

On Sufficient Statistics of Least-Squares Superposition of Vector Sets

ARUN S. KONAGURTHU¹, PARTHAN KASARAPU¹, LLOYD ALLISON¹,
JAMES H. COLLIER¹ and ARTHUR M. LESK^{2,3}

ABSTRACT

The problem of superposition of two corresponding vector sets by minimizing their sum-of-squares error under orthogonal transformation is a fundamental task in many areas of science, notably structural molecular biology. This problem can be solved exactly using an algorithm whose time complexity grows linearly with the number of correspondences. This efficient solution has facilitated the widespread use of the superposition task, particularly in studies involving macromolecular structures. This article formally derives a set of *sufficient statistics* for the least-squares superposition problem. These statistics are additive. This permits a highly efficient (*constant time*) computation of superpositions (and sufficient statistics) of vector sets that are composed from its constituent vector sets under addition or deletion operation, where the sufficient statistics of the constituent sets are already known (that is, the constituent vector sets have been previously superposed). This results in a drastic improvement in the run time of the methods that commonly superpose vector sets under addition or deletion operations, where previously these operations were carried out *ab initio* (ignoring the sufficient statistics). We experimentally demonstrate the improvement our work offers in the context of protein structural alignment programs that assemble a reliable structural alignment from well-fitting (substructural) fragment pairs. A C++ library for this task is available online under an open-source license.

Key words: RMSD, superposition, alignment, sufficient statistics.

1. INTRODUCTION

OPTIMAL SUPERPOSITION OF TWO CORRESPONDING VECTOR SETS is a commonly used method to measure spatial similarity of three-dimensional (3D) objects. In this method, treating both the vector sets as rigid bodies, one set is *rotated* and *translated* to fit on another. Such superposition of vector sets permits the evaluation of shape similarity both qualitatively and visually. A nearly universal *criterion of optimality* for the superposition problem is the one that *minimizes the sum of squared deviations* between the corresponding

¹Clayton School of Computer Science and Information Technology, Faculty of Information Technology, Monash University, Clayton, Australia.

²The Huck Institute of Genomics, Proteomics and Bioinformatics and ³Department of Biochemistry and Molecular Biology, Pennsylvania State University, University Park, Pennsylvania.

vectors after superposition (or least-squares superposition, in short). This yields a reliable quantitative measure of similarity, the *root mean square deviation* (or r.m.s.d.) between vector sets (Eidhammer et al., 2004; Lesk, 2000, 2001).

Many approaches to address the least-squares superposition problem in three dimensions (3D) have been proposed over the years (Cohen, 1997; Coutsias et al., 2004; Cox, 1967; Diamond, 1988; Kabsch, 1976, 1978; Kearsley, 1989; KenKnight, 1984; Koehl, 2001; Lesk, 1986; Mackay, 1984; McLachlan, 1972, 1982). It is indeed remarkable that this problem can be solved exactly and efficiently, and involves a computational effort that grows *linearly* with the number of correspondences in the vector sets. Noteworthy of the approaches to solve the superposition problem is the method by Kabsch (1976) that allows computing the optimal transformation via singular value decomposition of a covariance matrix derived from the coordinates of the corresponding vector sets. An equivalent, and more analytically elegant, approach for this problem proposed by Kearsley (1989) uses the algebra of *quaternions* (Hamilton and Hamilton, 1866). Quaternions are generalizations of complex numbers with direct applications to orthogonal transformations in three-dimensional (3D) space. Specifically, the space group corresponding to unit quaternions is equivalent to the group of all possible proper 3D rotations defined about an arbitrary origin. Any pure rotation in 3D by an angle θ about some normalized axis \hat{n} passing through the origin can be represented using a unit quaternion denoted by $[\cos(\frac{\theta}{2}), \hat{n} \sin(\frac{\theta}{2})]$. Among the key advantages of using Kearsley's method to solve the least-squares superposition problem are

1. the problem can be solved analytically as an eigenvalue and eigenvector problem in quaternion parameters, and
2. the method avoids problems with singularities (and rotary inversions) that can result from Kabsch's method, where these oddities are handled explicitly after the solution is found (Coutsias et al., 2004; Kearsley, 1989).

Structural molecular biology employs least-squares superposition to support a wide variety of tasks. An example is the role of superposition in programs for aligning protein 3D structures (Kolodny et al., 2005; Konagurthu et al., 2006; Lackner et al., 2000; Lesk, 1986; Shatsky et al., 2002, 2004; Shindyalov and Bourne, 1998; Vriend and Sander, 1991; Ye and Godzik, 2003). A structural alignment is the assignment of amino acid residue–residue correspondences between proteins based on the similarity of their structural contexts (Konagurthu et al., 2006). Among the commonly used heuristics to align protein structures are the ones that rely on identifying a library of well-fitting fragments (or contiguous substructures) within the protein structures being aligned. This library is then refined by *jointly superposing* the fragment pairs and determining the pairs that fit consistently, from which a structural alignment is finally *assembled*.

The joint superpositions in the current structural aligners are computed from scratch, even though previous superpositions of constituent fragment pairs provide a lot of information about the joint superposition. We can see that the *number* of joint superpositions grows quadratically in the size of the well-fitting fragment-pairs library, where each joint superposition takes a computational effort that is linear in the size of the combined vector sets. On average, when aligning a pair of homologous protein structures, there are well in excess of a thousand well-fitting fragment pairs, exhaustively requiring millions of joint superpositions before a structural alignment can be assembled. This exhaustive step poses a significant computational bottleneck for structural aligners, and, therefore, this problem is commonly mitigated by invoking joint superpositions rather restrictively, trading off potential structural alignment quality for speed.

In this article we explore the foundations of the orthogonal superposition problem and derive a set of statistics that are sufficient to compute the r.m.s.d. of the best superposition, and its corresponding rotation and translation parameters. We demonstrate that these *sufficient statistics* (Hogg and Craig, 1994) are additive. Therefore, these statistics can be used to compute new superpositions in *constant time* using the sufficient statistics of superpositions of constituent vector sets under addition and deletion operations. Using sufficient statistics results in a drastic speed up of tasks like joint superposition described above, compared to the current approaches that recompute these superpositions from scratch.

