

Normal Mode Analysis of Macromolecular Motions in a Database Framework: Developing Mode Concentration as a Useful Classifying Statistic

W. G. Krebs,¹ Vadim Alexandrov,¹ Cyrus A. Wilson,² Nathaniel Echols,¹ Haiyuan Yu,¹ and Mark Gerstein^{1*}

¹Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, Connecticut

²Department of Biochemistry, Stanford University, Stanford, California

ABSTRACT We investigated protein motions using normal modes within a database framework, determining on a large sample the degree to which normal modes anticipate the direction of the observed motion and were useful for motions classification. As a starting point for our analysis, we identified a large number of examples of protein flexibility from a comprehensive set of structural alignments of the proteins in the PDB. Each example consisted of a pair of proteins that were considerably different in structure given their sequence similarity. On each pair, we performed geometric comparisons and adiabatic-mapping interpolations in a high-throughput pipeline, arriving at a final list of 3,814 putative motions and standardized statistics for each. We then computed the normal modes of each motion in this list, determining the linear combination of modes that best approximated the direction of the observed motion. We integrated our new motions and normal mode calculations in the Macromolecular Motions Database, through a new ranking interface at <http://molmovdb.org>. Based on the normal mode calculations and the interpolations, we identified a new statistic, mode concentration, related to the mathematical concept of information content, which describes the degree to which the direction of the observed motion can be summarized by a few modes. Using this statistic, we were able to determine the fraction of the 3,814 motions where one could anticipate the direction of the actual motion from only a few modes. We also investigated mode concentration in comparison to related statistics on combinations of normal modes and correlated it with quantities characterizing protein flexibility (e.g., maximum backbone displacement or number of mobile atoms). Finally, we evaluated the ability of mode concentration to automatically classify motions into a variety of simple categories (e.g., whether or not they are “fragment-like”), in comparison to motion statistics. This involved the application of decision trees and feature selection (particular machine-learning techniques) to training and testing sets derived from merging the “list” of motions with manually classified ones. *Proteins* 2002;48:682–695. © 2002 Wiley-Liss, Inc.

INTRODUCTION

Protein motions play a key role in a wide range of biological phenomena, including chemical concentration

regulation, signal transduction, transport of metabolites, and cellular locomotion.^{1–3} Motion is typically the way a structure actually carries out a specific function; for this reason, motions are an essential link between function and structure.

We previously developed a database of macromolecular motions,^{1,4,5} which consisted of crystallographically documented protein motions. We also developed a morph server coupled to a collection of protein “morph” movies and related statistics.⁶ Here:

1. we identify ~4,000 putative new motions from automatic structural comparison on the PDB⁷;
2. we add these to our database and present the results in a new ranking interface;
3. we analyze the dynamics of these many motions, perform normal mode analysis on them, and calculate statistics to encapsulate the results of the normal mode analysis;
4. from the normal mode analysis and the interpolations, we assemble a corpus of statistics and perform datamining and feature extraction on this corpus; and
5. we identify a number of statistics, in particular, mode concentration, that we find useful.

Our work builds upon a rich literature in macromolecular motions.^{8–11} Motion related to proteins’ mechanical function has mainly been studied experimentally by X-ray crystallography. Traditional X-ray crystallography has provided key insights into the relationships between conformational change and macromolecular function; GroEL¹² and beta-actin¹³ are just two of many examples. Progress in the field of time-resolved X-ray crystallography^{14–16} has also enhanced the study of biologically significant protein conformational change. Recently, it has become possible to

W.G. Krebs’s current address is Department of Integrative Biosciences, San Diego Supercomputer Center MC 0505, The University of California, 9500 Gilman Drive, La Jolla, CA 92093-0505.

Grant sponsor: Keck Foundation; Grant sponsor: National Science Foundation; Grant number: DBI - 9723182.

*Correspondence to: Mark Gerstein, Dept of Molecular Biophysics and Biochemistry, Yale University, PO Box 208114, New Haven, CT 06520. E-mail: Mark.Gerstein@yale.edu

Received 10 September 2001; Accepted 18 March 2002

study larger protein conformational changes via NMR.¹⁷ Other approaches have focused on the use of computational methods.^{18–25} A systematic comparison of PDB-derived difference vectors has been published elsewhere on a much smaller scale.²⁶

Normal mode analysis is a computational approach that can be applied to protein conformational change. Widely used by spectroscopists for many years to associate IR and Raman experimental peaks with small molecule vibrational modes,²⁷ advances in computer technology over the last few decades has made normal mode analysis of proteins and other large molecules practical. This was first applied to proteins in the mid 1980s and has subsequently been scaled up.^{28–34} The concept of normal mode analysis is to find a set of basis vectors (normal modes) describing the molecule's concerted atomic motion and spanning the set of all $3N - 6$ degrees of freedom. For very large molecules, it is often of more interest to try to find a small subset of these normal modes that seem in some way especially important. By modeling the interatomic bonds as springs and analyzing the protein as a large set of coupled harmonic oscillators, one can calculate a frequency of periodic motion associated with each normal mode, and then attempt to find normal modes with low frequencies. The low-frequency normal modes of proteins are thought to correspond to the large-scale real-world vibrations of the protein, and can be used to deduce significant biological properties. There is evidence to suggest^{35–40} that proper, symmetric normal mode vibration of binding pockets is crucial to correct biological activity in some proteins.

The principle of normal mode analysis is to solve an eigenvalue equation of the form

$$\ddot{\mathbf{q}} + \mathbf{F} \cdot \mathbf{q} = \mathbf{0} \quad (1)$$

where the vector \mathbf{q} is a vector representing the displacements in three dimensions of the various atoms of the molecule, and \mathbf{F} is a matrix that can be computed from the system's mass and potential energy functions. Solutions to the above system are vectors of periodic functions (the normal modes) vibrating in unison at the characteristic frequency of the mode.

Normal modes have proved to be highly useful in both modeling protein motions and in interpretation of the experimental results.^{29,32,41–53} Macromolecular motions can be often characterized by a long (nanosecond or beyond) time-scale, and it has been suggested^{54,55} that it may be possible to identify one or a few low-frequency normal modes, which would connect conformational endpoints. However, in certain cases (e.g., calmodulin motion) the amplitudes for the actual (observed) motion and the normal modes displacement vectors may differ by several orders of magnitude. For these cases, our theory may only be valid in interpretation of the motion initiation stage and in analysis of facilitating factors causing the actual motion.

In this paper we apply normal mode analysis to the study of protein motions. Fundamentally, we chose normal modes over MD and other related computational tech-

niques because normal mode analysis gives a concise description of a motion (in terms of a small number of modes) that is ideal for subsequent statistical tabulation. Also, the application of normal mode analysis techniques to $\sim 4,000$ conformational changes is much less expensive than most of the competing techniques.

In this analysis, the question we are trying to answer is to what degree the direction of the observed motion (a set of vectors connecting the structure pair) occurs along with the displacement vectors of the lowest normal modes for the initial conformation. This may indirectly provide an insight about how much protein dynamics is dominated by anharmonic contributions, even though it was not a goal of this work to develop any such quantitative anharmonicity measure. Since the structure pairs may not always be available, one of the main motivations behind this work was to see if it were possible to develop an inexpensive motion analysis technique capable of assessing the direction of the actual protein motion.

Our normal mode analyses are related to the "Essential Dynamics" (ED) methods of Berdensen^{56,57} on normal modes, involving a singular value decomposition analysis of normal mode atomic displacements and how they relate to experimentally solved conformations. Essential Dynamics can also be applied to other dynamical approaches that generate displacements including techniques that do not make a harmonic assumption such as MD simulations or experimentally determined ensembles of structures.⁵⁶ However, our analysis is in many ways formally different, and we apply it within a database framework. Many of the problems customarily found in ED analyses also apply: e.g., the superfluous rotational and translational differences must be eliminated by superimposing the experimental structures to fix at least one domain; in the process, the motion's screw-axis may be characterized.⁵⁸ Previously, we developed web software tools to solve these problems in a different way using purely experimental information.⁶ Here, we analyze a comprehensive database of thousands of putative protein motions, whereas existing publications limit their scope to single proteins or databases specific to certain types of proteins.

