# THREE-DIMENSIONAL STRUCTURE DETERMINATION FROM COMMON LINES IN CRYO-EM BY EIGENVECTORS AND SEMIDEFINITE PROGRAMMING

A. SINGER* AND Y. SHKOLNISKY†

**Abstract.** The cryo-electron microscopy (EM) reconstruction problem is to find the three-dimensional structure of a macromolecule given noisy samples of its two-dimensional projection images at unknown random directions. Present algorithms for finding an initial 3D structure model are based on the "Angular Reconstitution" method in which a coordinate system is established from three projections, and the orientation of the particle giving rise to each image is deduced from common lines among the images. However, a reliable detection of common lines is difficult due to the low signal-to-noise ratio of the images. In this paper we describe two algorithms to find the unknown imaging directions of all projections by minimizing global self consistency errors. In the first algorithm, the minimizer is obtained by computing the three largest eigenvectors of a specially designed symmetric matrix derived from the common-lines, while the second algorithm is based on semidefinite programming (SDP). Compared with existing sequential algorithms, the advantages of our algorithms are four-fold: first, they correctly find all orientations at very low common-line detection rates; second, they are extremely fast, as they involve only the computation of a few top eigenvectors or a sparse SDP; third, they are non-sequential and use the information in all common-lines at once; finally, they are amenable to rigorous mathematical analysis using harmonic analysis and random matrix theory. As an example, we show a successful recovery from 500 projection images that only 10% of the common lines between them were correctly identified.

**Key words.** Cryo electron-microscopy, angular reconstitution, random matrices, semi-circle law, semidefinite programming, rotation group SO(3), tomography.

**1. Introduction.** Electron cryomicroscopy (cryo-EM) is a technique by which biological macromolecules are imaged in an electron microscope. The molecules are rapidly frozen in a thin ($\sim$ 100nm) layer of vitreous ice, trapping them in a nearly-physiological state [1, 2]. Cryo-EM images, however, have very low contrast, due to the absence of heavy-metal stains or other contrast enhancements, and have very high noise due to the small electron doses that can be applied to the specimen. Thus to obtain a reliable three-dimensional density map of a macromolecule, the information from thousands of images of identical molecules must be combined. When the molecules are arrayed in a crystal, the necessary signal-averaging of noisy images is straightforwardly performed. More challenging is the problem of single-particle reconstruction (SPR) where a three-dimensional density map is to be obtained from images of individual molecules present in random positions and orientations in the ice layer [1]. Because it does not require the formation of crystalline arrays of macromolecules, SPR is a very powerful and general technique, which has been successfully used for 3D structure determination of many protein molecules and complexes roughly 500 kDa or larger in size. In some cases, sufficient resolution ($\sim$ 0.4nm) has been obtained from SPR to allow tracing of the polypeptide chain and identification of residues in proteins [3, 4, 5]; however, even with lower resolutions many important features can be identified [6]. A particular challenge is the ab initio estimation of the 3D structure from a set of cryo-EM images. If no other information is available, the matching of common Fourier lines in three averaged images, a technique called "Angular Reconstitution" [7] allows a coordinate system to be established, and the orientation of the particles

*Department of Mathematics and PACM, Princeton University, Fine Hall, Washington Road, Princeton NJ 08544-1000 USA, email: amits@math.princeton.edu

†Department of Mathematics, Program in Applied Mathematics, Yale University, 10 Hillhouse Ave. PO Box 208283, New Haven, CT 06520-8283 USA. E-mail: yoel.shkolnisky@yale.edu

giving rise to other images are then deduced by further matching of common lines. This method fails however when the particles are too small or the signal-to-noise ratio is too low.

The common lines between three projections determine uniquely their relative orientations up to handedness (chirality). This is the basis of the angular reconstitution method of Van Heel [7], which was also developed independently by Vainshtein and Goncharov [8]. Farrow and Ottensmeyer [9] used quaternions to obtain the relative orientation of a new projection in a least square sense. The main problem with such sequential approaches is that they are sensitive to false detection of common lines which leads to the accumulation of errors. Penczek et al. [10] tried to obtain the rotations corresponding to all projections simultaneously by minimizing a global energy functional. Unfortunately, minimization of the energy functional requires a brute force search in a huge parametric space of all possible orientations for all projections. Mallick et al. [11] suggested an alternative Bayesian approach, in which the common line between a pair of projections can be inferred from their common lines with different projection triplets. The problem with this particular approach is that it requires too many (at least seven) common lines to be correctly identified simultaneously. Therefore, it is not suitable in cases where the detection rate of correct common lines is low. In [12] we introduced an improved Bayesian approach based on voting that requires only two common lines to be correctly identified simultaneously and performs at much lower detection rates.

In this paper we introduce two algorithms that find the unknown imaging directions of all projections in a globally consistent way. Both algorithms are based on relaxations of a global minimization problem of a particular self consistent error (SCE) that takes into account the matching of common lines between all pairs of images. A similar SCE was used in [9] to asses the quality of their angular reconstitution techniques. Our approach is different in the sense that we actually minimize the SCE in order to find the imaging directions. The precise definition of our global self consistency error is given in Section 2.

In Section 3 we present our first recovery algorithm in which the global minimizer is approximated by the top three eigenvectors of a specially designed symmetric matrix derived from the common-lines data. We describe how the unknown rotations are recovered from these eigenvectors.

In Section 4 we use a different relaxation of the global optimization, which leads to our second recovery method based on semidefinite programming (SDP) [13] drawing similarities to the Goemans-Williamson max-cut algorithm [14].

Compared with existing sequential algorithms, the main advantage of our methods is that they correctly find the orientations of all projections at amazingly low common-line detection rates as they take into account all the geometric information in all common-lines at once. In Section 5 we describe the results of several numerical experiments using the two algorithms, showing successful recoveries at very low common-line detection rates. For example, both algorithms successfully recover a meaningful ab-initio coordinate system from 500 projection images when only 10% of the common lines are correctly identified. The eigenvector method is extremely efficient, and the estimated 500 rotations were obtained in a manner of seconds on a standard laptop machine.