Section 2 gives the background to the orthogonal superposition problem using the least-squares criterion. Section 3 introduces the notion of sufficient statistics and derives the full set of sufficient statistics for the orthogonal superposition problem. Section 4 provides the update rules to update the superposition parameters in constant-time, building on the sufficient statistics of constituent superpositions. Section 5 describes an approach to speed up the diagonalization step used in the Kearsley approach. Finally, the article ends with an experimental evaluation of performance gain achieved using the results of this work.

2. ORTHOGONAL SUPERPOSITION

Formally let $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ and $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ denote two vector sets with one-to-one correspondence. In this article we consider vectors in three dimensions. Let the (x, y, z) components of each \mathbf{u}_i be represented here as $(\mathbf{u}_i(x), \mathbf{u}_i(y), \mathbf{u}_i(z))$. (Similar representation holds for \mathbf{v}_i or any other vector.)

The least-squares superposition problem is an optimization problem that involves finding the orthogonal rotation \mathbf{R} and translation \mathbf{t} with the optimality criterion defined as follows:

$$\mathcal{E} = \min \|\mathbf{R}\mathcal{U} + \mathbf{t} - \mathcal{V}\|^2 = \min \sum_{i=1}^n \|\mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i\|^2 = \min \sum_{i=1}^n \langle \mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i, \mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i \rangle$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product between the stated terms, \mathbf{R} is a 3×3 pure rotation matrix, and \mathbf{t} is a translation vector.

Under this least-squares criterion, the optimal translation can be made independent of the optimal rotation as follows. Differentiating \mathcal{E} with respect to \mathbf{t} and evaluating it at its extremum

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{t}} &= \frac{\partial}{\partial \mathbf{t}} \sum_{i=1}^n \langle \mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i, \mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i \rangle = 2 \sum_{i=1}^n \frac{\partial (\mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i)}{\partial \mathbf{t}} (\mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i) = 0 \\ &\Rightarrow \sum_{i=1}^n \mathbf{R}\mathbf{u}_i + \mathbf{t} - \mathbf{v}_i = 0 \Rightarrow \mathbf{t} = \frac{\sum_{i=1}^n \mathbf{v}_i}{n} - \mathbf{R} \frac{\sum_{i=1}^n \mathbf{u}_i}{n} = \mathbf{Centroid}(\mathcal{V}) - \mathbf{R} \mathbf{Centroid}(\mathcal{U}) \end{aligned}$$

It follows that moving each of the vector sets to an origin at its centroid, about which the rotation is defined, gives us a modified (but equivalent) objective, which is independent of the translation \mathbf{t} : $\mathcal{E} = \min \sum_{i=1}^n |\mathbf{R}\mathbf{u}'_i - \mathbf{v}'_i|^2$, where $\mathbf{u}'_i = \mathbf{u}_i - \mathbf{Centroid}(\mathcal{U})$ and $\mathbf{v}'_i = \mathbf{v}_i - \mathbf{Centroid}(\mathcal{V})$.

Kearsley (1989) proposed an elegant method that removes the non linear aspect to this least-squares problem and transforms it into an eigenvalue problem of the form $\mathbf{Q}\mathbf{q} = \lambda\mathbf{q}$, or in expanded terms

$$\begin{pmatrix} \sum (x_m^2 + y_m^2 + z_m^2) & \sum (y_p z_m - y_m z_p) & \sum (x_m z_p - x_p z_m) & \sum (x_p y_m - x_m y_p) \\ \sum (y_p z_m - y_m z_p) & \sum (x_m^2 + y_p^2 + z_p^2) & \sum (x_m y_m - x_p y_p) & \sum (x_m z_m - x_p z_p) \\ \sum (x_m z_p - x_p z_m) & \sum (x_m y_m - x_p y_p) & \sum (x_p^2 + y_m^2 + z_p^2) & \sum (y_m z_m - y_p z_p) \\ \sum (x_p y_m - x_m y_p) & \sum (x_m z_m - x_p z_p) & \sum (y_m z_m - y_p z_p) & \sum (x_p^2 + y_p^2 + z_m^2) \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{pmatrix} = \lambda \begin{pmatrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{pmatrix} \quad (1)$$

As seen above, \mathbf{Q} is a 4×4 square symmetric matrix and $\mathbf{q} = (q_1, q_2, q_3, q_4)^T$ are the unknown (to be solved) quaternion components associated with a 3D rotation, and λ is an (unknown) eigenvalue. In the eigenvalue problem defined in Equation 1, the notation x_m , a scalar quantity, denotes the component-wise difference $\mathbf{v}'_i(x) - \mathbf{u}'_i(x)$ (equivalent notations for y_m and z_m), and the scalar x_p denotes the component-wise sum $\mathbf{v}'_i(x) + \mathbf{u}'_i(x)$ (equivalently, y_p and z_p). From this point onwards, we use the term *quaternion matrix* to refer to the 4×4 matrix \mathbf{Q} shown in Equation 1.

Diagonalizing \mathbf{Q} yields four eigenvalues and (corresponding) eigenvectors. Kearsley (1989) shows that the eigenvector corresponding to the smallest eigenvalue, λ_{\min} , corresponds to the best rotation producing the least squares error, and the r.m.s.d. is computed as $\sqrt{\frac{\lambda_{\min}}{n}}$

The computational effort required to solve the rigid-body superposition problem using Kearsley's quaternion approach (or equivalently Kabsch's approach) grows linearly with the number of vectors being superimposed. In Kearsley's approach this is dominated by the computation of the quaternion matrix \mathbf{Q} where each of 10 distinct terms in the matrix requires $O(n)$ effort. The diagonalization of \mathbf{Q} is independent of n and shows a rapid convergence with numerical methods such as Jacobi's diagonalization algorithm (Jacobi, 1846).

3. SUFFICIENT STATISTICS

We note that this rigid-body superposition problem is a geometric instance of the general regression problem using total least-squares, where a regression line has to be determined that minimizes the sum-of-squares error with respect to the observed data points.