MATERIALS AND METHODS

Data Sources

Full outlier set

To identify a large dataset of proteins with conformational changes, Wilson et al.⁵⁹ performed automatic pairwise sequence, structure, and function comparisons on about 30,000 pairs of protein domains constructed according to scop fold classification.^{60–65} Using this set of alignments, we were able to identify $\sim 4,400$ pairs of likely protein motions. We call this set the "full outlier set" (the definitions of these terms are shown in Table IA). Its construction is described in detail in Figure 1. Basically, we plotted RMS structure alignment scores against sequence percent identity for the $\sim 30,000$ scop domain pairs aligned in Wilson et al.⁵⁹ We binned the plot into one-percent-wide bins. For each bin, we computed a mean RMS and standard deviation. Points lying more than two stan-

TABLE I. Definitions Table[†]

A. Term	Definition or URL location
Macromolecular motions database	http://bioinfo.mbb.yale.edu/MolMovDB Used for classification and annotation of motions in outlier database
SCOP database	http://scop.mrc-lmb.cam.ac.uk/scop/ Used for classification and annotation of motions via SCOP extension technique.
Wilson et al. ⁵⁹ set	As shown in Figure 1, a set of 30,000 of SCOP identifier pairs was constructed for Wilson et al. ⁵⁹ This was then separated into two sets: the 30,000 pair “Wilson et al.” set used in that paper, and the “Full Outlier Set” (described immediately below), which we use in this text. See the caption to Figure 1 for more information.
Full outlier set	Text file http://bioinfo.mbb.yale.edu/molmovdb/datasets/outliers.txt Pairs of proteins (SCOP domains) whose structural similarity score more than two standard deviations above the mean structural similarity for their sequence similarity. See the caption to Figure 1 for more information on the construction of this set.
Workable outlier set	This is the subset of the full outlier set on which both morph server processing and normal mode analysis were successful. It consists of 3,814 motion pairs.
Manual training set	This is the training set that was produced by examining the SCOP domains in the outlier set for matches against PDB IDs in the set of manually classified motions in the Database of Macromolecular Motions. ¹ Matches received the same classification as in the database, which were determined by manual examination of the specific literature. Thus, confidence in the accuracy of these classification is high.
Extended training set	The outlier set was searched for pairs that shared the same SCOP fold family as pairs classified in the Manual Training Set; these then received an identical classification. We found empirically that, because proteins that share the same SCOP fold often share similar mechanisms, proteins with the same SCOP fold have a high probability of undergoing similar conformation change and, hence, sharing the same motion size classification. Consequently, these classifications should be accurate but are less reliable than the classification in the Manual Training Set.
Classified set	This is simply the entire workable outlier set (minus those already classified in the extended training set) run through the automatic classifier defined by the decision tree, which we produced when we analyzed the extended training set.
B. Term	Definition
Mode concentration	This is discussed extensively in the text. It is a simple measure of how much the protein’s motion is concentrated into any single low-frequency normal mode.
No. of CAatoms	Number of C-alpha atoms in the protein
Residuals	This is the Euclidean length of the residual difference between the atomic displacements between protein pairs and the SVD fit of the normal modes to the atomic displacements (in Angstroms)
Norm0	Maximum Value of the SVD displacement vector (unitless)
Norm1	Mean of the SVD displacement vector (unitless)
Norm2	Root-mean-square of the SVD displacement vector (unitless)
Frequency	The frequency in relative units of the normal mode with the highest SVD coefficient.
Ranking overlap	Rank of the normal mode with the largest overlap (unitless). Overlap is defined in the caption to Figure 2.
Maximum overlap	Value of the largest overlap (unitless quantity). Overlap is defined in the caption to Figure 2.
Size of 2nd core	This is the number of residues in the 2nd core (the 2ndCoreCAs key in the database). This is typically related to the size of the protein, although in poorly matched protein pairs the number can be less.
Trimmed RMS	This is the trimmed RMS score, as defined in Wilson et al. ⁵⁹ and Gerstein and Krebs. ¹
Maximum CA movement	This is the largest movement (in Angstroms) of any residue during the course of the motion, as computed by the Morph Server.
Number of atoms	This is the number of atoms in the protein as computed by the Morph Server. (Atoms in non-standard amino acids are excluded). This is a measure of the size of the protein.
Energy of frames	The Morph Server computes energies for the various intermediate structures. These show a strong relationship to the sequence similarity between the two structures, and are indicators of how “good” a given morph is. The relationship of intermediate energies (energy of 4th frame, for example) with endpoint frames (energy of 8th frame, for example) can sometimes provide a rough sense of activation energies.
Translation	In hinge motions, the approximate translation (in Angstroms) the moving domains undergoes in the course of the motion, as automatically computed by the morph server. (This number is also computed for non-hinge motions, where it is less meaningful.)
Hinge rotation	In hinge motions, the rotation (in degrees) of the moving domain around the screw axis in the course of the motion, as automatically computed by the morph server. (This number tends to be small in non-hinge motions.)
Number of hinges	The number of putative hinges, or flexible linkages involved in the motion, as determined by the Morph Server
Traditional RMS	This is simply the traditional RMS score between the domains.
Rank of Norm0 mode	This is a software index that identifies the normal mode contributing the most to the motion as computed within our SVD framework. (The same normal mode that sets norm 0.)

[†]Section A lists the various data sources used in this paper, giving the location of each, along with a brief explanation of its use or importance. Section B lists definitions of the key statistics and other terms used in subsequent tables as well as in the text of the paper.

RMS vs Sequence Identity

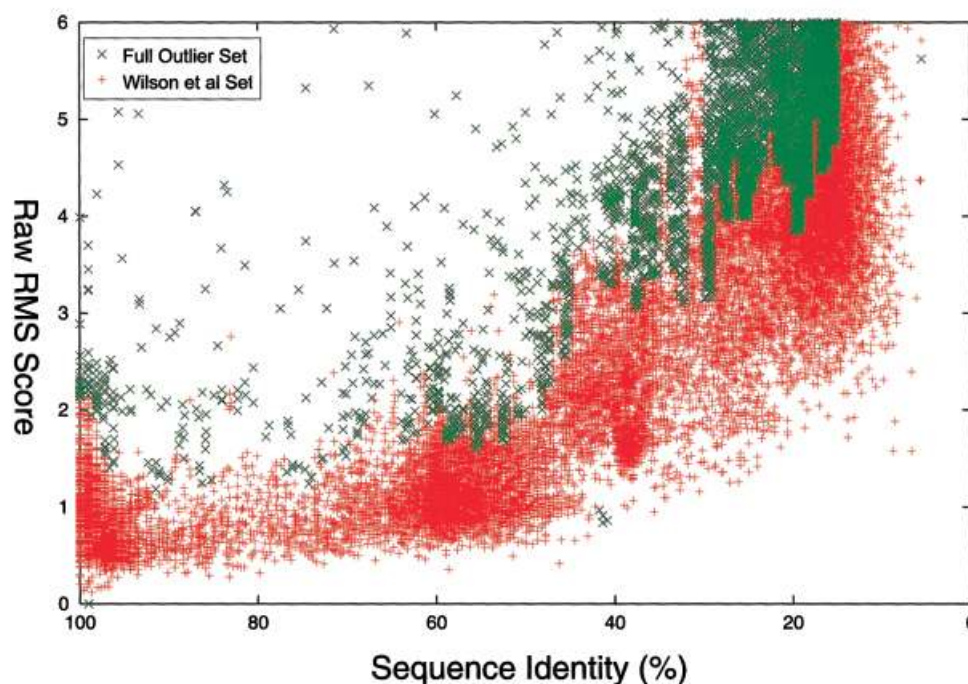


Fig. 1. Construction of full outlier set. The crosses on this page illustrate motion pairs plotted in terms of RMS structure alignment scores against sequence percent identity for the 30,000 SCOP domain pairs Wilson et al.⁵⁹ identified from the PDB. Data points were binned into one-percent-wide bins, and the mean RMS and standard deviation in each one-percent-bin was computed. Points more than two standard deviations above the mean were removed from the original 30,000 pair dataset (red crosses) and used to compose the full outlier set (green crosses), which ultimately consisted of 4,400 such pairs.

standard deviations above the mean were removed from the dataset and used to generate a new dataset, the full outlier dataset, which ultimately consisted of 4,400 such pairs. This set is intended as a comprehensive sample of protein flexibility in the PDB.

Workable outlier set

We ran the full outlier set through our protein morphing server.⁶ We placed the resulting database of pre-processed PDB files, morph statistics, and movies, on the World Wide Web, organized by their scop fold classification. The new automated approach was able to process and generate several thousand new morph movies. As described below, the morph server acted as a filter, eliminating about 600 pairs in the full outlier set that corresponded to non-physical motions. Next, we applied the normal mode analysis described below on the successfully morphed pairs, to produce a set of about 3,800 motion pairs, the “outlier set”. In this paper, we concentrate exclusively on this new “workable outlier set” data.

Manual set

In order to perform feature analysis, we classified two subsets of the workable outlier set (the “manual set” and the “extended set”) into the classification schema of the Database of Macromolecular Motions¹ (“fragment,” “domain,” “subunit,” “complex” on the basis of size and

“hinge,” “shear,” “neither hinge nor shear,” and “unclassifiable” on the basis of packing). Further details about this classification may be found in our previous paper.¹

For the “manual set,” we performed a database merge of the “outlier set” against the previously published set of manually classified motions in the Database of Macromolecular Motions,¹ the “1998 motions.” The PDB identifiers in each motion pair in the outlier set were checked for matches against the PDB identifiers associated with the 1998 motions. When a match was found (meaning the protein had been manually classified), the motion pair was given the same classification as its constituent protein had been given in the database. Two hundred and forty-five motion pairs met this criterion and were classified accordingly. Classifications in this manual training are expected to be accurate. (There was, however, one issue in applying this merge: GroEL is classified both as a subunit and a fragment motion. Because the Morph server analyzes single domains, not entire subunits, the fragment classification was used in this isolated case.)

Extended set

To enlarge the training data for the supervised machine learning analysis, we constructed a second, larger training set (the “extended set”). For a variety of physical reasons, proteins sharing the same fold family generally share a similar motion classification—in particular, we have ob-

served this in our manual surveys of motions.^{1,6,60,66,67} Consequently, we constructed this set under the assumption that domains sharing a fold usually share a motion classification. The outlier set is constructed in such a way that both pairs always belong to the same fold family. It was, therefore, necessary only to determine the scop fold classification^{60,65} for each of the 245 motion pairs in manual training set and then assign the classification in the manual set to the entire scop fold family. Pairs in the outlier set belonging to this scop fold family then simply received the family's classification. In this way we identified a set of 1,670 motions, which we call the "extended training set." This set of classifications, although potentially less accurate than the manual training set, is still quite useful.