In Section 6 we show that in the limit of infinite number of projection images, the symmetric matrix that we design converges to a convolution integral operator on the rotation group SO(3). This observation explains many of the spectral properties

that the matrix exhibits. Moreover, in Section 7 we show that the effect of the misidentified common-lines is equivalent to a random matrix perturbation. Thus, using random matrix theory, we demonstrate that the top three eigenvalues are stable as long as the detection rate of common lines exceeds $\frac{6\sqrt{2}}{5\sqrt{N}}$, where $N$ is the number of images. From the practical point of view, this result implies that 3D reconstruction is possible even at extreme levels of noise, provided that enough projections were taken. Even if reconstruction is not possible from the raw noisy images, the eigenvector method would allow reconstructions from noisy class averages obtained from averaging together fewer projections.

**2. The global self consistency error.** Suppose we collect $N$ two-dimensional digitized projection images $P_1, \ldots, P_N$ of a 3D object taken at unknown random orientations. To each projection image $P_i$ $(i = 1, \ldots, N)$ there corresponds a $3 \times 3$ unknown rotation matrix $R_i$ describing its orientation (see Figure 2.1). Excluding the contribution of noise, the pixel intensities correspond to line integrals of the electric potential induced by the molecule along the path of the imaging electrons, that is,

$$P_i(x, y) = \int_{-\infty}^{\infty} \phi_i(x, y, z)\, dz, \tag{2.1}$$

where $\phi(x, y, z)$ is the electric potential of the molecule in some fixed 'laboratory' coordinate system and $\phi_i(r) = \phi(R_i^{-1} r)$ with $r = (x, y, z)$. The projection operator (2.1) is also known as the X-ray transform [15]. Our goal is to find all rotation matrices $R_1, \ldots, R_N$ given the dataset of noisy images.
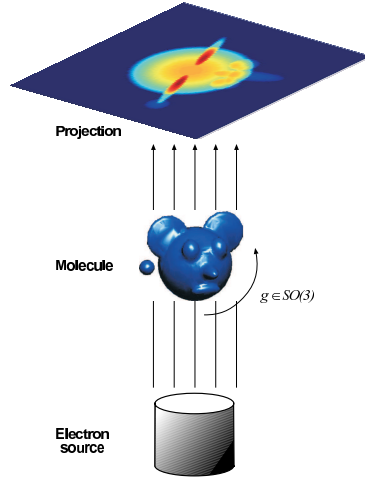


FIG. 2.1. *Schematic drawing of the imaging process: every projection image corresponds to some unknown 3D rotation of the unknown molecule.*

The Fourier-projection slice theorem (see, e.g., [15, p. 11]) says that the 2D Fourier transform of a projection image, denoted $\hat{P}$, is the restriction of the 3D Fourier transform of the projected object $\hat{\phi}$ to the central plane (i.e., going through the origin) $\theta^\perp$ perpendicular to the imaging direction, that is,

$$\hat{P}(\eta) = \hat{\phi}(\eta), \quad \eta \in \theta^\perp. \tag{2.2}$$

As every two non-parallel planes intersect at a line, it follows from the Fourier-projection slice theorem that any two projection images have a common line of inter-

section in the Fourier domain. Therefore, if $\hat{P}_i$ and $\hat{P}_j$ are the two-dimensional Fourier transforms of projections $P_i$ and $P_j$, then there must be a central line in $\hat{P}_i$ and a central line in $\hat{P}_j$ on which the two transforms agree (see Figure 2.2). This pair of lines is known as the common-line. We parameterize the common-line by $(\omega x_{ij}, \omega y_{ij})$ in $\hat{P}_i$ and by $(\omega x_{ji}, \omega y_{ji})$ in $\hat{P}_j$, where $\omega \in \mathbb{R}$ is the radial frequency and $(x_{ij}, y_{ij})$ and $(x_{ji}, y_{ji})$ are two unit vectors for which

$$\hat{P}_i(\omega x_{ij}, \omega y_{ij}) = \hat{P}_j(\omega x_{ji}, \omega y_{ji}), \text{ for all } \omega \in \mathbb{R}. \tag{2.3}$$

It is instructive to consider the unit vectors $(x_{ij}, y_{ij})$ and $(x_{ji}, y_{ji})$ as three-dimensional vectors by zero-padding: Let $c_{ij}$ and $c_{ji}$ be unit vectors in the direction of the common line between $\hat{P}_i$ and $\hat{P}_j$, respectively, given coordinate-wise by

$$c_{ij} = (x_{ij}, \ y_{ij}, \ 0)^T, \tag{2.4}$$
$$c_{ji} = (x_{ji}, \ y_{ji}, \ 0)^T. \tag{2.5}$$

Being the common-line of intersection, the mapping of $c_{ij}$ by $R_i$ must coincide with the mapping of $c_{ji}$ by $R_j$:

$$R_i c_{ij} = R_j c_{ji}, \text{ for } 1 \leq i < j \leq N. \tag{2.6}$$

These can be viewed as $\binom{N}{2}$ linear equations for the $6N$ variables corresponding to the first two columns of the rotation matrices (as $c_{ij}$ and $c_{ji}$ have a zero third entry, the third column of each rotation matrix does not contribute in (2.6)). Such overdetermined systems of linear equations are usually solved by the least squares method [10]. Unfortunately, the least squares approach is inadequate in our case due to the typically large proportion of falsely detected common lines that will dominate the sum of squares error in

$$\min_{R_1, \ldots, R_N} \sum_{i \neq j} \|R_i c_{ij} - R_j c_{ji}\|^2. \tag{2.7}$$

Moreover, the global least squares problem (2.7) is extremely difficult to solve if one requires the matrices $R_i$ to be rotations, that is, when adding the constraints

$$R_i R_i^T = I, \ \det(R_i) = 1, \text{ for } i = 1, \ldots, N, \tag{2.8}$$

where $I$ is the $3 \times 3$ identity matrix. A relaxation method that neglects the constraints (2.8) will simply collapse to the trivial solution $R_1 = \ldots = R_N = 0$ which obviously does not satisfy the constraint (2.8). Such a collapse is easily prevented by fixing one of the rotations, for example, by setting $R_1 = I$, but this would not make the robustness problem of the least squares method to go away. We therefore take a different approach for solving the global optimization problem.

