU 205 A ~EDO +-178            ASP 291 A +-290</p>	
<p><b>ALPHA-AMYLASE I</b>            EC 3.2.1.1            Organism THERMOACTINOMYCES VULGARIS (BACTERIA)            pdb <a href="#">1ii1</a> chain a            links <a href="#">pdbSum</a>  <a href="#">SCOP</a></p>	1.087	<p>TYR 223 A            ARG 354 A +-356            ASP 356 A +-354            GLU 396 A +-354            ASP 472 A +-</p>	

**Figure 4. GASPSdb Group description page, partial.**

Showing the first three motifs from a group in the GO and SCOP combination groupings. The full page shows additional entries for each remaining motif in the group.

# GASPSdb Search Results Page

UCSF University of California, San Francisco | About UCSF | UCSF Medical Center

GASPS RBVI

Main Search Get Results Browse Motifs Downloads Help References

### Matches to 1rvk among motifs from SCOP superfamilies

(code: 1rvk\_VuLN0v\_scop3)

#	E	RMSD	Group	Domain	G-score	Residues	Matches
1	2.4e-04	0.24	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	1.0551	ASP, GLU, PRO	A206, A232, A234
2	7.2e-03	0.43	<a href="#">d.54.1</a>	<a href="#">d1ec7a1</a>	0.8710	LEU, ASP, HIS, GLY	A70, A110, A49, A288
3	1.2e-01	0.48	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	0.8047	LEU, ASP, HIS, GLY	A70, A110, A49, A288
4	2.0e-01	1.80	<a href="#">d.157.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
5	2.5e-01	1.82	<a href="#">d.157.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
6	8.5e-01	0.60	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
7	2.0e+00	1.17	<a href="#">c.2.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
8	2.1e+00	1.06	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
9	2.4e+00	1.07	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
10	4.0e+00	1.57	<a href="#">b.47.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
11	4.1e+00	1.12	<a href="#">c.1.11</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
12	4.5e+00	1.16	<a href="#">b.11.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
13	5.2e+00	1.58	<a href="#">b.29.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
14	5.4e+00	1.26	<a href="#">c.46.1</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
15	6.5e+00	1.22	<a href="#">a.4.5</a>	<a href="#">d1ec7a1</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
16	6.5e+00	1.28	<a href="#">c.46.1</a>	<a href="#">d1qb0a</a>	0.8710	LEU, ASP, HIS, GLY	A70, A110, A49, A288
17	6.7e+00	0.82	<a href="#">c.47.1</a>	<a href="#">d1e6ba2</a>	0.0857	VAL, TYR, ASP	A161, A339, A341
18	7.3e+00	0.75	<a href="#">e.3.1</a>	<a href="#">d1ci9a</a>	0.6977	SER, LYS, LEU	A136, A169, A250
19	8.5e+00	1.53	<a href="#">b.38.1</a>	<a href="#">d1n9ra</a>	0.8988	LEU, VAL, GLY, ASP, LEU	A81, A40, A48, A77, A111
20	9.9e+00	1.09	<a href="#">c.1.9</a>	<a href="#">d1o12a2</a>	0.9667	HIS, HIS, ASP	A24, A26, A20
21	1.3e+01	0.86	<a href="#">c.14.1</a>	<a href="#">d1tyfa</a>	0.3069	GLY, ALA, GLY, VAL	A113, A292, A287, A283
22	1.5e+01	1.37	<a href="#">c.46.1</a>	<a href="#">d1gmxa</a>	0.7905	LEU, ASP, HIS, GLY	A70, A110, A49, A288
23	1.7e+01	1.39	<a href="#">c.46.1</a>	<a href="#">d1rhs 2</a>	0.8047	LEU, ASP, HIS, GLY	A70, A110, A49, A288
24	1.7e+01	1.29	<a href="#">a.4.5</a>	<a href="#">d1ku9a</a>	0.1211	LEU, ASP, LEU, VAL	A367, A110, A293, A318
25	1.8e+01	1.60	<a href="#">b.38.1</a>	<a href="#">d1mgqa</a>	0.9030	LEU, VAL, GLY, ASP, LEU	A70, A40, A48, A77, A111
26	1.8e+01	0.68	<a href="#">c.66.1</a>	<a href="#">d1h1da</a>	0.1080	LEU, ALA, LEU	A123, A320, A293
27	1.8e+01	0.61	<a href="#">c.36.1</a>	<a href="#">d1poxa2</a>	0.3923	ALA, PRO, ALA	A112, A119, A321
28	1.8e+01	1.19	<a href="#">c.69.1</a>	<a href="#">d1qlwa</a>	0.5592	GLY, SER, GLY, HIS	A135, A136, A331, A171
29	1.8e+01	2.14	<a href="#">d.157.1</a>	<a href="#">d1k07a</a>	1.0278	LYS, HIS, ASP, GLY, HIS	A323, A209, A206, A23, A24
30	1.8e+01	0.50	<a href="#">d.54.1</a>	<a href="#">d1ec7a2</a>	0.6731	GLY, ALA, ASP	A45, A107, A110
31	1.9e+01	0.56	<a href="#">c.67.1</a>	<a href="#">d1mdoa</a>	0.2148	ILE, GLY, THR	A205, A331, A172
32	2.0e+01	1.18	<a href="#">b.38.1</a>	<a href="#">d1m5q1</a>	0.8927	VAL, GLY, ASP, LEU	A40, A48, A77, A111
33	2.0e+01	1.18	<a href="#">c.1.9</a>	<a href="#">d1p1ma2</a>	0.9756	HIS, HIS, ASP	A24, A26, A20
34	2.3e+01	1.20	<a href="#">c.1.9</a>	<a href="#">d1j5sa</a>	0.9904	HIS, HIS, ASP	A24, A26, A20

**c.1.11 : Enolase C-terminal domain-like**

Structure: d1ec7a1  
G-Score: 1.0551

ASP 235 A -MG +  
GLU 260 A -MG +  
PRO 262 A

**Figure 5.** Search results table for search of 1rvk against SCOP superfamily motifs. The popup window is displaying the first matched motif because the pointer is hovering over its group id.