Preprocessing With Morph Server

We analyzed 3,814 proteins using this method from the full outlier set. Previously,⁶ we modified the X-PLOR package⁶⁸ to homogenize the stored coordinates, a non-trivial problem.^{69,70} Filling-in of missing non-hydrogen coordinates was necessary for the energy minimization subsystems to work robustly with a large number of PDB files and ensured consistent numbering of atoms so the PDB files for the starting and ending conformations had to be pre-processed ("homogenized") by the Morph Server.⁶ Only pairs of protein conformations for which the Morph Server had successfully produced a movie were considered; this had the effect of filtering out pairs unlikely to involve a true motion, although no doubt some pairs that did not represent a true biological motion nevertheless did generate a plausible morph. The Morph Server also removes overall rotation and translation motions from the input structure.

High-Throughput Normal Mode Analysis of the Outlier Set

We used MMTK⁷¹ to carry out normal mode analysis on the pre-processed PDB file pairs. The numerical Python module⁷² made the linear algebra computations. A master Perl⁷³ script fed database information to the slave Python MMTK module. The results reported here were performed by computing the normal modes of the starting structures in each pair. Reversing the calculations by computing the normal modes of the ending structures did not appreciably alter the results.

Finding the normal modes themselves dominated the time and memory requirements of our analyses. In order to process the larger proteins in our database, we approximated each residue as a single, virtual atom centered at its C- α coordinate and selected the corresponding standard force field in MMTK.⁷¹ This made the memory requirements of the normal mode analysis tractable on our systems. To further accelerate the computations, we restricted MMTK to compute only the twenty lowest-frequency normal modes.

We used the MMTK deformation force field model. In this model, the energy is computed as the difference between some displaced model and the experimental structure using the formula:

$$E_1 = \frac{1}{2} \sum_{j=1}^N k(\mathbf{R}_{ij}^{(0)}) [|\mathbf{R}_{ij}^{(0)} + \mathbf{d}_i - \mathbf{d}_j| - |\mathbf{R}_{ij}^{(0)}|]^2 \quad (2)$$

where k is a constant, $\mathbf{R}_{ij}^{(0)}$ is the vector from atom i to atom j in the experimental structure, \mathbf{d}_i is the vector between the atom i in the displaced structure and the same atom in the ground-state experimental structure.

Each calculation averaged 20 seconds per protein pair on a 450-Mhz Pentium III processor with 0.7 Gigabytes of RAM running the Red Hat Linux operating system. An average analysis took about 100 Megabytes of memory to invert the matrix.

Theoretical Approach for Analysis of Normal Mode Statistics

We computed a number of key statistics on the normal modes (Table IB), which we describe here.

Analysis of observed motion

The lowest frequency normal modes determined by Normal Mode Analysis may be represented as an $m \times n$ matrix A , where m is three times the number of atoms in the system (one entry for each Cartesian axis), and n is the number of normal modes of interest. In this paper, n is twenty.

Imagine a vector $\bar{\mathbf{v}}$ of length n , specifying some interesting linear combination of normal modes. Then $A\bar{\mathbf{v}}$ is a vector of length m , representing a trajectory of atoms. If we let the vectors c_i and c_f be the vectors of length m giving the positions of the $m/3$ atoms in conformations C_i (starting) and C_f (ending), respectively. We determined these from our database of motion, which has such data, chiefly derived from experimental sources such as X-ray crystallography.

If we now define a new vector $\mathbf{b} = c_f - c_i$, or the differences between the ending and starting positions of each of the atoms of the structure along all three Cartesian axes, then we can find optimal \mathbf{v} so that

$$A\bar{\mathbf{v}} = \mathbf{b} \quad (3)$$

In the normal case where $\dim \bar{\mathbf{v}} < 3N - 6$, this represents an over-determined system of linear equations, and may be solved by an appropriate numerical technique for solving linear least squares, such as Single Value Decomposition (SVD).⁷⁴ In practice, this is a very quick calculation, nearly instantaneous to the user.

Analytic Measures

Overlap of each mode with direction of motion

For every motion pair, we computed the overlap of each normal mode against the vectors giving the differences between the structures corresponding to the motions. For one particular atom, we define the "overlap" O_{ij} as the cosine of the angle between the mode and the direction of motion,

$$O_{ij} = \frac{\bar{\mathbf{b}}_i \cdot \bar{\mathbf{f}}_{ij}}{|\bar{\mathbf{b}}_i| \cdot |\bar{\mathbf{f}}_{ij}|} \quad (4)$$

In the above formula O_{ij} is represented as a normalized dot product between some reference vector $\bar{\mathbf{b}}_i$ (in this case,

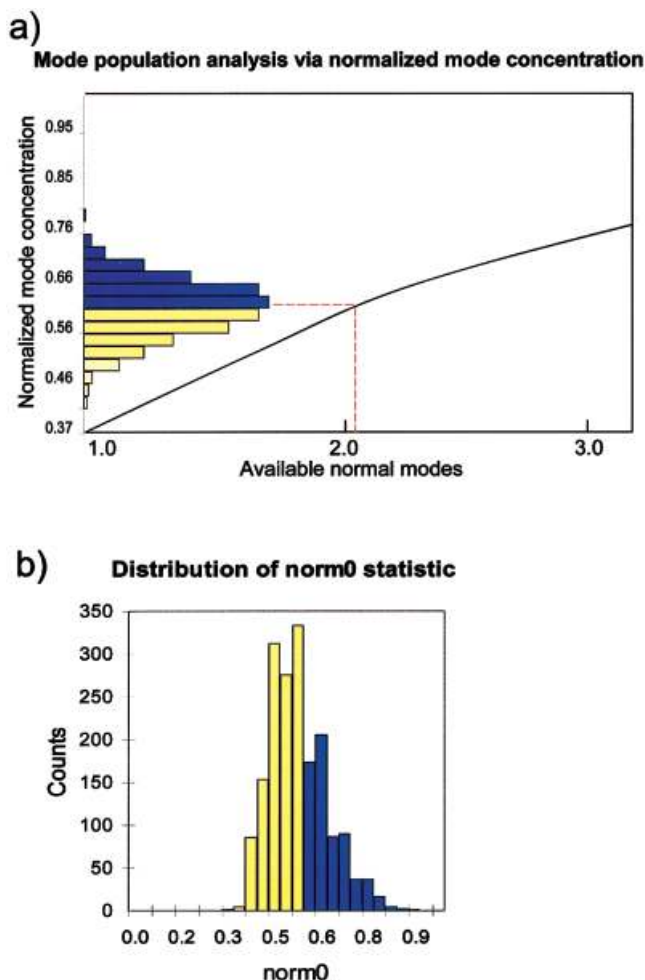


Fig. 2. **a**: Analysis of the normalized mode-concentration statistic to assess the normal modes populations. The center of the normalized mode-concentration histogram is traced to the number of available states (modes) using the Boltzmann logarithmic dependence relation. **b**: Histogram of norm0 statistic calculated over all entries in our database. The plot clearly shows that the large contributions (over 50%) from a single normal mode are not uncommon.

the displacement between the PDB structures of the motion pair in question) and \mathbf{f}_{ij} , the j th normal mode displacement vector for the same atom.

For the ensemble of atoms in a structure, we can define “average overlap” O_j as the mean overlap averaged over all N atoms in the structure, i.e.,

$$O_j \equiv \frac{1}{n} \sum_{i=1}^n O_{ij}. \quad (5)$$

We can also calculate an average absolute value of the cosine $1/n \sum_{i=1}^n |O_{ij}|$, which provides a quantitative measure of the first-order overall deviation for a particular normal mode from the observed motion. The larger values of this quantity indicate that a given mode’s atomic displacement vectors are more similar in directionality to the vectors giving the differences between the PDB files. The mode of “maximum overlap” is the mode with the

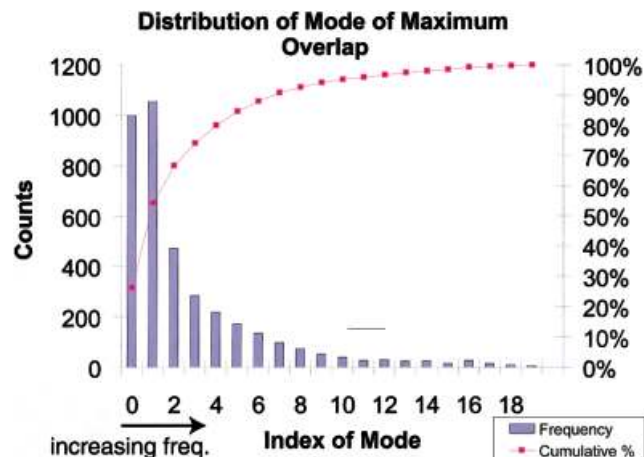


Fig. 3. Our software places the twenty lowest-frequency normal modes in an array, thereby assigning each normal mode an index, from zero to nineteen. Increasing index numbers identify higher-frequency normal modes. We computed the overlap of each normal mode and recorded the index of the normal mode of greatest overlap. We plotted the number of times each index had greatest overlap in this histogram.

greatest “absolute average overlap” and most matches the protein motion’s directionality.