Since $\|c_{ij}\| = \|c_{ji}\| = 1$ are three-dimensional unit vectors, their rotations are also unit vectors, that is, $\|R_i c_{ij}\| = \|R_j c_{ji}\| = 1$. It follows that the minimization problem (2.7) is equivalent to the maximization problem of the sum of dot products

$$\max_{R_1, \ldots, R_N} \sum_{i \neq j} R_i c_{ij} \cdot R_j c_{ji}, \tag{2.9}$$

subject to the constraints (2.8). For the true assignment of rotations, the dot product $R_i c_{ij} \cdot R_j c_{ji}$ equals 1 whenever the common line between images $i$ and $j$ was correctly
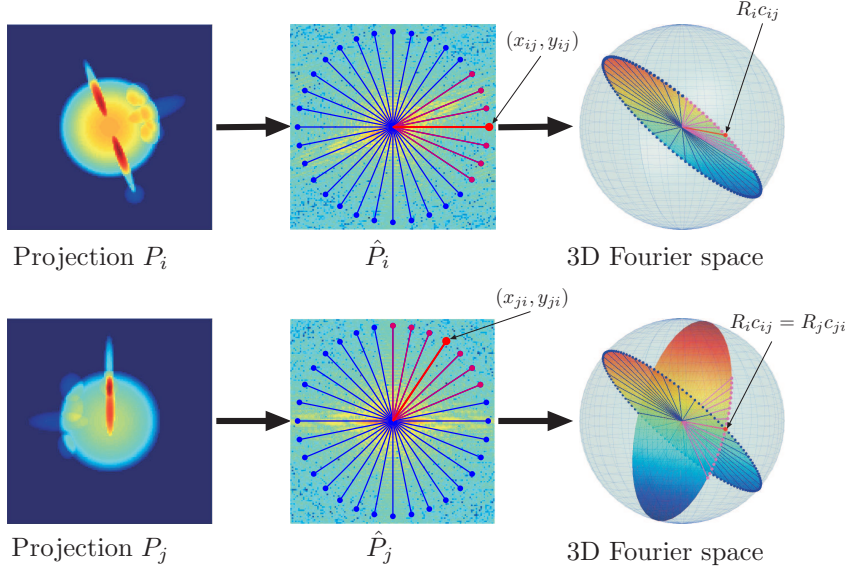
FIG. 2.2. *Fourier projection-slice theorem*

detected. Dot products corresponding to misidentified common lines can take any value between $-1$ to $1$, and if we assume that such misidentified lines have random directions, then such dot products can be considered as identically independently distributed (i.i.d) zero-mean random variables taking values in $[-1, 1]$. The objective function in (2.9) is the summation over all possible dot products. Summing up dot products that correspond to misidentified common lines results in many cancelations, whereas summing up dot products of correctly identified common lines is simply a sum of ones. We may consider the contribution of the falsely detected common lines as a random walk on the real line, where steps to the left and to the right are equally probable. From this interpretation it follows that the total contribution of the misidentified common lines to the objective function (2.9) is proportional to the square root of the number of misidentifications, whereas the contribution of the correctly identified common lines is linear. This square-root diminishing effect of the misidentifications makes the global optimization (2.9) extremely robust compared with the least squares approach, which is much more sensitive because its objective function is dominated by the misidentifications.

These intuitive arguments regarding the statistical attractiveness of the optimization problem (2.9) will be later put on a firm mathematical ground using random matrix theory as elaborated in Section 7. Still, in order for the optimization problem (2.9) to be of any practical use, we must show that its solution can be efficiently computed. We note that our objective function is closely related to the SCE of Farrow and Ottensmeyer [9, p. 1754, eq. (6)], given by

$$SCE = \sum_{i \neq j} \arccos\left(R_i c_{ij} \cdot R_j c_{ji}\right). \qquad (2.10)$$

This SCE was introduced and used in [9] to measure the success of their quaternion-based sequential iterative angular reconstitution methods. At the little price of deleting the well-behaved monotonic nonlinear arccos function in (2.10) we arrive at (2.9),

which, as we will soon show, has the great advantage of being amenable to efficient global non-sequential optimization by either spectral or semidefinite programming relaxations.

**3. Eigenvector relaxation.** The objective function in (2.9) is quadratic in the unknown rotations $R_1, \ldots, R_N$, which means that if the constraints (2.8) are properly relaxed, then the solution to the maximization problem (2.9) would be related to the top eigenvectors of the matrix defining the quadratic form. In this section we give a precise definition of that matrix, and show how the unknown rotations can be recovered from its top three eigenvectors.

We first define the following four $N \times N$ matrices $S^{11}, S^{12}, S^{21}$, and $S^{22}$ using all available common-line data (2.4)-(2.5) as follows:

$$S_{ij}^{11} = x_{ij}x_{ji}, \quad S_{ij}^{12} = x_{ij}y_{ji}, \quad S_{ij}^{21} = y_{ij}x_{ji}, \quad S_{ij}^{22} = y_{ij}y_{ji}, \tag{3.1}$$

for $1 \le i \ne j \le N$, while their diagonals are set to zero

$$S_{ii}^{11} = S_{ii}^{12} = S_{ii}^{21} = S_{ii}^{22} = 0, \quad i = 1, \ldots, N.$$

Clearly, $S^{11}$ and $S^{22}$ are symmetric matrices ($S^{11} = S^{11^T}$ and $S^{22} = S^{22^T}$), while $S^{12} = S^{21^T}$. It follows that the $2N \times 2N$ matrix $S$ given by

$$S = \begin{pmatrix} S^{11} & S^{12} \\ S^{21} & S^{22} \end{pmatrix} \tag{3.2}$$

is symmetric ($S = S^T$) and storing all available common line information. More importantly, the top eigenvectors of $S$ will reveal all rotations in a manner we describe below.