The error terms of this regression problem are assumed to be normally distributed as $\mathcal{N}(0, \sigma)$, where the mean μ is 0 and σ is the standard deviation, which is minimized by the problem. In fact, the least squares estimator of σ is also its maximum likelihood estimator.

More formally, consider the standard normal distribution of some random variable x :

$$\mathcal{N}(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

This normal density can be reparameterized into a general form denoting the family of exponential distributions: $f(x|\boldsymbol{\eta}) = h(x)g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{U}(x))$ where $h(x) = \frac{1}{\sqrt{\pi}}$, $g(\boldsymbol{\eta}_2) = \sqrt{-\boldsymbol{\eta}_2} \exp\left(\frac{\boldsymbol{\eta}_1^2}{4\boldsymbol{\eta}_2}\right)$, $\boldsymbol{\eta}^T = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)$, and $\mathbf{U}^T(x) = (x, x^2)$.

This transformation can be used to show certain important properties that allows efficient computation of maximum likelihood estimators of μ and σ .

Considering a sample set of observations that are normally distributed $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$. The likelihood for these samples is given by:

$$f(\mathbf{X}|\boldsymbol{\eta}) = \left(\prod_{i=1}^n h(x_i)\right) (g(\boldsymbol{\eta}))^n \exp\left(\boldsymbol{\eta}^T \sum_{i=1}^n \mathbf{U}(x_i)\right)$$

Taking natural logarithms on both sides gives us the log likelihood:

$$\log(f(\mathbf{X}|\boldsymbol{\eta})) = \kappa + n \log(g(\boldsymbol{\eta})) + \boldsymbol{\eta}^T \sum_{i=1}^n \mathbf{U}(x_i)$$

where $\kappa = \sum_{i=1}^n \log(h(x_i))$ is a term independent of $\boldsymbol{\eta}$.

To find the maximum likelihood estimators $\hat{\boldsymbol{\eta}}$, take the gradient with respect to $\boldsymbol{\eta}$ and set to 0. This results in:

$$n \nabla_{\boldsymbol{\eta}}[\log(g(\hat{\boldsymbol{\eta}}))] + \sum_{i=1}^n \mathbf{U}(x_i) = 0 \Rightarrow -\nabla_{\boldsymbol{\eta}}[\log(g(\hat{\boldsymbol{\eta}}))] = \frac{1}{n} \sum_{i=1}^n \mathbf{U}(x_i) = \frac{-1}{g(\hat{\boldsymbol{\eta}})} \nabla_{\boldsymbol{\eta}} g(\hat{\boldsymbol{\eta}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{U}(x_i)$$

Notice that maximum likelihood estimate $\hat{\boldsymbol{\eta}}$ depends on the statistic $\sum_{i=1}^n \mathbf{U}(x_i)$ rather than the individual data.

This suggests that to obtain the maximum likelihood estimate we do not need the data explicitly as it can be derived from that statistic. This sufficiency to derive the maximum likelihood estimator without explicit consideration of data makes $\sum_{i=1}^n \mathbf{U}(x_i)$ a *sufficient statistic* for the exponential family of functions. For normal distribution, we saw earlier that $\mathbf{U}(x_i) = (x_i, x_i^2)$ gives the sufficient statistics of $\sum_{i=1}^n x_i$ and $\sum_{i=1}^n x_i^2$ (Hogg and Craig, 1994).

3.1. Sufficient statistics for orthogonal superposition

For the orthogonal superposition problem, each error term, $\varepsilon_i = \mathbf{R}\mathbf{u}_i' - \mathbf{v}_i'$, is also assumed to be normally distributed, that is, $\varepsilon_i \sim \mathcal{N}(\mu=0, \sigma)$. We now derive the sufficient statistics for σ of ε_i terms, which is equivalent to the r.m.s.d. after least-squares superposition. The likelihood of the observed normally distributed errors after superposition, $\mathbf{E} = \{\varepsilon_1, \dots, \varepsilon_n\}$, can be written as:

$$f(\varepsilon_1, \dots, \varepsilon_n|\sigma) = \prod_{i=1}^n (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{R}\mathbf{u}_i' - \mathbf{v}_i'\|^2\right) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \|\mathbf{R}\mathbf{u}_i' - \mathbf{v}_i'\|^2\right)$$

Let's examine the decomposition of any ε_i term:

$$\varepsilon_i^2 = \|\mathbf{R}\mathbf{u}_i' - \mathbf{v}_i'\|^2 = \|\mathbf{u}_i'\|^2 + \|\mathbf{v}_i'\|^2 - 2\mathbf{v}_i'^T \mathbf{R}\mathbf{u}_i'. \quad (2)$$

From Equation 1, the matrix \mathbf{Q} is made up of terms of the form $A_m = \mathbf{v}_i'(A) - \mathbf{u}_i'(A)$ and $A_p = \mathbf{v}_i'(A) + \mathbf{u}_i'(A)$ where each A and B take the values x , y , or z denoting the vector components. Rewriting, we get $\mathbf{v}_i'(A) = \frac{A_p + A_m}{2}$ and $\mathbf{u}_i'(A) = \frac{A_p - A_m}{2}$. The first two terms on the right-hand side of Equation 2 can be expanded as follows:

$$\begin{aligned}\|\mathbf{u}_i'\|^2 + \|\mathbf{v}_i'\|^2 &= (u_i'(x)^2 + u_i'(y)^2 + u_i'(z)^2) + (v_i'(x)^2 + v_i'(y)^2 + v_i'(z)^2) \\ &= \frac{1}{2}(x_m^2 + x_p^2 + y_m^2 + y_p^2 + z_m^2 + z_p^2) = \frac{1}{2} \sum_{A \in \{x, y, z\}} A_m^2 + \frac{1}{2} \sum_{A \in \{x, y, z\}} A_p^2\end{aligned}\quad (3)$$