S-correlation

A means of quantifying the similarity of the displacement between the PDB structures and the normal mode displacement vectors can be also achieved by calculating the following quantity,

$$s = \sqrt{\sum_{j=1}^n j^2 O_j^2 - \left(\sum_{j=1}^n j O_j \right)^2} \quad (6)$$

where O_j is defined as above. This formula*, directly adapted from Hinsen’s work³⁹ with a lowering of dimensionality, gives the s-correlation between the reference vector and the set of normal mode displacement vectors. This may be used to provide an overall quantitative measure of the similarity in directionality between the observed displacements and those of the various normal modes. Thus, the convention used to number the modes does not affect s-correlation in a meaningful way.

In the present work, we also utilize an interesting mathematical property of this statistic: its positive definite values imply that the displacement vectors from only the lowest two normal modes may coincide with the direction of the observed motion.

Mode concentration

Based on the fit of the modes to the observed motion, we calculate a number of statistics that show the degree to

*This formula is identical to the one in Hinsen’s work. However, one may also find useful a related statistic, a direct analog of the second order momentum:

$$\sqrt{\sum_{j=1}^n j^2 O_j^2 - \left(\sum_{j=1}^n j O_j \right)^2}$$

Maximum Overlap vs. Size

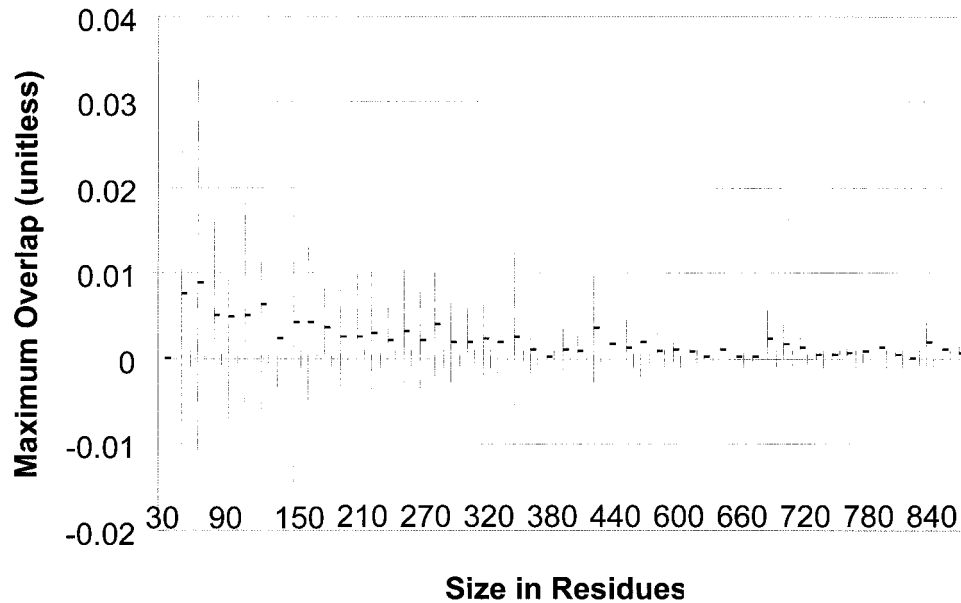


Fig. 4. Relationship between protein size and maximum overlap. To make the effect clearer, the y-values were binned into groups of 15 residues. The mean and standard deviation were computed for the values in each bin, with the results plotted. Each heavy horizontal bar indicates the mean in each bin, while the vertical bars indicate two standard deviations above and below the mean.

Frequency of Max Overlap vs. Size

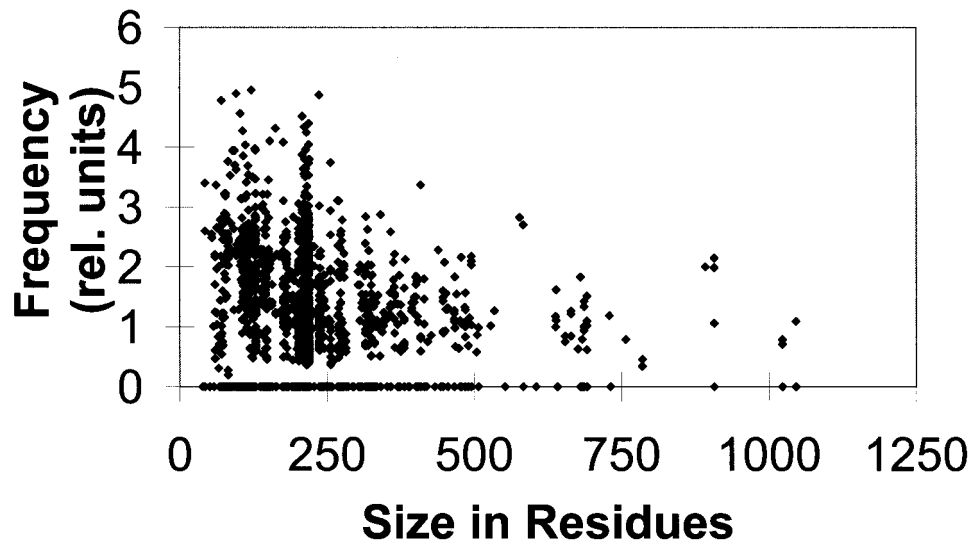


Fig. 5. Correlation between the frequency of the mode of maximum overlap and protein size.

which the fit is dominated by a single mode. We define norm zero (“norm0”) as simply the weight of the largest component (i.e., the largest value in the vector \mathbf{v}), the one norm as the average component (“norm1”), and the two norm as simply the Euclidean mean (“norm2”) of the component’s weight.

All of these statistics give a measure of the degree to which the vector \mathbf{v} is dominated by a single component. In somewhat more sophisticated fashion, we can measure this using information theory approaches.

In coding theory, information content is related to the negative entropy of a physical system. It specifies how much

TABLE II. Summary of New Statistics Added to Morph Server[†]

Key	No. of CAatoms	Residuals	Norm1	Norm2	Frequency	Ranking overlap	Maximum overlap
Mean	220	480	-0.001	540	3.1	2.7	0.0031
Std. dev.	110	660	0.051	360	0.89	3.6	0.005
Minimum	39	0.23	-0.14	15	4.2E-08	0	4.7E-5
Maximum	1,000	8,800	0.15	2,700	8.6	19	0.11
Median	210	330	0.00093	520	3.1	1	0.0017

[†]This table presents mean, standard deviation, minimum, maximum, and median values for the new statistics that were added to the database following normal mode analysis of approximately 3,800 motion pairs in the database. The statistics are defined in Table IB.

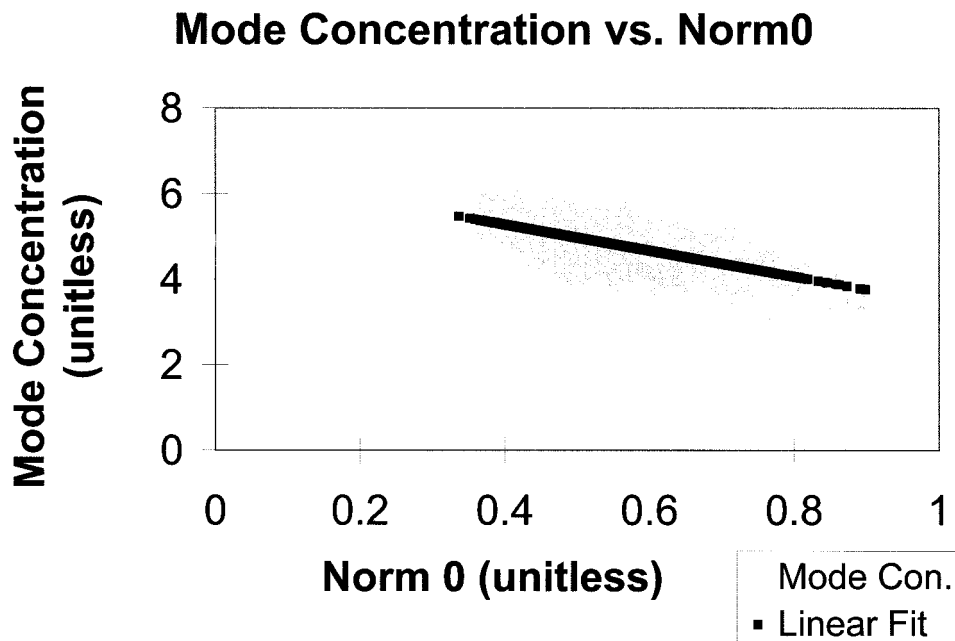


Fig. 6. Relationship between mode concentration and norm0 (concentration of motion in the mode with greatest concentration).

information is stored in a given set of numbers, and is typically used to compare the efficiencies of compression techniques. Therefore, once \bar{v} has been obtained, a statistic may be computed to summarize the information contained in the vector \bar{v} :

$$I = \sum_{i=1}^n -|v_i| \ln |v_i| \quad (7)$$

This statistic specifies how much movement is concentrated in any given mode, hence its name, “mode concentration.”