## ***GASPSdb Help Page***

### **About GASPSdb**

#### **Where do these motifs come from?**

These motifs were generated using an entirely automated method given the acronym GASPS. For specific details please refer to the GASPS reference on the References Page. In short, GASPS chooses a motif from a single structure that best separates related structures from all other structures. The GASPS score or G score measures the degree of this separation. For the motifs here, we used several different systems to define related structures. In all cases, the set of structures was reduced to exclude mutants as well as sequence-redundant structures at the level of 40% or 25% identity.

#### **SCOP version 1.65**

Members of the the same superfamily or family are grouped together.

#### **Gene Ontology**

Structures are grouped according to their molecular function terms, automatically assigning parent terms where appropriate so that groups can be defined at any level in the GO hierarchy. Terms that give redundant groups to terms lower in the hierarchy are ignored as well as groups that appear to be too general (those with more than 50 structures.)

#### **SCOP superfamilies subdivided by GO annotations**

In an attempt to get isofunctional, homologous groups, SCOP superfamilies were subdivided by all assigned and implied GO molecular function terms. Terms that give redundant groups were ignored.

GASPS uses SPASM (see References) to identify matches, so that each residue is modeled as two points in space: one representing the alpha carbon, and one representing the side chain centroid.

### **How does the search work?**

We are indebted to others who have published and made available their motif searching tools for our use. Specifically Gerard Kleywegt and his RIGOR tool (see References).

Our search receives your structure file in PDB format and finds all motifs in the specified library that have a close match in the PDB file.

### **How do I interpret the search results?**

The search results on the search page are ranked according to an expectation value (E).

The expectation value is computed according to the model generated by Stark et al. (see References), and is based on the RMSD as well as the type and number of residues in the motif. The G score may also help decide whether a match represents a significant similarity.

### **What is the G score?**

Each motif is given a G score by GASPS. This is the score that GASPS tries to maximize as it constructs motifs. In short, a G score indicates how well conserved the motif is across the group, and how unique it is among unrelated proteins. This score has a theoretical range of 0-1.1, though any score near 1.0 is highly significant, and scores below about 0.4 are highly suspect. In cases with marginal E values, the G score may provide additional support.

In slightly more detail the G score is the sum of two components, the largest is the normalized area under an ROC style plot to a false positive rate of approximately 0.001, and the other component is the relative distance between true positive and false positive RMSD distributions. This latter component accounts for only 0 to 0.1 of the total G score,

so that most G scores above 1.0 imply perfect separation in an ROC style plot (ROC area = 1.0).

### **What do the motif images show?**

The motif images attempt to show the relative orientation of motifs and the local secondary structure and ligands in the protein. Residues in the motif are drawn as lopsided barbells. The smaller white sphere represents the alpha carbon and the larger colored sphere represents the side chain centroid. Side chains are colored according to residue type based on Bob Fletterick's 'shapely models' color scheme. Local secondary structure and ligands are drawn transparently to not mask the motifs.

These images are generated automatically with no effort to find a decent viewable orientation. With over 12,000 images, doing so by hand was not feasible.

## ***GASPSdb References Page***

### **References**

Several tools provided by us and others have been instrumental in getting this resource online. Where possible, the tool names link to the relevant web sites.

#### **GASPS**

Genetic Algorithm Search for Patterns in Structures Responsible for generation of all motifs:

Polacco, Benjamin J. and Patricia C. Babbitt (2006) Automated discovery of 3D motifs for protein function annotation." *Bioinformatics* 22(6), 723-30.

#### **SPASM**

Provided by Gerard Kleywegt. This motif search tool was instrumental in calculating the scoring function for the above Genetic Algorithm.

Kleywegt, G. J. (1999). "Recognition of spatial motifs in protein structures." *J Molecular Biology* 285(4), 1887-97.



## **RIGOR**

Sister program to SPASM. Simply does the reverse search of the above. Our search feature relies on RIGOR. Refer to citation above.

## **Raster3D**

For final rendering of motif images (and GASPSdb logo!)

Merritt, Ethan A. and Bacon, David J. (1997). "Raster3D: Photorealistic Molecular Graphics" *Methods in Enzymology* 277, 505-524.

## **Molscript**

Generated ribbon descriptions for rendering by Raster3D

Kraulis, Per J. (1991). "MOLSCRIPT: A Program to Produce Both Detailed and Schematic Plots of Protein Structures." *Journal of Applied Crystallography* 24, 946-950.

## **Statistical Model**

The model described in the reference below was used to calculate the expectation values of matches to the motifs.

Stark, A., S. Sunyaev, et al. (2003). "A model for statistical significance of local similarities in structure." *J Molecular Biology* 326(5): 1307-16.

**Publishing Agreement**

*It is the policy of the University to encourage the distribution of all theses and dissertations. Copies of all UCSF theses and dissertations will be routed to the library via the Graduate Division. The library will make all theses and dissertations accessible to the public and will preserve these to the best of their abilities, in perpetuity.*

***Please sign the following statement:***

*I hereby grant permission to the Graduate Division of the University of California, San Francisco to release copies of my thesis or dissertation to the Campus Library to provide access and preservation, in whole or in part, in perpetuity.*

*Benjamin Plummer*      *June 29, 2007*  
Author Signature      Date