The last term on the right-hand side of Equation 2 can be expanded as $\mathbf{v}_i'^T \mathbf{R} \mathbf{u}_i' = \mathbf{v}_i'^T [\mathbf{r}_1 \mathbf{r}_2 \mathbf{r}_3] \mathbf{u}_i'$ where \mathbf{r}_1 , \mathbf{r}_2 , \mathbf{r}_3 are column vectors of the 3×3 rotation matrix \mathbf{R} . Therefore,

$$\mathbf{v}_i'^T \mathbf{R} \mathbf{u}_i' = (\mathbf{v}_i' \cdot \mathbf{r}_1) u_i'(x) + (\mathbf{v}_i' \cdot \mathbf{r}_2) u_i'(y) + (\mathbf{v}_i' \cdot \mathbf{r}_3) u_i'(z) \quad (4)$$

Take the first term on the right-hand side of Equation 4. This can be expanded as:

$$\begin{aligned}(\mathbf{v}_i' \cdot \mathbf{r}_1) u_i'(x) &= r_{11} v_i'(x) u_i'(x) + r_{12} v_i'(y) u_i'(x) + r_{13} v_i'(z) u_i'(x) \\ &= \frac{r_{11}}{4} (x_p + x_m)(x_p - x_m) + \frac{r_{12}}{4} (y_p + y_m)(x_p - x_m) + \frac{r_{13}}{4} (z_p + z_m)(x_p - x_m) \\ &= \frac{r_{11}}{4} (x_p^2 - x_m^2) + \frac{r_{12}}{4} (y_p x_p - y_p x_m + y_m x_p - y_m x_m) \\ &\quad + \frac{r_{13}}{4} (z_p x_p - z_p x_m + z_m x_p - z_m x_m)\end{aligned}$$

where r_{11} , r_{12} , r_{13} are the terms in the \mathbf{r}_1 column vector in \mathbf{R} . More generally,

$$(\mathbf{v}_i' \cdot \mathbf{r}_1) u_i'(x) = c_1 A_p^2 + c_2 A_m^2 + c_3 A_p B_p + c_4 A_m B_m + c_5 A_m B_p \quad (5)$$

where c_k are constants in terms of components of \mathbf{r}_1 .

Similarly, $(\mathbf{v}_i' \cdot \mathbf{r}_2) u_i'(y)$ and $(\mathbf{v}_i' \cdot \mathbf{r}_3) u_i'(z)$ can be expanded as above and will have the same form as Equation 5, but with different constants. Therefore, combining Equations 3 and 4, Equation 2 can be expressed as

$$\varepsilon_i^2 = \zeta_1 \sum_A A_p^2 + \zeta_2 \sum_A A_m^2 + \zeta_3 \sum_{\forall A \neq B} A_p B_p + \zeta_4 \sum_{\forall A \neq B} A_m B_m + \zeta_5 \sum_{\forall A \neq B} A_m B_p$$

where ζ_k are constants. Hence, the likelihood function can be written as

$$f(\varepsilon_1, \dots, \varepsilon_n | \sigma) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \mathbf{U}\right) \quad (6)$$

where $\mathbf{U} = \sum_{i=1}^n \left(\zeta_1 \sum_A A_p^2 + \zeta_2 \sum_A A_m^2 + \zeta_3 \sum_{\forall A \neq B} A_p B_p + \zeta_4 \sum_{\forall A \neq B} A_m B_m + \zeta_5 \sum_{\forall A \neq B} A_m B_p \right)$ and $A, B \in \{x, y, z\}$.

Using Equation 6, the negative log-likelihood is given as:

$$\mathcal{L}(\varepsilon_1, \dots, \varepsilon_n | \sigma) = \frac{n}{2} \log(2\pi) + n \log \sigma + \frac{1}{2\sigma^2} \mathbf{U} \quad (7)$$

The maximum likelihood estimate $\hat{\sigma}$ can be determined by minimizing Equation 7 and evaluating the corresponding σ , that is, $\frac{\partial \mathcal{L}}{\partial \sigma} = 0 \Rightarrow \hat{\sigma}^2 = \frac{\mathbf{U}}{n}$.

\mathbf{U} involves statistics that do not take into account the data (the coordinates of vector sets) explicitly, and are sufficient to estimate σ (or r.m.s.d.). Therefore, the set of *sufficient statistics* for the least-squares superposition problem can be defined as:

$$\Psi = \left\{ \sum_{i=1}^n A_m B_m, \quad \sum_{i=1}^n A_m B_p, \quad \sum_{i=1}^n A_p B_p \right\} \quad (8)$$

where A and B take the values $\{x, y, z\}$, $A_m = \mathbf{v}_i'(A) - \mathbf{u}_i'(A)$ is the component-wise difference (similarly B_m), and $A_p = \mathbf{v}_i'(A) + \mathbf{u}_i'(A)$ is the component-wise sum (similarly B_p). Altogether, the set Ψ consists of 18 distinct statistics. Since we translated the original vector sets, these terms are sufficient to compute the best rotation for the least-squares superposition problem. Further, the sufficient statistics to compute the translation are simply the component-wise sums for each of the two (untranslated) vector sets \mathcal{U} and \mathcal{V} .

4. UPDATING SUFFICIENT STATISTICS

4.1. Updating sufficient statistics under an addition operation

Consider two pairs of corresponding vector sets: $\mathcal{Q} \leftrightarrow \mathcal{R}$ containing n_1 correspondences and $\mathcal{S} \leftrightarrow \mathcal{T}$ containing n_2 correspondences. Let \mathcal{U} be defined as a combination of vectors \mathcal{Q} and \mathcal{S} , and similarly \mathcal{V} as a combination of \mathcal{R} and \mathcal{T} . Let Ψ_1 denote the sufficient statistics of the first pair after superposition and Ψ_2 denote the same for the second pair. Define these as:

$$\Psi_1 = \left\{ \sum_{i=1}^{n_1} C_m D_m, \sum_{i=1}^{n_1} C_m D_p, \sum_{i=1}^{n_1} C_p D_p \right\}, \Psi_2 = \left\{ \sum_{i=1}^{n_2} E_m F_m, \sum_{i=1}^{n_2} E_m F_p, \sum_{i=1}^{n_2} E_p F_p \right\} \quad (9)$$

Where $C, D, E,$ and F take the values $x, y,$ or $z,$ denoting the respective components of vectors in the sets. Consistent with the previous notation (see Eq. 8), C_p and C_m (similarly D_p and D_m) are the component-wise sums and differences between corresponding vectors in \mathcal{Q} and \mathcal{R} . The same definitions hold for E_m (and E_p) and F_m (and F_p), with respect to corresponding vectors in \mathcal{S} and \mathcal{T} .