We denote the columns of the rotation matrix $R_i$ by $R_i^1$, $R_i^2$ and $R_i^3$, and write the rotation matrices as

$$R_i = \begin{pmatrix} | & | & | \\ R_i^1 & R_i^2 & R_i^3 \\ | & | & | \end{pmatrix} = \begin{pmatrix} x_i^1 & x_i^2 & x_i^3 \\ y_i^1 & y_i^2 & y_i^3 \\ z_i^1 & z_i^2 & z_i^3 \end{pmatrix}, \quad i = 1, \ldots, N. \tag{3.3}$$

Only the first two columns of the $R_i$'s need to be recovered, because the third columns are given by the cross product: $R_i^3 = R_i^1 \times R_i^2$. We therefore need to recover the six $N$-dimensional coordinate vectors $x^1, y^1, z^1, x^2, y^2, z^2$ that are defined by

$$x^1 = (x_1^1 \ x_2^1 \ \cdots \ x_N^1)^T, \quad y^1 = (y_1^1 \ y_2^1 \ \cdots \ y_N^1)^T, \quad z^1 = (z_1^1 \ z_2^1 \ \cdots \ z_N^1)^T, \tag{3.4}$$

$$x^2 = (x_1^2 \ x_2^2 \ \cdots \ x_N^2)^T, \quad y^2 = (y_1^2 \ y_2^2 \ \cdots \ y_N^2)^T, \quad z^2 = (z_1^2 \ z_2^2 \ \cdots \ z_N^2)^T. \tag{3.5}$$

Alternatively, we need to find the following three $2N$-dimensional vectors $x$, $y$ and $z$

$$x = \begin{pmatrix} x^1 \\ x^2 \end{pmatrix}, \quad y = \begin{pmatrix} y^1 \\ y^2 \end{pmatrix}, \quad z = \begin{pmatrix} z^1 \\ z^2 \end{pmatrix}. \tag{3.6}$$

Using this notation we rewrite the objective function (2.9) as

$$\sum_{i \ne j} R_i c_{ij} \cdot R_j c_{ji} = x^T S x + y^T S y + z^T S z, \tag{3.7}$$

6

which is a result of the following index manipulation

$$\sum_{i \neq j} R_i c_{ij} \cdot R_j c_{ji} = \sum_{i \neq j} x_{ij} x_{ji} R_i^1 \cdot R_j^1 + x_{ij} y_{ji} R_i^1 \cdot R_j^2 + y_{ij} x_{ji} R_i^2 \cdot R_j^1 + y_{ij} y_{ji} R_i^2 \cdot R_j^2$$

$$= \sum_{i \neq j} S_{ij}^{11} R_i^1 \cdot R_j^1 + S_{ij}^{12} R_i^1 \cdot R_j^2 + S_{ij}^{21} R_i^2 \cdot R_j^1 + S_{ij}^{22} R_i^2 \cdot R_j^2 \qquad (3.8)$$

$$= \sum_{i,j} S_{ij}^{11} (x_i^1 x_j^1 + y_i^1 y_j^1 + z_i^1 z_j^1) + S_{ij}^{12} (x_i^1 x_j^2 + y_i^1 y_j^2 + z_i^1 z_j^2) +$$

$$S_{ij}^{21} (x_i^2 x_j^1 + y_i^2 y_j^1 + z_i^2 z_j^1) + S_{ij}^{22} (x_i^2 x_j^2 + y_i^2 y_j^2 + z_i^2 z_j^2)$$

$$= x^{1^T} S^{11} x^1 + y^{1^T} S^{11} y^1 + z^{1^T} S^{11} z^1 +$$

$$x^{1^T} S^{12} x^2 + y^{1^T} S^{12} y^2 + z^{1^T} S^{12} z^2 +$$

$$x^{2^T} S^{21} x^1 + y^{2^T} S^{21} y^1 + z^{2^T} S^{21} z^1 +$$

$$x^{2^T} S^{22} x^2 + y^{2^T} S^{22} y^2 + z^{2^T} S^{22} z^2$$

$$= x^T S x + y^T S y + z^T S z. \qquad (3.9)$$

The equality (3.7) shows that the maximization problem (2.9) is equivalent to the maximization problem

$$\max_{R_1, \dots, R_N} x^T S x + y^T S y + z^T S z, \qquad (3.10)$$

subject to the constraints (2.8). In order to make this optimization problem tractable, we relax the constraints and look for the solution of the proxy maximization problem

$$\max_{\|x\|=1} x^T S x. \qquad (3.11)$$

The connection between the solution to (3.11) and that of (3.10) will be made shortly. Since $S$ is a symmetric matrix it has a complete set of orthonormal eigenvectors $\{v^1, \dots, v^{2N}\}$ satisfying

$$S v^n = \lambda_n v^n, \quad n = 1, \dots, 2N,$$

with real eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{2N}.$$

The solution to the maximization problem (3.11) is therefore given by the top eigenvector $v^1$ with largest eigenvalue $\lambda_1$

$$v^1 = \operatorname*{argmax}_{\|x\|=1} x^T S x. \qquad (3.12)$$

If the unknown rotations are sampled from the uniform distribution (Haar measure) over SO(3), that is, when the molecule has no preferred orientation, then the largest eigenvalue should have multiplicity three, corresponding to the vectors $x$, $y$ and $z$, as the symmetry of the problem in this case suggests that there is nothing special about $x$. We therefore expect to recover the first two columns of the rotation matrices $R_1, \dots, R_N$ from the top three eigenvectors $v^1, v^2, v^3$ of $S$. This recovery is performed

by constructing for every $i = 1, \ldots, N$ a $3 \times 3$ matrix $A_i = \begin{pmatrix} | & | & | \\ A_i^1 & A_i^2 & A_i^3 \\ | & | & | \end{pmatrix}$ whose columns are given by

$$A_i^1 = \begin{pmatrix} v_i^1 \\ v_i^2 \\ v_i^3 \end{pmatrix}, \quad A_i^2 = \begin{pmatrix} v_{N+i}^1 \\ v_{N+i}^2 \\ v_{N+i}^3 \end{pmatrix}, \quad A_i^3 = A_i^1 \times A_i^2. \tag{3.13}$$

The matrix $A_i$ is not necessarily a rotation, so we recover $R_i$ as its closest rotation matrix via the well know procedure [16]: $R_i = U_i V_i^T$, where $A_i = U_i \Sigma_i V_i^T$ is the singular value decomposition of $A_i$.