We can normalize I to unity by dividing it over its maximal value, corresponding to the uniform movement distribution over all available modes, and obtain the “percentage mode concentration” statistic \bar{I} , that specifies the degree to which a given motion is localized within a few modes relative to the uniform distribution (maximal disorder). As mentioned above, one can also directly relate information content (and, thus, also our normalized information content) to the well-known Boltzmann formula $S = k \ln N$ for the entropy (measure of the system

disorder in statistical mechanics) expressed through the number of states N available to the system, i.e.,

$$\bar{I} \sim \ln N \quad (8)$$

The normalization ensures that \bar{I} approaches zero if all movement is concentrated in only one normal mode ($N = 1$), whereas the value of $\bar{I} = 1$ corresponds to the even distribution of motion over all available normal modes (i.e., to the maximal value of I computed from Eq. (7)).

RESULTS

Application of These Statistics to the Outlier Dataset

Figures 2 through 6 illustrate some properties of the above statistics on the outlier dataset.

Figure 2 shows distributions of the normalized mode concentration and norm0 statistics. Using the logarithmic dependence Eq.(8) of the normalized mode concentration with respect to the number of available modes, one can arrive at the number of most heavily involved modes. This would be the value of N , for which the value of \bar{I} is most frequently observed. The observed peak in the normalized

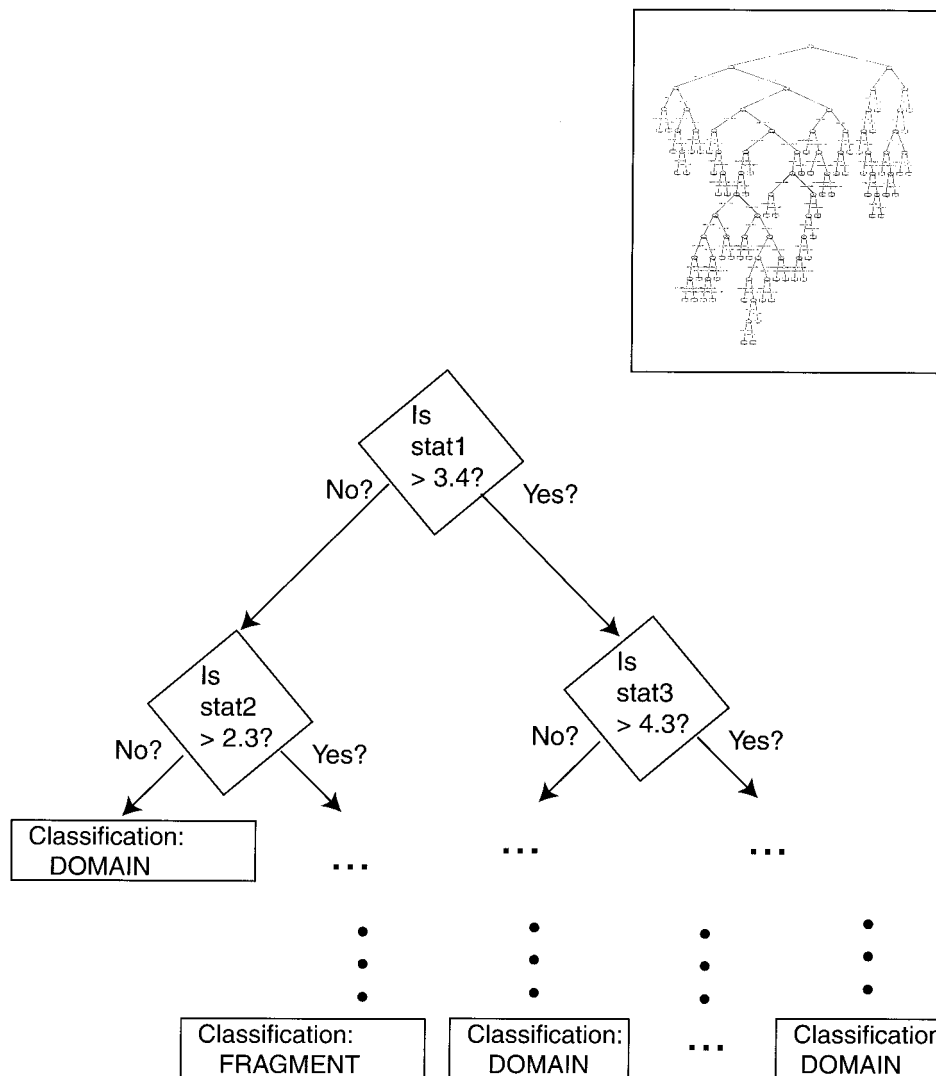


Fig. 7. Decision tree concepts. Two decision trees (not shown here) were generated by S-Plus (MathSoft, Inc.) using default parameters from the 245-element manual training set and the 1,670-element extend training set (defined in Table IA). These trees classify motions as "fragment," "domain," or "subunit." The decision tree associated with the extended training set defined an automatic classifier (implemented in Perl by examination of the tree) that produced the "classified set." This figure shows the conceptual operation of decision trees. At each node, the classifier chooses either the right or left branch, respectively, depending on whether or not the node's associated statistic is greater than the value associated with the node. **Inset:** Structure of an actual decision tree is shown in miniature. The classifier follows the decision tree until it reaches one of the terminal leaves, where a classification is made. A "training set" providing a set of examples and associated "correct" classifications is run through the S-Plus program, which generates a decision tree that can classify the training set correctly.

mode-concentration histogram at 0.6 [Fig. 2(a)] suggests that the actual direction of the motion lies most often along the direction of two modes. Analysis of norm0 histogram [Fig. 2(b)] further confirms this finding: the most commonly observed weight of the major contributing mode lies within the range 0.5–0.6 (i.e., there is usually one mode that dominates the motion fit) whereas the normal mode approximations with values of norm0 below 0.4 are quite rare (the latter would imply that there are usually more than two mostly contributing normal modes exist for each normal mode fit).

Figure 3 shows that most often the low-frequency modes tend to be the ones with maximum overlap with the actual direction of motion (Fig. 3). There is also a relationship between protein size (measured in number of residues), mode frequency, and maximum overlap (Figs. 4 and 5).

Protein size (measured in number of residues) is negatively correlated to maximum overlap (Fig. 4). Larger proteins have additional fragments that can be involved in a motion and, hence, additional degrees of freedom, decreasing the overlap between the tested normal modes and the observed motion. (An alternative explanation for this

TABLE III. Comparison of the Percentages and Absolute Counts of Domain, Fragment, and Subunit Motions in the Classified, Extended, and Manual Training Sets[†]

Motion size	Predicted		Observed			
	Classified set		Extended set		Manual set	
	Count	Percent	Count	Percent	Count	Percent
Domain	2,165	95	1549	93	180	73
Fragment	94	4	107	6	50	20
Subunit	14	1	14	1	15	6
Totals	2,273	100	1670	100	245	100

[†]Definitions of the different sets in the header are given in the text as well as Table IA. “Count” gives the number of times the particular motion size classification (Domain, Fragment, and Subunit) occurs in that dataset. “Percent” is the percentage out of the total number (“Total”) of domain, fragment, and subunit motions in the dataset. The two columns on the left for the auto-classified set (“count” and “percent”) represent a prediction made by an auto-classifier; the remaining columns represent observations.

TABLE IV. [†] Automatic Ranking of Database Statistics via Decision Tree Feature Extraction

Database statistic	Depth in tree built upon extended set	Depth in tree built upon manual set
Size of 2nd core	1	1
Trimmed RMS	3	2
Maximum CA movement	5	2
Number of atoms	4	3
Mode concentration	6	4
Energy of 2nd frame	6	4
Translation	4	5
Hinge rotation (degrees)	4	6
Number of hinges		6
Energy of 3rd frame		6
Norm0 (maximum value)	5	9
Energy of 9th frame	3	
Number of residues	5	
Frequency	5	
Residuals	6	
Norm1 (average norm)	6	
Rank of Norm0 mode	7	
Traditional RMS	8	
Norm2 (Euclidean norm)	8	
Energy of 4th frame	9	
Energy of 9th frame	9	
Energy of 8th frame	13	

[†]This table indicates the earliest depth of the supervised machine learning decision tree each statistic first occurs, thus quantifying the relevance of each statistic to the particular motion property at hand (“fragment,” “domain,” or “subunit” motion, in this case).

observation is that the various approximations used in normal modes approximation work less well for larger proteins.) Maximum overlap decreases with protein size, but the effect is not dramatic, so it should be possible to design a standard analysis that works well on proteins comparable to those in our database.

Increasing protein size (in residues) corresponds to modes of maximum overlap of decreasing frequency (Fig. 5). A standard analysis concerned with larger proteins may need to consider more low-frequency normal modes than would suffice for smaller proteins. It would be

desirable, given a protein of specific size, to deduce a frequency cut-off value, above which normal modes could be expected to be less useful in an analysis of motion. Analyses of individual proteins in the literature support the existence of such a cutoff^{46,75} showing a slight dependency on the force field used. Our results show that it is possible to determine such a cut-off frequency statistically from our database (Fig. 5) and thereby empirically deduce a reasonable number of normal modes to use in a given type of analysis. Researchers using an identical force field to the one used in this study may consult Figure 5 directly to determine the appropriate cut-off for their particular protein; researchers using slightly different force fields or dynamical methods may wish to obtain access to the database to compute a cut-off value appropriate for their specific dynamical analysis.

Validation of Mode Concentration With Feature Extraction Techniques

The physical and information theory basis of the mode concentration statistic suggested it might be useful in classification problems. Subsequent analysis via machine learning techniques (below) supports this.