We want to use Ψ_1 and Ψ_2 to compute a new set of sufficient statistics Ψ (defined in Eq. 8) for the superposition of vector sets $\mathcal{U} = \mathcal{Q} + \mathcal{S}$ with $\mathcal{V} = \mathcal{R} + \mathcal{T}$. Below we derive the construction of the new sufficient statistics. The statistics involved in computing the new centroids of the sets \mathcal{U} and \mathcal{V} , $\sum_{i=1}^{n=n_1+n_2} \mathbf{u}_i(A)$ and $\sum_{i=1}^{n=n_1+n_2} \mathbf{v}_i(A)$, can be trivially updated using the statistics $\sum_{i=1}^{n_1} \mathbf{q}_i(C), \sum_{i=1}^{n_1} \mathbf{r}_i(D), \sum_{i=1}^{n_2} \mathbf{s}_i(E),$ and $\sum_{i=1}^{n_2} \mathbf{t}_i(F)$. To compute the remaining statistics in Ψ , we define the following vectors:

$$\begin{aligned} \alpha_1 &= \mathbf{Centroid}(\mathcal{U}) - \mathbf{Centroid}(\mathcal{Q}) & \beta_1 &= \mathbf{Centroid}(\mathcal{V}) - \mathbf{Centroid}(\mathcal{R}) \\ \alpha_2 &= \mathbf{Centroid}(\mathcal{U}) - \mathbf{Centroid}(\mathcal{S}) & \beta_2 &= \mathbf{Centroid}(\mathcal{V}) - \mathbf{Centroid}(\mathcal{T}) \end{aligned}$$

Also define $\Delta_m^C = \beta_1(C) - \alpha_1(C), \Delta_m^D = \beta_1(D) - \alpha_1(D), \Delta_m^E = \beta_2(E) - \alpha_2(E),$ and $\Delta_m^F = \beta_2(F) - \alpha_2(F)$, where $A = C = E \in \{x, y, z\}$ and $B = D = F \in \{x, y, z\}$. These vectors define the corrections (or updates) required to the constituent centroids so that the new centroid can be constructed. Using these definitions, the lemma and corollaries below show the computation of Ψ from Ψ_1 and Ψ_2 .

Lemma 1.
$$\sum_{i=1}^{n=n_1+n_2} A_m B_m = \sum_{i=1}^{n_1} C_m D_m + n_1 \Delta_m^C \Delta_m^D + \sum_{i=1}^{n_2} E_m F_m + n_2 \Delta_m^E \Delta_m^F$$

Proof.

$$\begin{aligned} \text{Updated} \left(\sum_{i=1}^{n_1} C_m D_m \right) &= \sum_{i=1}^{n_1} [(\mathbf{r}'(C) + \beta_1(C)) - (\mathbf{q}'(C) + \alpha_1(C))] [(\mathbf{r}'(D) + \beta_1(D)) - (\mathbf{q}'(D) + \alpha_1(D))] \\ &= \sum_{i=1}^{n_1} [(\mathbf{r}'(C)\mathbf{r}'(D) - \mathbf{q}'(C)\mathbf{q}'(D) - \mathbf{r}'(C)\mathbf{q}'(D) + \mathbf{q}'(C)\mathbf{r}'(D))] \\ &+ \sum_{i=1}^{n_1} [(\beta_1(C)\mathbf{r}'(D) - \alpha_1(C)\mathbf{r}'(D) - \beta_1(C)\mathbf{q}'(D) + \alpha_1(C)\mathbf{q}'(D))] \\ &+ \sum_{i=1}^{n_1} [(\mathbf{r}'(C)\beta_1(D) - \mathbf{r}'(C)\alpha_1(D) - \mathbf{q}'(C)\beta_1(D) + \mathbf{q}'(C)\alpha_1(D))] \\ &+ \sum_{i=1}^{n_1} [(\beta_1(C)\beta_1(D) - \alpha_1(C)\beta_1(D) - \beta_1(C)\alpha_1(D) + \alpha_1(C)\alpha_1(D))] \\ &= \sum_{i=1}^{n_1} C_m D_m + \sum_{i=1}^{n_1} \Delta_m^C D_m + \sum_{i=1}^{n_1} \Delta_m^D C_m + \sum_{i=1}^{n_1} \Delta_m^C \Delta_m^D \\ &= \sum_{i=1}^{n_1} C_m D_m + \Delta_m^C \sum_{i=1}^{n_1} D_m + \Delta_m^D \sum_{i=1}^{n_1} C_m + n_1 \Delta_m^C \Delta_m^D \\ &= \sum_{i=1}^{n_1} C_m D_m + n_1 \Delta_m^C \Delta_m^D \quad (\because \sum_{i=1}^{n_1} C_m = \sum_{i=1}^{n_1} D_m = 0) \end{aligned}$$

Similarly, we can show that Updated $(\sum_{i=1}^{n_2} E_m F_m) = \sum_{i=1}^{n_2} E_m F_m + n_2 \Delta_m^E \Delta_m^F$. Adding the two updated statistics, the lemma above follows. ■

$$\text{Corollary 1.} \quad \sum_{i=1}^{n=n_1+n_2} A_p B_p = \left[\sum_{i=1}^{n_1} C_p D_p + n_1 \Delta_p^C \Delta_p^D \right] + \left[\sum_{i=1}^{n_2} E_p F_p + n_2 \Delta_p^E \Delta_p^F \right]$$

$$\text{Corollary 2.} \quad \sum_{i=1}^{n=n_1+n_2} A_m B_p = \left[\sum_{i=1}^{n_1} C_m D_p + n_1 \Delta_m^C \Delta_p^D \right] + \left[\sum_{i=1}^{n_2} E_m F_p + n_2 \Delta_m^E \Delta_p^F \right]$$

4.2. Updating sufficient statistics under a deletion operation

Let us consider the case where we want to find a superposition under a deletion operation. Let $\mathcal{U} \leftrightarrow \mathcal{V}$ and $\mathcal{Q} \subset \mathcal{U} \leftrightarrow \mathcal{R} \subset \mathcal{V}$ denote two pairs of vector sets that are in correspondence. We want to compute the superposition of the vector sets $\mathcal{S} = \mathcal{U} \setminus \mathcal{Q}$ and $\mathcal{T} = \mathcal{V} \setminus \mathcal{R}$.