From the computational point of view, we note that a simple way of computing the top three eigenvectors is using the iterative power method, where three initial randomly chosen vectors are repeatedly multiplied be the matrix $S$ and then orthonormalized by the Gram-Schmidt (QR) procedure until convergence. The number of iterations required by such a procedure is determined by the spectral gap between the third and forth eigenvalues. The spectral gap is further discussed in Sections 5 and 7. In practice, for large values of $N$ we use MATLAB's eigs function to compute the few top eigenvectors, while for small $N$ we compute all eigenvectors using MATLAB's eig function. We remark that the computational bottleneck for large $N$ is the storage of the $2N \times 2N$ matrix $S$ rather than the time complexity of computing the top eigenvectors.

**4. Relaxation by a semidefinite program.** In this Section we present an alternative relaxation of (2.9) using semidefinite programming (SDP) [13], which draws similarities with the Goemans-Williamson SDP for finding the maximum cut in a weighted graph [14]. The relaxation of the SDP is tighter than the eigenvector relaxation and does not require the assumption that the rotations are uniformly sampled over SO(3).

The SDP formulation begins with the introduction of two $3 \times N$ matrices $R^1$ and $R^2$ defined by concatenating the first columns and second columns of the $N$ rotation matrices, respectively,

$$R^1 = \begin{pmatrix} | & | & & | \\ R_1^1 & R_2^1 & \cdots & R_N^1 \\ | & | & & | \end{pmatrix}, \quad R^2 = \begin{pmatrix} | & | & & | \\ R_1^2 & R_2^2 & \cdots & R_N^2 \\ | & | & & | \end{pmatrix}. \tag{4.1}$$

We also concatenate $R^1$ and $R^2$ to define a $3 \times 2N$ matrix $R$ given by

$$R = (R^1 \ R^2) = \begin{pmatrix} | & | & & | & | & | & & | \\ R_1^1 & R_2^1 & \cdots & R_N^1 & R_1^2 & R_2^2 & \cdots & R_N^2 \\ | & | & & | & | & | & & | \end{pmatrix}. \tag{4.2}$$

The Gram matrix $G$ for the matrix $R$ is a $2N \times 2N$ matrix of inner products between the three-dimensional column vectors of $R$, that is,

$$G = R^T R. \tag{4.3}$$

Clearly, $G$ is a rank-3 semidefinite positive matrix ($G \succeq 0$), which can be conveniently written as a block matrix

$$G = \begin{pmatrix} G^{11} & G^{12} \\ G^{21} & G^{22} \end{pmatrix} = \begin{pmatrix} R^{1^T} R^1 & R^{1^T} R^2 \\ R^{2^T} R^1 & R^{2^T} R^2 \end{pmatrix}. \tag{4.4}$$

The orthogonality of the rotation matrices ($R_i^T R_i = I$) implies that

$$G_{ii}^{11} = G_{ii}^{22} = 1, \quad i = 1, 2, \ldots, N, \tag{4.5}$$

and

$$G_{ii}^{12} = G_{ii}^{21} = 0, \quad i = 1, 2, \ldots, N. \tag{4.6}$$

From equation (3.8) it follows that the objective function (2.9) is the trace of the matrix product $SG$:

$$\sum_{i \neq j} R_i c_{ij} \cdot R_j c_{ji} = \text{trace}(SG). \tag{4.7}$$

A natural relaxation of the optimization problem (2.9) is thus given by the SDP

$$\max_{G \in \mathbb{R}^{2N \times 2N}} \text{trace}(SG) \tag{4.8}$$

$$\text{s.t. } G \succeq 0, \tag{4.9}$$

$$G_{ii}^{11} = G_{ii}^{22} = 1, \quad G_{ii}^{12} = G_{ii}^{21} = 0, \quad i = 1, 2, \ldots, N. \tag{4.10}$$

The only constraint missing in this SDP formulation is the non-convex rank-3 constraint on the Gram matrix $G$. The matrix $R$ is recovered from the Cholesky decomposition of the solution $G$ of the SDP (4.8)-(4.10). If the rank of $G$ is greater than 3, then we project the rows of $R$ onto the subspace spanned by the top three eigenvectors of $G$, and recover the rotations using the procedure that was detailed in the previous section in (3.13). We note that except for the orthogonality constraint (4.6), the semidefinite program (4.8)-(4.10) is identical to the Goemans-Williamson SDP for finding the maximum cut in a weighted graph [14].

From the complexity point of view, SDP can be solved in a polynomial time to any given precision, but even the most sophisticated SDP solvers that exploit the sparsity structure of the max cut problem are not competitive with the much faster eigenvector method. At first glance it may seem that the SDP (4.8)-(4.10) should outperform the eigenvector method in terms of producing more accurate rotation matrices. However, our simulations show that the accuracy of both methods is almost identical when the rotations are uniformly sampled over SO(3). As the eigenvector method is much faster, it should also be the method of choice, whenever the rotations are a-priori known to be uniformly sampled.

**5. Numerical simulations.** We performed several numerical experiments that illustrate the robustness of the eigenvector and the SDP methods to false identifications of common-lines. All simulations were performed in MATLAB on a Lenovo Thinkpad X300 laptop with Intel(R) Core(TM)2 CPU L7100 1.2GHz with 4GB RAM running Windows Vista.