Artificial intelligence feature analysis techniques, particularly supervised machine learning, provide one way of validating the usefulness of our mode concentration statistic. In general, the concept of supervised machine learning is that the system is “taught” to classify a given set of inputs by being given a “training set” that matches a sample set of inputs to a correct set of outputs.⁷⁶

As described above, we created the manual and extended data sets as training sets to perform feature analysis. Using supervised machine learning techniques,^{76,77} we constructed two decision trees in S-Plus (MathSoft, Inc.) using the software’s default parameters^{77–79} (one for each of the two training sets) to classify the statistics, including the new ones (Table II), in the morph server.⁶ The use of S-Plus to construct decision trees from a specific training data set is a straightforward operation.

Decision trees, a form of supervised machine learning, attempt to partition the examples in the training set based

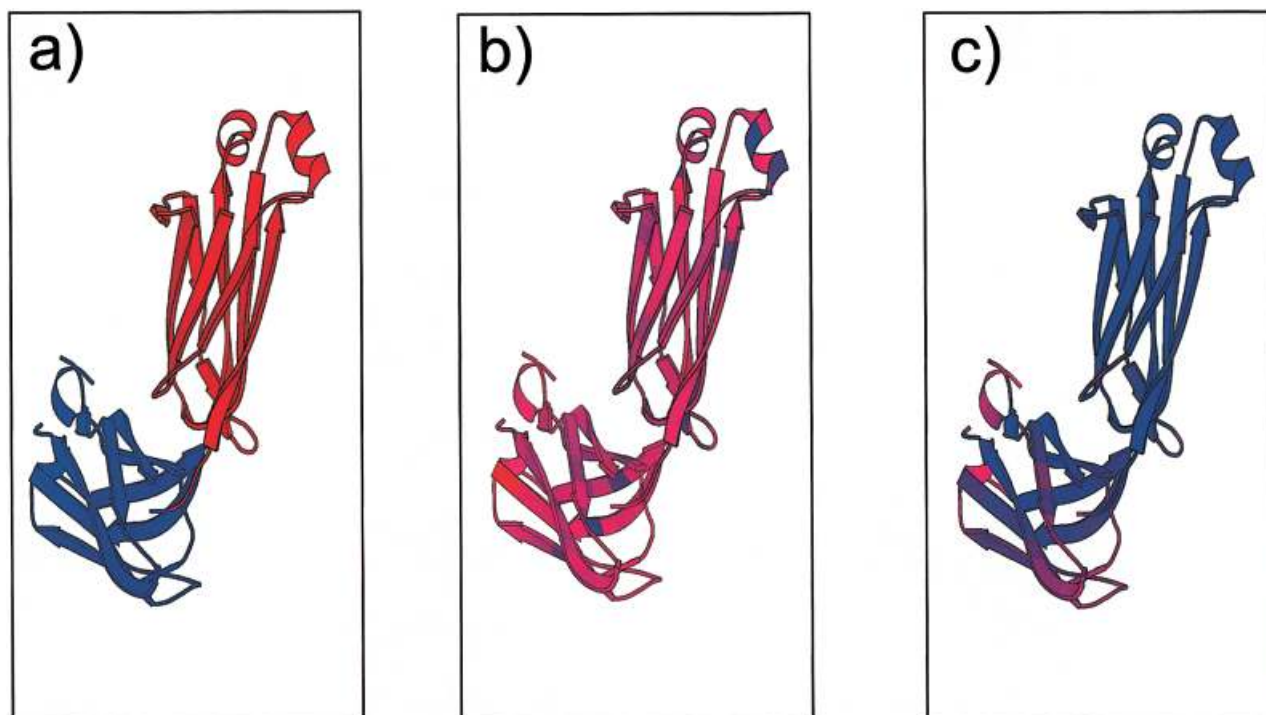


Fig. 8. Output of new set of Web tools associated with normal mode analysis that the user may request on any protein for which a PDB structure file is available. The URL for this server is <http://www.molmovdb.org>; these features may be accessed by browsing to a specific movie and selecting one of these analyses from the menu. **a:** The parts of the protein that actually move, as calculated from comparison of the starting and ending PDB structures for the motion. Areas that move are colored in red, while areas that remain stationary are colored in blue. The user may compare these three panels to deduce structural information. Hinge locations involved in the motion may be deduced, as these are highly flexible regions (as identified by a and b) located near the moving domains (show in red in c). The specific protein example shown is that of an immunoglobulin elbow joint motion (morph ID d2fb411-dlafv11). **b:** Performs a normal mode flexibility analysis on the structure. Regions that are more flexible are colored in red, while less flexible regions are colored in blue. **c:** Similar information, using experimental temperature factors supplied in the PDB file, if available.

on the values of individual statistics (Fig. 7). In the actual decision tree, each statistic used in the classification decision appears in at least one branch junction. Features more relevant to the classification problem tend to appear earlier in the decision-making process, corresponding to a higher-level branch in the trees. By recording the depth any statistic first appears in the decision tree, decision trees may be used for feature analysis (Table III). Mode concentration ranks prominently with a low depth, indicating that it appears high in the tree and is, therefore, useful for classifying motions.

Using appropriate, simple physical and mathematical concepts (normal mode analysis, singular value decomposition), we postulated several statistics (mode concentration and the various analytic norm measures) and confirmed our initial hypotheses using artificial intelligence techniques. These culled the morph server's⁶ output of 36 physically-motivated statistics down to a set of nine "essential" statistics that proved most useful in this particular classification problem (Table IV), which agree roughly with our own sense of the statistics most related to motion size. Similar databases of heterogeneous biological statistics may be "distilled" from a larger body of experimental data with these and similar techniques. In this case, the automatic classification features of the decision trees are only a side benefit. Feature analysis confirmed our earlier

intuition that mode concentration can be useful for classifying motions.

Depending on the supervised machine learning technique used (decision trees), larger training sets can sometimes produce a more accurate automatic classifier than a smaller classifier. For this reason, it is possible that an automatic classifier produced from the larger extended training set may classify more accurately than one produced from the smaller, more accurate manual set, although this may seem counterintuitive. Comparing the results produced by the manual and the extended training sets thus will serve as a useful check.

Web and Database Integration

We used the results of our decision tree analysis (Table III) to improve the ordering and presentation of statistics in Macromolecular Motions. Database web reports (<http://molmovdb.org>). In addition, a new web tool (Fig. 8) on this site graphically depicts output from the normal mode analysis as well as older experimental information.

The new data from normal mode analysis have been integrated into both the Macromolecular Motions Database and the Partslist Database (<http://partslist.org/>) as well.⁸⁰ This allows comparison by fold of motion and other data by a number of techniques, including regression analysis. Interactive users can test a number of statistics

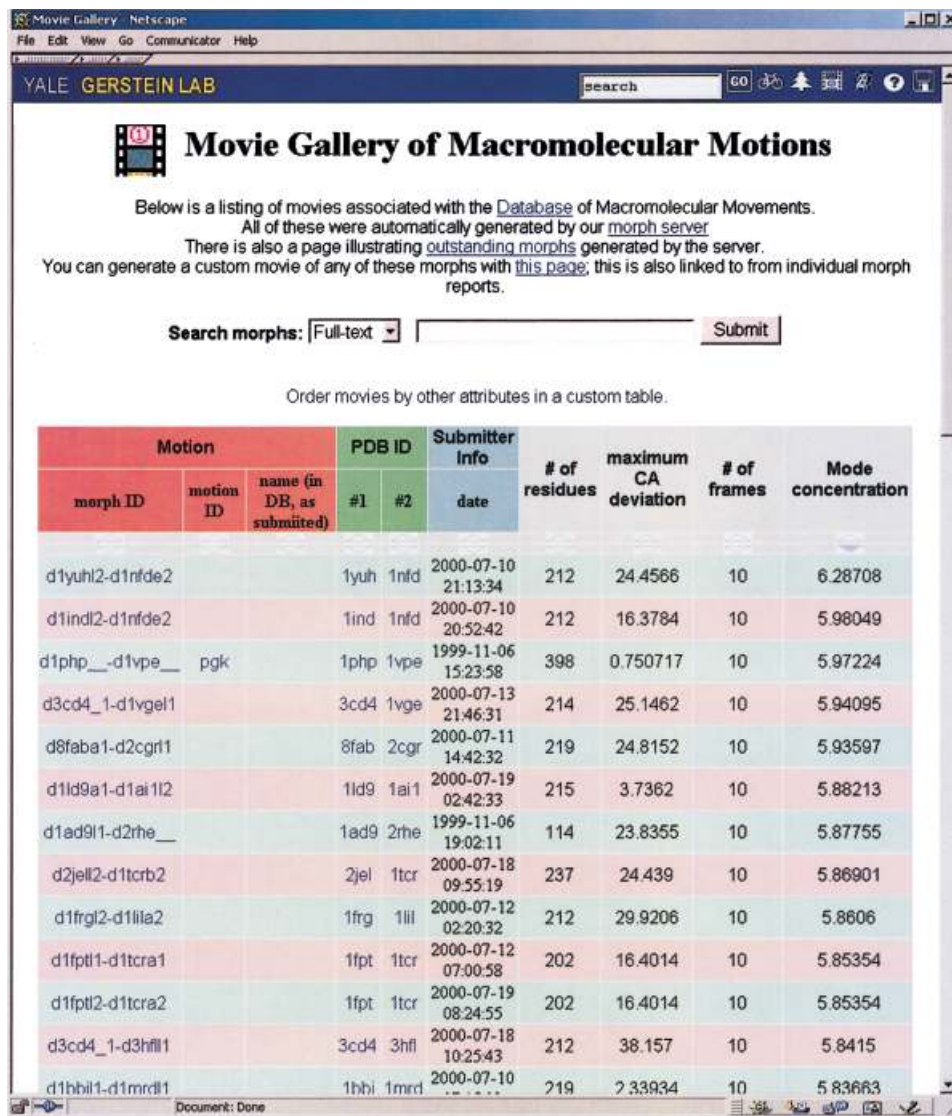


Fig. 9. Screenshot of the Movie Gallery web page. This shows a Movie Gallery page in the Macromolecular Motions Database that ranks different motions according to the average mode concentration.

for correlation against the new data, as well as identify outlying folds that do not maintain the normal regression pattern by mouse over. Figure 9 gives a screen shot of motions ranked by average mode concentration in the Movie Gallery page of the Macromolecular Motions Database, which will show the animation of the corresponding motion on click.