Using the same notations as in the previous section, it is straightforward to rearrange Lemma 1 and Corollaries 1 and 2 to derive the sufficient statistics Ψ_2 of the superposition of \mathcal{S} with \mathcal{T} using the sufficient statistics Ψ (of \mathcal{U} with \mathcal{V}) and Ψ_1 (of \mathcal{Q} with \mathcal{R}) as follows.

$$\text{Corollary 3.} \quad \sum_{i=1}^{n_2} E_m F_m = \sum_{i=1}^{n=n_1+n_2} A_m B_m - \sum_{i=1}^{n_1} C_m D_m - n_1 \Delta_m^C \Delta_m^D - n_2 \Delta_m^E \Delta_m^F$$

$$\text{Corollary 4.} \quad \sum_{i=1}^{n_2} E_p F_p = \sum_{i=1}^{n=n_1+n_2} A_p B_p - \sum_{i=1}^{n_1} C_p D_p - n_1 \Delta_p^C \Delta_p^D - n_2 \Delta_p^E \Delta_p^F$$

$$\text{Corollary 5.} \quad \sum_{i=1}^{n_2} E_m F_p = \sum_{i=1}^{n=n_1+n_2} A_m B_p - \sum_{i=1}^{n_1} C_m D_p - n_1 \Delta_m^C \Delta_p^D - n_2 \Delta_m^E \Delta_p^F$$

5. COMPUTING THE R.M.S.D. FROM UPDATED SUFFICIENT STATISTICS AND OTHER OPTIMIZATIONS

It is easy to see that Kearsley's 4×4 symmetric quaternion matrix \mathbf{Q} given in Equation 1 can be constructed using the updated sufficient statistics derived from the constituent sufficient statistics of previously superposed vector sets under addition and deletion operations in constant time.

In practice, \mathbf{Q} is diagonalized and its eigenvalues and eigenvectors are computed using the numerical approach of Jacobi, where rotations are applied to \mathcal{Q} iteratively, where each rotation annihilates (that is, sets to zero) a symmetric pair of off-diagonal entries in \mathcal{Q} . This approach has a fast convergence, and the smallest diagonal element (or eigenvalue) is used to derive the r.m.s.d. term, as described in section 2.

However, in some cases, for instance, when a current superposition is extended by just one correspondence, \mathcal{Q} changes marginally. Therefore, diagonalization of \mathcal{Q} can build on the previous solution. Let \mathbf{Q} denote the quaternion matrix corresponding to the superposition of corresponding vector sets \mathcal{U} and \mathcal{V} . From the eigenvalue decomposition theorem, we have $\mathbf{Q} = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^{-1}$, where \mathbf{S} is the matrix of eigenvectors and $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues. Also note that \mathbf{Q} is a positive semidefinite matrix with the property $\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T$. This implies that all the eigenvectors are orthogonal to each other. This further simplifies the decomposition to $\mathbf{Q} = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^T$. Also, since \mathbf{S} is an orthogonal matrix, $\mathbf{Q} = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^T \Rightarrow \mathbf{\Lambda} = \mathbf{S}^T \mathbf{Q} \mathbf{S}$.

Now, assume that the corresponding vector sets are augmented from \mathcal{U} and \mathcal{V} to \mathcal{U}' and \mathcal{V}' , resulting in an updated Kearsley's matrix \mathbf{Q}' . We want to diagonalize this matrix into $\mathbf{S}' \mathbf{\Lambda}' \mathbf{S}'^T$. Instead of starting the Jacobi's iterative process from scratch, we will use the previously computed eigenvectors (before the vector sets were augmented), \mathbf{S} , and compute $\tilde{\mathbf{A}}$ as $\tilde{\mathbf{A}} = \mathbf{S}^T \mathbf{Q}' \mathbf{S}$. Notice that if the augmentation does not include drastic changes, then $\tilde{\mathbf{A}}$ is nearly diagonal (that is, $\tilde{\mathbf{A}} \approx \mathbf{\Lambda}'$), thus requiring very few iterations to fully diagonalize $\tilde{\mathbf{A}}$. This provides a further optimization to the diagonalization step under update operations on vector sets.

6. RESULTS

6.1. Source code

We wrote a C++ library that allows superposition of vector sets. Superpositions can be carried out from scratch (using the raw coordinates), or, alternatively, using the sufficient statistics of constituent vector sets

by operating on them using set addition or deletion operations. The source code is available freely under a GNU public license online (visit corresponding author's webpage).

6.2. Consistency of superpositions using sufficient statistics

To validate the consistency of superpositions generated using sufficient statistics (under both addition and deletion operations discussed in section 4) we undertake the following experiment: 8992 ASTRAL SCOP (Chandonia et al., 2004; Murzin et al., 1995) domains were used as the source structures from which superposable fragments were randomly sampled. The general procedure of sampling is as follows. From the list of source structures, randomly choose any structure. Within this structure choose two random fragments of lengths l_1 and l_2 , where each fragment has between 10 and 40 residues. These chosen fragments form the sets Q and S . Yet another structure is again chosen randomly from our source list, and two fragments are randomly extracted from it with exactly the same length, l_1 and l_2 . These form the sets R and T . Assuming one-to-one correspondence between $Q \leftrightarrow R$, we compute the sufficient statistics Ψ_1 of their orthogonal superposition. Similarly the sufficient statistics Ψ_2 is computed for the orthogonal superposition between $S \leftrightarrow T$. Define $U = Q + S$ and $V = R + T$. Iterating this random sampling 100 million times, we computed:

1. the r.m.s.d. of superposition of $U \leftrightarrow V$ from scratch (using the raw coordinates in the vector sets); denote this r.m.s.d. as ρ_1
2. the r.m.s.d. of superposition of $U \leftrightarrow V$, but using the sufficient statistics Ψ_1 and Ψ_2 of superpositions of constituent vector sets $Q \leftrightarrow R$, and $S \leftrightarrow T$; denote this r.m.s.d. as ρ_2 .