**5.1. Experiments with simulated rotations.** In the first series of simulations we tried to imitate the experimental setup by using the following procedure. In each simulation, we randomly sampled $N$ rotations from the uniform distribution over SO(3). This was done by randomly sampling $N$ vectors in $\mathbb{R}^4$ whose coordinates are i.i.d Gaussians, followed by normalizing these vectors to the unit three-dimensional sphere $S^3 \subset \mathbb{R}^4$. The normalized vectors are viewed as unit quaternions which we converted into $3 \times 3$ rotation matrices $R_1, \ldots, R_N$. We then computed all pairwise

common-line vectors $c_{ij} = R_i^{-1} \frac{R_i^3 \times R_j^3}{\|R_i^3 \times R_j^3\|}$ and $c_{ji} = R_j^{-1} \frac{R_i^3 \times R_j^3}{\|R_i^3 \times R_j^3\|}$ (see also the discussion following (6.2)). For each pair of rotations, with probability $p$ we kept the values of $c_{ij}$ and $c_{ji}$ unchanged, while with probability $1 - p$ we replaced $c_{ij}$ and $c_{ji}$ by two random vectors that were sampled from the uniform distribution over the unit circle in the plane. The parameter $p$ ranges from 0 to 1 and indicates the proportion of the correctly detected common lines. For example, $p = 0.1$ means that only 10% of the common lines are identified correctly, and all other 90% entries of the matrix $S$ are filled in with random entries corresponding to some randomly chosen unit vectors.

Figure 5.1 shows the distribution of the eigenvalues of the matrix $S$ for two different values of $N$ and four different values of the probability $p$. It took a matter of seconds to compute each of the eigenvalue histograms shown in Figure 5.1. Evident from the eigenvalue histograms is the spectral gap between the three largest eigenvalues and the remaining eigenvalues, as long as $p$ is not too small. As $p$ decreases, the spectral gap narrows down, until it completely disappears at some critical value $p_c$, which we call the threshold probability. Figure 5.1 indicates that the value of the critical probability for $N = 100$ is somewhere between 0.1 and 0.25, whereas for $N = 500$ it is bounded between 0.05 and 0.1. As a result, the algorithm can cope with a higher percentage of misidentifications by using more images (larger $N$). In particular, for $N = 500$ it can deal with as many as 90% misidentified common lines.

When $p$ decreases, not only does the gap narrows, but also the histogram of the eigenvalues becomes smoother. The smooth part of the histogram seems to follows the semi-circle law of Wigner [17, 18], as illustrated in Figure 5.1. The support of the semi-circle gets larger as $p$ decreases. In the next sections we will provide a mathematical explanation for the numerically observed eigenvalue histograms and for the emergence of Wigner's semi-circle.
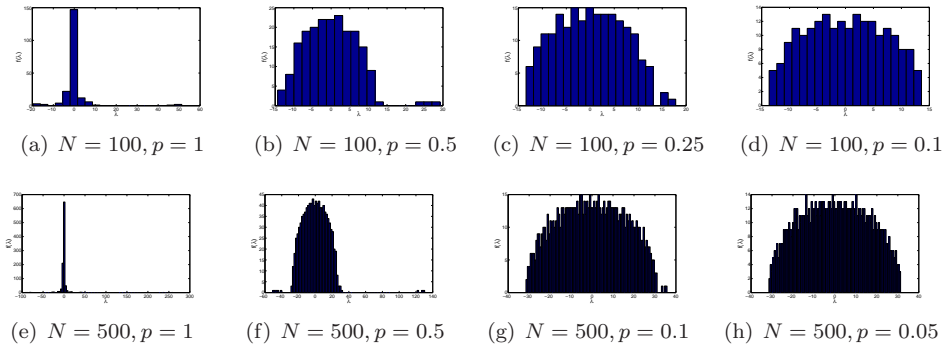


(a) $N = 100, p = 1$    (b) $N = 100, p = 0.5$    (c) $N = 100, p = 0.25$    (d) $N = 100, p = 0.1$

(e) $N = 500, p = 1$    (f) $N = 500, p = 0.5$    (g) $N = 500, p = 0.1$    (h) $N = 500, p = 0.05$

FIG. 5.1. *Eigenvalue histograms for the matrix S for different values of N and p.*

A further investigation into the results of the numerical simulations also revealed that the rotations that were recovered by the top three eigenvectors were highly correlated with the sampled rotations, as long as $p$ was above the threshold probability $p_c$. The correlation between the rotations estimated by the eigenvector method, denoted $\hat{R}_1, \ldots, \hat{R}_N$, and the true sampled rotations $R_1, \ldots, R_N$, is defined in the following manner. First, note that (2.6) implies that the true rotations can be recovered only up to a fixed $3 \times 3$ orthogonal transformation $O$, since if $R_i c_{ij} = R_j c_{ji}$, then also $OR_i c_{ij} = OR_j c_{ji}$. In other words, a completely successful recovery is a recovery for which $\hat{R}_i^{-1} R_i = O$, for all $i = 1, \ldots, N$ for some fixed orthogonal matrix $O$. The

10

success of the recovery procedure can therefore be measured by the closeness of the matrix

$$\hat{O} = \frac{1}{N} \sum_{i=1}^{N} \hat{R}_i^T R_i \qquad (5.1)$$

to being an orthogonal matrix, and we define the $3 \times 3$ "correlation" matrix $\rho_+$ by

$$\rho_+ = \hat{O}\hat{O}^T = \frac{1}{N^2} \left( \sum_{i=1}^{N} \hat{R}_i^T R_i \right) \left( \sum_{i=1}^{N} \hat{R}_i^T R_i \right)^T . \qquad (5.2)$$

The recovery is considered successful if $\rho_+$ is close to the $3 \times 3$ identity matrix $I$.

However, there is another degree of freedom which amounts to a global reflection, which the cryo-EM expert may immediately recognize as the chirality ambiguity of the reconstructed molecule. Indeed, since the third coordinate of all $c_{ij}$'s is zero, multiplying the third column of the $R_i$'s by $-1$ leads to yet another valid solution of the maximization problem solved by the eigenvector method. We therefore define a second set of matrices $\tilde{R}_1, \ldots, \tilde{R}_N$ by

$$\tilde{R}_i = \hat{R}_i \begin{pmatrix} 1 & & \\ & 1 & \\ & & -1 \end{pmatrix}, \quad i = 1, \ldots, N, \qquad (5.3)$$

and their corresponding correlation matrix $\rho_-$ by

$$\rho_- = \frac{1}{N^2} \left( \sum_{i=1}^{N} \tilde{R}_i^T R_i \right) \left( \sum_{i=1}^{N} \tilde{R}_i^T R_i \right)^T . \qquad (5.4)$$

Indeed, in all of our successful experiments one of the correlation matrices, either $\rho_+$ or $\rho_-$ was close to the identity matrix and so we define the correlation matrix $\rho$ as

$$\rho = \operatorname*{argmax}_{\rho_+, \rho_-} \{\operatorname{trace}(\rho_+), \operatorname{trace}(\rho_-)\}. \qquad (5.5)$$

Table 5.1 compares the correlation matrices that were obtained by the eigenvector method with the ones obtained by the SDP method for $N = 100$ with the same common lines input data. Table 5.2 gives a similar comparison for $N = 500$. The SDP was solved using SDPT3 [19, 20] in MATLAB.