DISCUSSION

Comprehensive structural studies within a database framework, such as the one described here, can complement more traditional computational studies of single molecules in a number of ways. The most immediate benefit is that a database study makes more data available to researchers, and can sometimes make more general statements about trends and patterns in the results than would be possible from similar studies on a smaller sample

of macromolecules. A disadvantage of the type of database study performed here is that they require greater computational resources than equivalent studies on single macromolecules. Also, the implementation of automatic methods to handle a large class of macromolecular may require somewhat greater algorithmic sophistication since steps requiring manual processing are less desirable when dealing with a large number of structures.

Researchers who have developed their own, novel computational structural studies may expand their computations from analyses of single molecules to a comprehensive study of an entire structural database, such as the Database of Macromolecular Motions. The results of such structural studies constitute databases in their own right. Artificial intelligence techniques can then be applied to such derived databases to append additional, useful statistics, "distill" a derived database down to a set of "essential"

statistics, as well as construct automatic classifiers. This has obvious practical applications; e.g., pharmaceutical companies might mine existing biological databases and apply existing or new algorithmic techniques (e.g., variants on normal mode analysis) to generate derived databases describing potential drug targets within a statistical framework. Artificial intelligence techniques can be used to extract key features and empirically assess the validity of new statistical models.

CONCLUSIONS

We have developed a framework that allows for a statistical study, in combination with our Database of Macromolecular Motions, of the importance of normal mode vibrations in biologically significant macromolecular motions. A statistic calculated from our analysis of normal mode displacements, mode concentration, is corroborated by feature selection as a useful statistic in classification. Feature selection techniques can be used to "summarize" databases of experimentally derived statistics into an especially salient set of "essential" statistics.

Examining the relationship between the aggregate directionality of the normal modes and structures' conformational change through a statistic such as mode concentration can be used to classify the motion ("fragment," "domain," or "subunit"). Normal modes have already been used⁵⁸ to identify dynamic protein domains. An analysis of the distribution of low-frequency normal mode trajectories should provide information about the type of protein motion and size of the domains involved in the motion. Our data empirically support earlier results⁴⁶ that analysis of only a small number of low-frequency modes should suffice for qualitative analysis of protein dynamics. The database can also be used to determine statistically the cut-off for normal modes computed using different force fields.

In addition to being made available through the Macromolecular Motions Database, our new data sets are integrated into the external Partslist database.⁸⁰ We have provided additional web tools associated with this paper that allow molecular biologists to perform flexibility analysis on structures with putative motions, thereby identify key residues involved in the motion, and compare the results with similar analysis on the over 4,000 new motions now available in the database.

ACKNOWLEDGMENTS

We thank Dr. Yuval Kluger for his help with machine learning, and Dr. Jiang Qian for the Partslist integration of the data. Numerous people have also either contributed entries or information to the database and morph server or have given us feedback on what the user community wants. The authors also thank Informix Software, Inc., for providing a grant of its database software.

REFERENCES

- Gerstein M, Krebs W. A database of macromolecular motions. *Nucleic Acids Res* 1998;26:4280–4290.
- Debrunner PG, Frauenfelder H. Dynamics of proteins. *Annual Rev of Phys Chem* 1982;33:283.
- Lipscomb WN. Acceleration of reactions by enzymes. *Accounts Chem Res* 1982;15:232.
- Gerstein MB, Jansen R, Johnson T, Park B, Krebs W. Studying macromolecular motions in a database framework: from structure to sequence. In Thorpe MF, and Duxbury PM., editors. *Rigidity theory and applications*. Kluwer Academic/Plenum press; 1999. p. 401–442.
- Gerstein M. A protein motions database. *PDB Q Newsletter* 1995;73:2 (July).
- Krebs WG, Gerstein M. The morph server: a standardized system for analyzing and visualizing macromolecular motions in a database framework. *Nucleic Acids Res* 2000;28:1665–1675.
- Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The protein data bank. *Nucleic Acids Res* 2000;28:235–242.
- Israelewitz B, Gao M, Schulten K. Steered molecular dynamics and mechanical functions of proteins. *Curr Opin Struct Biol* 2001;11:224–230.
- Young M, Kirshenbaum K, Dill KA, Highsmith S. Predicting conformational switches in proteins. *Protein Sci* 1999;8:1752–1764.
- Shaknovich R, Shue G, Kohtz DS. Conformational activation of a basic helix-loop-helix protein (MyoD1) by the C-terminal region of murine HSP90 (HSP84). *Mol Cell Biol* 1992;12:5059–5068.
- Dixon MM, Nicholson H, Shewchuk L, Baase WA, Matthews BW. Structure of a hinge-bending bacteriophage T4 lysozyme mutant, Ile3→Pro. *J Mol Biol* 1992;227:917–933.
- Xu Z, Sigler PB. GroEL/GroES: structure and function of a two-stroke folding machine. *J Struct Biol* 1998;124:129–141.
- Chik JK, Lindberg U, Schutt CE. The structure of an open state of beta-actin at 2.65 Å resolution. *J Mol Biol* 1996;263:607–623.
- Oka T, Yagi N, Fujisawa T, Kamikubo H, Tokunaga F, Kataoka M. Time-resolved x-ray diffraction reveals multiple conformations in the M-N transition of the bacteriorhodopsin photocycle. *Proc Natl Acad Sci USA* 2000;97:14278–14282.
- Genick UK, Borgstahl GE, Ng K, Ren Z, Pradervand C, Burke PM, Srajer V, Teng TY, Schildkamp W, McRee DE, Moffat K, Getzoff ED. Structure of a protein photocycle intermediate by millisecond time-resolved crystallography. *Science* 1997;275:1471–1475.
- Schlichting I, Almo SC, Rapp G, Wilson K, Petratos K, Lentfer A, Wittinghofer A, Kabsch W, Pai EF, Petsko GA, Goody RS. Time-resolved X-ray crystallographic study of the conformational change in Ha-Ras p21 protein on GTP hydrolysis. *Nature* 1990;345:309.
- Volkman BF, Lipson D, Wemmer DE, Kern D. Two-state allosteric behavior in a single-domain signaling protein. *Science* 2001;291:2429–2433.
- Tsai J, Levitt M, Baker D. Hierarchy of structure loss in MD simulations of src SH3 domain unfolding. *J Mol Biol* 1999;291:215–225.
- Tang YZ, Chen WZ, Wang CX. Molecular dynamics simulations of the gramicidin A-dimyristoylphosphatidylcholine system with an ion in the channel pore region. *Eur Biophys J* 2000;29:523–534.
- Van Belle D, De Maria L, Iurcu G, Wodak SJ. Pathways of ligand clearance in acetylcholinesterase by multiple copy sampling. *J Mol Biol* 2000;298:705–726.
- Wlodek ST, Shen T, McCammon JA. Electrostatic steering of substrate to acetylcholinesterase: analysis of field fluctuations. *Biopolymers* 2000;53:265–271.
- Daggett V, Levitt M. Realistic simulations of native-protein dynamics in solution and beyond. *Annu Rev Biophys Biomol Struct* 1993;22:353–380.
- Berneche S, Roux B. Molecular dynamics of the KcsA K(+) channel in a bilayer membrane. *Biophys J* 2000;78:2900–2917.
- Gilson MK, Straatsma TP, McCammon JA, Ripoll DR, Faerman CH, Axelsen PH, Silman I, Sussman JL. Open "back door" in a molecular dynamics simulation of acetylcholinesterase. *Science* 1994;263:1276–1278.
- Wriggers W, Schulten K. Investigating a back door mechanism of actin phosphate release by steered molecular dynamics. *Proteins* 1999;35:262–273.
- van Aalten DM, Conn DA, de Groot BL, Berendsen HJ, Findlay JB, Amadei A. Protein dynamics derived from clusters of crystal structures. *Biophys J* 1997;73:2891–2896.
- Wilson EB, Decius JC, Cross PC. *Molecular vibrations*. New York: McGraw-Hill; 1955 (New York: Dover; 1980). 388 p.
- Levy RM. Computer simulations of macromolecular dynamics: models for vibrational spectroscopy and X-ray refinement. *Ann N Y Acad Sci* 1986;482:24–43.
- Levitt M, Sander C, Stern PS. Protein normal-mode dynamics;