We measure the difference between the two r.m.s.d. values, $\Delta\rho = \rho_1 - \rho_2$. Over the 100 million samples, the mean and standard deviation of $\Delta\rho$ was found to be zero to a very high precision ($< 10^{-17}$).

We repeat the same experiment to validate superpositions under deletion operation using sufficient statistics, where, in each iteration, we compute the superposition of vector sets $S \leftrightarrow T$, with $S = U \setminus Q$ and $T = V \setminus R$. This experiment again confirms the same consistency as observed in the test on addition operation.

6.3. Measuring the performance gain using sufficient statistics for superpositions

We demonstrated in section 4 that superpositions under addition and deletion can be updated in constant time, building on the sufficient statistics of the constituent sets. This was also validated empirically using the experiments above. We will now measure the gain in performance using this approach by comparing it with superpositions built from scratch.

Figures 1a–c show the runtime plots of three sets of randomly chosen 10 million joint superpositions carried out from scratch (blue line) and compared against the same superpositions updated using sufficient statistics (green line). We note that these three sets vary in the (average) size of the joint superpositions being carried out, as indicated in the plot titles.

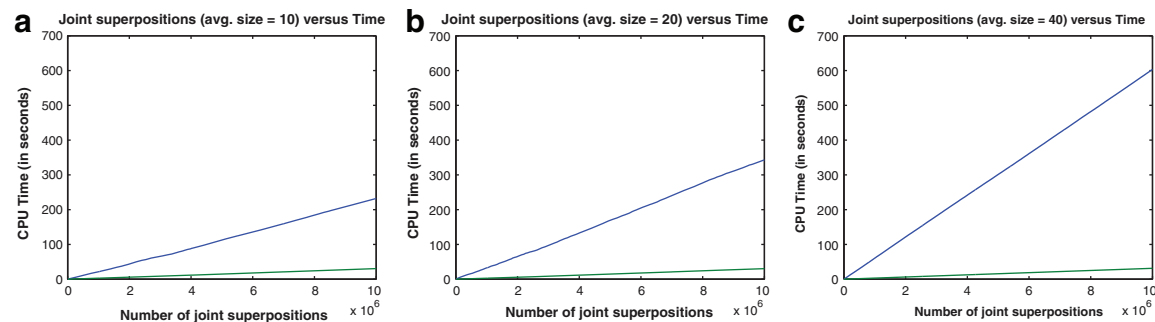


FIG. 1. The CPU times (in seconds) performing joint superpositions from scratch (blue line) compared to the same using sufficient statistics (green line) over 10 million random fragment data sets derived from ASTRAL SCOP domains. The three plots vary in the average superposition size as indicated in the title.

TABLE 1. TIME TAKEN TO PERFORM EXHAUSTIVE JOINT SUPERPOSITIONS ON A LIBRARY OF WELL-FITTING FRAGMENT PAIRS BETWEEN TWO STRUCTURES FROM DIFFERENT FAMILIES.

<i>Protein Family</i>	<i>Structural pair wwPDB IDs</i>	<i>Number of joint superpositions</i>	<i>Average size of superpositions</i>	<i>Time in seconds (from scratch)</i>	<i>Time in seconds (sufficient stats)</i>
Serine proteinases	3EST vs. 2PKA	18,486,240	14	419.6	56.0
Calmodulin-like	1NCX vs. 2SAS	67,820,481	18	1618.4	178.2
Serine proteinases	3EST vs. 2SNV	71,025,321	16	1328.0	187.8
Globins	1HHOA vs. 1HHOB	74,890,441	20	1923.3	194.6

Notice that when joint superpositions are carried out from scratch ignoring the sufficient statistics, the average size of the superpositions introduces a constant factor to the run time. This is expected as each superposition is linear in the size of the vector sets being superposed. Consequently, the slopes of those blue lines across the three plots in Figure 1 become steeper with the increase in the superposition size. On the other hand, when the joint superpositions are updated in constant time, the updates are independent of the superposition size. This is because any update involves recomputing only a small (fixed) number of sufficient statistics. This is clearly reflected in the slopes of the green lines being unchanged across the three plots in Figure 1.

More formally, if $|J|$ is the number of joint superpositions and l is the (average) number of vectors being superposed, then the first method (blue line in Fig. 1) grows as $O(l|J|)$. Since $l \ll |J|$ we see a linear trend (with a steeper gradient accounting for the multiplier l in the complexity term). In comparison, the results with sufficient statistics (green line in Fig. 1) grow simply as $O(|J|)$, independent of the superposition size.

6.4. Using sufficient statistics of superposition in the setting of structural alignments

As discussed in the introduction, a common heuristic employed to compute a structural alignment between pairs of structures involves collecting a library of well-fitting fragments. This library is refined by jointly superposing pairs from this library, and a final structural alignment is assembled from these results.

To test the potential performance gain by using sufficient statistics, we computed the time taken to undertake an exhaustive joint superposition on libraries of well-fitting fragments corresponding to a small collection of structural pairs. For each pair of structures chosen, the well-fitting fragments are identified as those that superpose *maximally* within an r.m.s.d. threshold of 2 Å. By maximal, we mean those fragment pairs that cannot be extended any further without violating the r.m.s.d. threshold.

Table 1 shows the run times of joint superpositions performed exhaustively on a small set of protein structural pairs. As can be seen from the table, using sufficient statistics for superpositions results in up to an order of magnitude improvement in the run time to carry out these superpositions exhaustively. Since performing this task *without* sufficient statistics creates a computational bottleneck, existing structural alignment programs attempt to drastically restrict the number of superpositions, often trading off structural alignment quality for speed. We note that the improvements gained from using sufficient statistics for superpositions will allow these restrictions to be generously relaxed without any effect on the current run times, but potentially improving the structural alignment quality.