**5.2. Experiments with simulated noisy projections.** In the second series of experiments, we tested the eigenvector method on simulated noisy projection images of a ribosome, for different numbers of projections ($N = 100, 500, 1000$) and different levels of noise. For each $N$, we generated $N$ noise-free centered projections of the ribosome, whose corresponding rotations were uniformly distributed on SO(3). Each projection was of size $129 \times 129$ pixels. Next, we fixed a signal-to-noise ratio (SNR), and added to each clean projection additive Gaussian white noise of the prescribed SNR. The SNR in all our experiments is defined by

$$\mathrm{SNR} = \frac{\mathrm{Var}(Signal)}{\mathrm{Var}(Noise)}, \qquad (5.6)$$

where Var is the variance (energy), $Signal$ is the clean projection image and $Noise$ is the noise realization of that image. Figure 5.2 shows one of the projections at

11

| $p$ | $\rho_{eig}$ | $\rho_{sdp}$ |
|---|---|---|
| 1 | $\begin{pmatrix} 0.9960 & 0.0011 & 0.0014 \\ 0.0011 & 0.9926 & 0.0011 \\ 0.0014 & 0.0011 & 0.9930 \end{pmatrix}$ | $\begin{pmatrix} 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.0000 & 0.0000 \\ 0.0000 & 0.0000 & 1.0000 \end{pmatrix}$ |
| 0.5 | $\begin{pmatrix} 0.9692 & -0.0016 & -0.0001 \\ -0.0016 & 0.9717 & -0.0003 \\ -0.0001 & -0.0003 & 0.9712 \end{pmatrix}$ | $\begin{pmatrix} 0.9723 & -0.0015 & -0.0004 \\ -0.0015 & 0.9718 & -0.0012 \\ -0.0004 & -0.0012 & 0.9721 \end{pmatrix}$ |
| 0.25 | $\begin{pmatrix} 0.6796 & 0.0088 & -0.0356 \\ 0.0088 & 0.6538 & -0.0296 \\ -0.0356 & -0.0296 & 0.6356 \end{pmatrix}$ | $\begin{pmatrix} 0.7108 & 0.0029 & -0.0427 \\ 0.0029 & 0.6419 & -0.0323 \\ -0.0427 & -0.0323 & 0.6414 \end{pmatrix}$ |
| 0.15 | $\begin{pmatrix} 0.0868 & 0.0077 & 0.0290 \\ 0.0077 & 0.1220 & -0.0561 \\ 0.0290 & -0.0561 & 0.1470 \end{pmatrix}$ | $\begin{pmatrix} 0.0657 & 0.0155 & 0.0001 \\ 0.0155 & 0.0739 & -0.0772 \\ 0.0001 & -0.0772 & 0.1648 \end{pmatrix}$ |

TABLE 5.1

*Correlation matrices for $N = 100$ and different values of p using the eigenvector method and the SDP method.*

| $p$ | $\rho_{eig}$ | $\rho_{sdp}$ |
|---|---|---|
| 1 | $\begin{pmatrix} 0.9994 & 0.0000 & -0.0001 \\ 0.0000 & 0.9996 & -0.0001 \\ -0.0001 & -0.0001 & 0.9995 \end{pmatrix}$ | $\begin{pmatrix} 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 1.0000 & 0.0000 \\ 0.0000 & 0.0000 & 1.0000 \end{pmatrix}$ |
| 0.5 | $\begin{pmatrix} 0.9946 & -0.0000 & -0.0003 \\ -0.0000 & 0.9945 & -0.0002 \\ -0.0003 & -0.0002 & 0.9947 \end{pmatrix}$ | $\begin{pmatrix} 0.9952 & -0.0001 & -0.0001 \\ -0.0001 & 0.9951 & -0.0002 \\ -0.0001 & -0.0002 & 0.9950 \end{pmatrix}$ |
| 0.1 | $\begin{pmatrix} 0.5782 & -0.0188 & -0.0082 \\ -0.0188 & 0.5808 & 0.0023 \\ -0.0082 & 0.0023 & 0.5481 \end{pmatrix}$ | $\begin{pmatrix} 0.6194 & -0.0247 & -0.0078 \\ -0.0247 & 0.6175 & 0.0097 \\ -0.0078 & 0.0097 & 0.5880 \end{pmatrix}$ |
| 0.05 | $\begin{pmatrix} 0.0047 & -0.0057 & -0.0069 \\ -0.0057 & 0.0197 & -0.0066 \\ -0.0069 & -0.0066 & 0.0307 \end{pmatrix}$ | $\begin{pmatrix} 0.0097 & -0.0068 & -0.0014 \\ -0.0068 & 0.0215 & 0.0019 \\ -0.0014 & 0.0019 & 0.0003 \end{pmatrix}$ |

TABLE 5.2

*Correlation matrices for $N = 500$ and different values of p using the eigenvector method and the SDP method.*

different SNR levels. The SNR values used throughout this experiment were $2^{-k}$ with $k = 0, \ldots, 9$. Clean projections were generated by setting SNR $= 2^{20}$.

We computed the 2D Fourier transform of all projections on a polar grid discretized into $L = 72$ central lines, corresponding to an angular resolution of $360°/72 = 5°$. We constructed the matrix $S$ according to (3.1)-(3.2) by comparing all $\binom{N}{2}$ pairs of projection images; for each pair we detected the common line by computing all $L^2/2$ possible different correlations between their Fourier central lines, of which the pair of central lines having the maximum correlation was declared as the common-line. Table 5.3 shows the proportion $p$ of the correctly detected common lines as a function of the SNR (we consider a common line as correctly identified if each of the estimated direction vectors $(x_{ij}, y_{ij})$ and $(x_{ji}, y_{ji})$ is within $10°$ of its true direction). As expected, the proportion $p$ is a decreasing function of the SNR.