- trypsin inhibitor, crambin, ribonuclease, and lysozyme. *J Mol Biol* 1985;181:423–447.
30. van der Spoel D, de Groot BL, Hayward S, Berendsen HJ, Vogel HJ. Bending of the calmodulin central helix: a theoretical study. *Protein Sci* 1996;5:2044–2053.
 31. Ma J, Sigler PB, Xu Z, Karplus M. A dynamic model for the allosteric mechanism of GroEL. *J Mol Biol* 2000;302:303–313.
 32. Brooks B, Karplus M. Normal modes for specific motions of macromolecules: Application to the hinge-bending mode of lysozyme. *Proc Natl Acad Sci USA* 1985;82:4995–4999.
 33. Duncan BS, Olson AJ. Approximation and visualization of large-scale motion of protein surfaces. *J Mol Graph* 1995;13:250–257.
 34. Hinsen K, Thomas A, Field MJ. Analysis of domain motions in large proteins. *Proteins* 1999;34:369–382.
 35. Miller DW, Agard DA. Enzyme specificity under dynamic control: a normal mode analysis of alpha-lytic protease. *J Mol Biol* 1999;286:267–278.
 36. Thomas A, Hinsen K, Field MJ, Perahia D. Tertiary and quaternary conformational changes in aspartate transcarbamylase: a normal mode study. *Proteins* 1999;34:96–112.
 37. Thomas A, Field MJ, Perahia D. Analysis of the low-frequency normal modes of the R state of aspartate transcarbamylase and a comparison with the T state modes. *J Mol Biol* 1996;261:490–506.
 38. Thomas A, Field MJ, Mouawad L, Perahia D. Analysis of the low frequency normal modes of the T-state of aspartate transcarbamylase. *J Mol Biol* 1996;257:1070–1087.
 39. Hinsen K. Analysis of domain motions by approximate normal mode calculations. *Proteins* 1998;33:417–429.
 40. Marques O, Sanejouand YH. Hinge-bending motion in citrate synthase arising from normal mode calculations. *Proteins* 1995;23:557–560.
 41. Smith JC, Cusack S, Pezzeca U, Brooks B, Karplus M. Inelastic neutron scattering analysis of low frequency motion in proteins: a normal mode study of the bovine pancreatic trypsin inhibitor. *J Chem Phys* 1986;85:3636–3654.
 42. Smith J, Cusack S, Tidor B, Karplus M. Inelastic neutron scattering analysis of low-frequency motions in proteins: harmonic and damped harmonic models of bovine pancreatic trypsin inhibitor. *J Chem Phys* 1990;93:2974–2991.
 43. Noguity T, Go N. Collective variable description of small-amplitude conformational fluctuations in a globular protein. *Nature* 1982;296:776.
 44. Brooks B, Karplus M. Harmonic dynamics of proteins: Normal modes and fluctuations in bovine pancreatic trypsin inhibitor. *Proc Natl Acad Sci USA* 1983;80:6571.
 45. Go N, Noguti T, Nishikawa T. Dynamics of a small globular protein, in terms of low-frequency vibrational modes. *Proc Natl Acad Sci USA* 1983;80:3696.
 46. Levy R, Srinivasan A, Olson W, McCammon J. Quasi-harmonic method for studying very low frequency modes in proteins. *Biopolymers* 1984;23:1099–1112.
 47. Levy RM, Rojas OdL, Friesner RA. Quasi-harmonic method for calculating vibrational spectra from classical simulations on multidimensional anharmonic potential surfaces. *J Phys Chem* 1984;88:4233.
 48. Henry ER. Molecular dynamics simulations of cooling in laser-excited heme proteins. *Proc Natl Acad Sci USA* 1986;83:8982–8986.
 49. Gibrat JF. Normal mode analysis of human lysozyme: study of the relative motion of the two domains and characterization of the harmonic motion. *Proteins* 1990;8:258–279.
 50. Cusack S. Temperature dependence of the low frequency dynamics of myoglobin. Measurement of the vibrational frequency distribution by inelastic neutron scattering. *Biophys J* 1990;58:243–251.
 51. Ahn JS, Kanematsu Y, Enomoto M, Kushida T. Determination of weighted density states of vibrational modes in Zn-substituted myoglobin. *Chem Phys Lett* 1993;215:336–340.
 52. Alden R, Schneebeck M, Ondrias M, Courtney S, Friedman J. Mode-specific relaxation dynamics of photoexcited Fe(II) protoporphyrin IX in hemoglobin. *J Raman Spectrosc* 1992;23:569–574.
 53. Miller RJD. Energetics and dynamics of deterministic protein motion. *Acc Chem Res* 1994;27:145–150.
 54. Durand P, Trinquier G, Sanejouand Y-H. A new approach for determining low-frequency normal modes in macromolecules. *Biopolymers* 1994;34:759–771.
 55. Perahia D, Mouawad L. Computation of low-frequency normal modes in macromolecules: improvements to the method of diagonalization in a mixed basis and application to hemoglobin. *Comput Chem* 1995;19:241–246.
 56. Amadei A, Linssen AB, Berendsen HJ. Essential dynamics of proteins. *Proteins* 1993;17:412–425.
 57. de Groot BL, Vriend G, Berendsen HJ. Conformational changes in the chaperonin GroEL: new insights into the allosteric mechanism. *J Mol Biol* 1999;286:1241–1249.
 58. Hayward S, Kitao A, Berendsen HJ. Model-free methods of analyzing domain motions in proteins from simulation: a comparison of normal mode analysis and molecular dynamics simulation of lysozyme. *Proteins* 1997;27:425–437.
 59. Wilson CA, Kreychman J, Gerstein M. Assessing annotation transfer for genomics: quantifying the relations between protein sequence, structure and function through traditional and probabilistic scores. *J Mol Biol* 2000;297:233–249.
 60. Hubbard TJP, Murzin AG, Brenner SE, Chothia C. SCOP: a structural classification of proteins database. *Nucleic Acids Res* 1997;25:236–239.
 61. Brenner S, Chothia C, Hubbard TJP, Murzin AG. Understanding protein structure: using SCOP for fold interpretation. *Methods Enzymol* 1996;266:635–642.
 62. Dubchak I, Muchnik I, Kim SH. Protein folding class predictor for SCOP: approach based on global descriptors. *Proc ISMB97*. 1997; 5:104–107.
 63. Gerstein M, Levitt M. Comprehensive assessment of automatic structural alignment against a manual standard, the Scop classification of proteins. *Protein Sci* 1998;7:445–456.
 64. Gerstein M, Levitt M. A structural census of the current population of protein sequences. *Proc Natl Acad Sci USA* 1997;94:11911–11916.
 65. Murzin A, Brenner SE, Hubbard T, Chothia C. SCOP: A structural classification of proteins for the investigation of sequences and structures. *J Mol Biol* 1995;247:536–540.
 66. Janin J, Wodak S. Structural domains in proteins and their role in the dynamics of protein function. *Prog Biophys Mol Biol* 1983;42: 21–78.
 67. Gerstein M, Lesk AM, Chothia C. Structural mechanisms for domain movements. *Biochemistry* 1994;33:6739–6749.
 68. Brünger AT. X-PLOR 3.1, A System for X-ray crystallography and NMR. New Haven: Yale University Press; 1993. 382 p.
 69. Hogue CW, Ohkawa H, Bryant SH. A dynamic look at structures: WWW-Entrez and the molecular modeling database. *Trends Biochem Sci* 1996;21:226–229.
 70. Ohkawa H, Ostell J, Bryant S. MMDB: an ASN.1 specification for macromolecular structure. *Ismb* 1995;3:259–267.
 71. Hinsen K. The molecular modeling toolkit: a new approach to molecular simulations. *J Comp Chem* 2000;79–85.
 72. Ascher D, Dubois PF, Hinsen K, Hugunin J, Oliphant T. Numerical Python. Livermore, CA: Lawrence Livermore National Laboratory; 2000.
 73. Wall L, Christiansen D, Schwartz R. Programming Perl. Sebastopol, CA: O'Reilly and Associates; 1996. 645 p.
 74. Press WH, Flannery BP, Teukolsky SA, Vetterling WT. Numerical recipes in C. Cambridge, MA: Cambridge University Press.
 75. Levy R, Perahia D, Karplus M. Molecular dynamics of an ff-helical polypeptide: temperature dependence and deviation from harmonic behavior. *Proc Natl Acad Sci USA* 1982;79:1346–1350.
 76. Christendat D, Yee A, Dharamsi A, Kluger Y, Savchenko A, Cort JR, Booth V, Mackereth CD, Saridakis V, Elkiel I, Kozlov G, Maxwell KL, Wu N, McIntosh LP, Gehring K, Kennedy MA, Davidson AR, Pai EF, Gerstein M, Edwards AM, Arrowsmith CH. Structural proteomics of an archaeon. *Nat Struct Biol* 2000;7:903–909.
 77. Ripley BD. Pattern recognition and neural networks. Cambridge, MA: Cambridge University Press; 1996.
 78. Venables WN, Ripley BD. Modern applied statistics with S-PLUS. New York: Springer; 1997.
 79. Krause A, Olson M. The basics of S and S-Plus. New York: Springer; 2000.
 80. Qian J, Stenger B, Wilson CA, Lin J, Jansen R, Teichmann SA, Park J, Krebs W, Yu H, Alexandrov V, Echols N, Gerstein M. PartsList: a web-based system for dynamically ranking protein folds based on disparate attributes, including whole-genome expression and interaction information. *Nucleic Acids Res* 2001;29: 1750–1764.