All the above experiments were carried out on a standard laptop with 2.2GHz Intel[®] CPU and 4GB RAM.

7. CONCLUSIONS

Optimal superpositions of vector sets are central to identify similarities and differences between spatial objects. We derived a set of sufficient statistics for the orthogonal superposition problem under the least squares criterion. These statistics provide an efficient way to operate (via addition and deletion of vectors) on previously computed superpositions. Our results demonstrate a drastic improvement in the computational effort required to compute r.m.s.d. based on sufficient statistics. These results are relevant to many

analyses involving structural data. These include the plethora of algorithms to construct pairwise and multiple protein structural alignments by assembling fragment pairs. Source code (written in C++) to undertake superpositions implementing our work can be downloaded online.

AUTHOR DISCLOSURE STATEMENT

The authors declare that no competing financial interests exist.

REFERENCES

- Chandonia, J.-M., Hon, G., Walker, N. S., et al. 2004. The ASTRAL compendium in 2004. *Nucleic Acids Res.* 32, D189–D192.
- Cohen, G. 1997. ALIGN: a program to superimpose protein coordinates, accounting for insertions and deletions. *Journal of Applied Crystallography* 30, 1160–1161.
- Coutsias, E.A., Seok, C., and Dill, K.A. 2004. Using quaternions to calculate RMSD. *J. Comput. Chem.* 25, 1849–1857.
- Cox, J.M. 1967. Mathematical methods used in the comparison of the quaternary structures. *J. Mol. Bio.* 28, 151–156.
- Diamond, R. 1988. A note on the rotational superposition problem. *Acta Crystallographica Section A: Foundations of Crystallography* 44, 211–216.
- Eidhammer, I., Jonassen, I., and Taylor, W.R. 2004. *Protein Bioinformatics: An algorithmic Approach to Sequence and Structure Analysis*. J. Wiley & Sons, Hoboken, New Jersey.
- Hamilton, W.R., and Hamilton, W.E. 1866. *Elements of Quaternions*. Longmans, Green, & Company, New York.
- Hogg, R.V., and Craig, A. 1994. *Introduction to Mathematical Statistics*. Prentice Hall, Upper Saddle River, New Jersey.
- Jacobi, C.G.J. 1846. Über ein leichtes Verfahren, die in der Theorie der Säkularstörungen vorkommenden Gleichungen numerisch aufzulösen. *Journal für die Reine und Angewandte Mathematik* 30, 51–95.
- Kabsch, W. 1976. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 32, 922–923.
- Kabsch, W. 1978. A discussion of the solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 34, 827–828.
- Kearsley, S.K. 1989. On the orthogonal transformation used for structural comparisons. *Acta Crystallographica Section A: Foundations of Crystallography* 45, 208–210.
- KenKnight, C. 1984. Comparison of methods of matching protein structures. *Acta Crystallographica Section A: Foundations of Crystallography* 40, 708–712.
- Koehl, P. 2001. Protein structure similarities. *Current Opinion in Structural Biology* 11, 348–353.
- Kolodny, R., Koehl, P., and Levitt, M. 2005. Comprehensive evaluation of protein structure alignment methods: scoring by geometric measures. *J. Mol. bio.* 346, 1173–1188.
- Konagurthu, A.S., Whisstock, J.C., Stuckey, P.J., and Lesk, A.M. 2006. MUSTANG: a multiple structural alignment algorithm. *Proteins: Structure, Function, and Bioinformatics* 64, 559–574.
- Lackner, P., Koppensteiner, W.A., Sippl, M.J., and Domingues, F.S., 2000. Prosup: a refined tool for protein structure alignment. *Protein Engineering* 13, 745–752.
- Lesk, A. 1986. A toolkit for computational molecular biology. II. On the optimal superposition of two sets of coordinates. *Acta Crystallographica Section A: Foundations of Crystallography* 42, 110–113.
- Lesk, A.M. 2000. The unreasonable effectiveness of mathematics in molecular biology. *The Mathematical Intelligencer* 22, 28–37.
- Lesk, A.M. 2001. *Introduction to protein architecture: the structural biology of proteins*. Oxford University Press, New York.
- Mackay, A.L. 1984. Quaternion transformation of molecular orientation. *Acta Crystallographica Section A: Foundations of Crystallography* 40, 165–166.
- McLachlan, A.D. 1972. A mathematical procedure for superimposing atomic coordinates of proteins. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 28, 656–657.
- McLachlan, A.D. 1982. Rapid comparison of protein structures. *Acta Crystallographica Section A: Crystal Physics, Diffraction, Theoretical and General Crystallography* 38, 871–873.
- Murzin, A.G., Brenner, S.E., Hubbard, T., and Chothia, C. 1995. SCOP: a structural classification of proteins database for the investigation of sequences and structures. *J. Mol. Bio.* 247, 536–540.
- Shatsky, M., Nussinov, R., and Wolfson, H.J. 2002. Flexible protein alignment and hinge detection. *Proteins: Structure, Function, and Bioinformatics* 48, 242–256.

- Shatsky, M., Nussinov, R., and Wolfson, H.J. 2004. A method for simultaneous alignment of multiple protein structures. *Proteins: Structure, Function, and Bioinformatics* 56, 143–156.
- Shindyalov, I.N. and Bourne, P.E. 1998. Protein structure alignment by incremental combinatorial extension (ce) of the optimal path. *Protein Eng.* 11, 739–747.
- Vriend, G., and Sander, C. 1991. Detection of common three-dimensional substructures in proteins. *Proteins: Structure, Function, and Bioinformatics* 11, 52–58.
- Ye, Y., and Godzik, A. 2003. Flexible structure alignment by chaining aligned fragment pairs allowing twists. *Bioinformatics* 19, ii246–ii255.

Address correspondence to:

*Dr. Arun Konagurthu
Clayton School of Computer Science and
Information Technology
Faculty of Information Technology
Monash University
Clayton VIC 3800
Australia*

E-mail: arun.konagurthu@monash.edu