We used MATLAB's eig function to compute the eigenvalue histograms of all $S$ matrices as shown in Figures 5.3-5.5. There is a clear resemblance between the
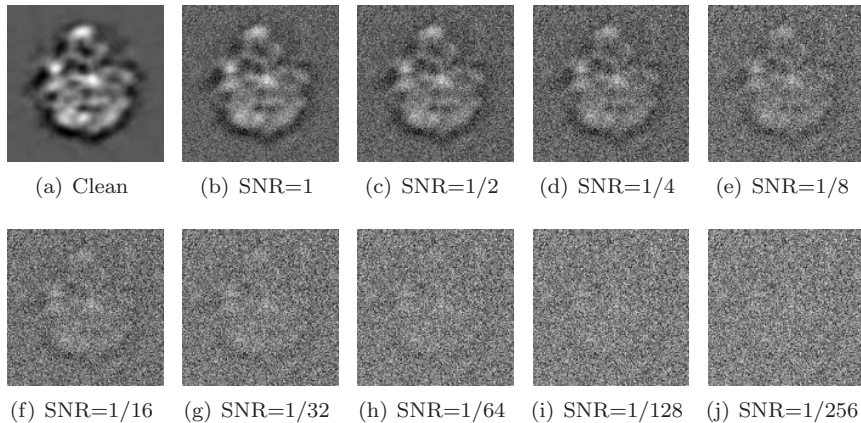
(a) Clean     (b) SNR=1     (c) SNR=1/2     (d) SNR=1/4     (e) SNR=1/8

(f) SNR=1/16     (g) SNR=1/32     (h) SNR=1/64     (i) SNR=1/128     (j) SNR=1/256

FIG. 5.2. *Simulated projection with various levels of additive Gaussian white noise.*

(a) $N = 100$

| $\log_2(\text{SNR})$ | $p$ |
|---|---|
| 20 | 0.997 |
| 0 | 0.968 |
| -1 | 0.930 |
| -2 | 0.828 |
| -3 | 0.653 |
| -4 | 0.444 |
| -5 | 0.247 |
| -6 | 0.108 |
| -7 | 0.046 |
| -8 | 0.023 |
| -9 | 0.017 |

(b) $N = 500$

| $\log_2(\text{SNR})$ | $p$ |
|---|---|
| 20 | 0.997 |
| 0 | 0.967 |
| -1 | 0.922 |
| -2 | 0.817 |
| -3 | 0.639 |
| -4 | 0.433 |
| -5 | 0.248 |
| -6 | 0.113 |
| -7 | 0.046 |
| -8 | 0.023 |
| -9 | 0.015 |

(c) $N = 1000$

| $\log_2(\text{SNR})$ | $p$ |
|---|---|
| 20 | 0.997 |
| 0 | 0.966 |
| -1 | 0.919 |
| -2 | 0.813 |
| -3 | 0.638 |
| -4 | 0.437 |
| -5 | 0.252 |
| -6 | 0.115 |
| -7 | 0.047 |
| -8 | 0.023 |
| -9 | 0.015 |

TABLE 5.3

*The proportion p of correctly detected common lines as a function of the SNR. As expected, p is not a function of the number of images N.*

eigenvalue histograms of the noisy $S$ matrices shown in Figure 5.1 and those shown in Figures 5.3-5.5. One noticeable difference is that the top three eigenvalues in Figures 5.3-5.5 tend to spread (note, for example, the spectral gap between the top three eigenvalues in Figure 5.3(e)), whereas in Figure 5.1 they tend to stick together. We attribute this spreading effect to the fact that the model used in Section 5.1 is too simplified. First, the detection of common lines tends to be more successful when computed between projections that have more pronounced signal features. This means that the assumption that each common line is detected correctly with a fixed probability $p$ is too restrictive. Second, falsely detected common lines are far from being random. The correct common line is often confused with a Fourier central line that is similar to it; it is not just confused with any Fourier central line with equal probability. Still, despite the simplified assumptions that were made in Section 5.1 to model the matrix $S$, the resulting eigenvalue histograms are very similar. We note that a related spreading effect of eigenvalues was recently examined in [21], where it was shown that the top singular values of low-rank matrices undergo similar spreading
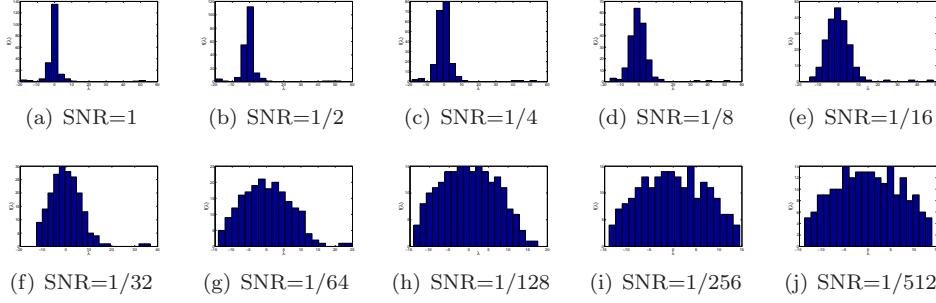
(a) SNR=1    (b) SNR=1/2    (c) SNR=1/4    (d) SNR=1/8    (e) SNR=1/16

(f) SNR=1/32    (g) SNR=1/64    (h) SNR=1/128    (i) SNR=1/256    (j) SNR=1/512

FIG. 5.3. *Eigenvalue histograms of $S$ for $N = 100$ and different levels of noise.*



(a) SNR=1    (b) SNR=1/2    (c) SNR=1/4    (d) SNR=1/8    (e) SNR=1/16

(f) SNR=1/32    (g) SNR=1/64    (h) SNR=1/128    (i) SNR=1/256    (j) SNR=1/512

FIG. 5.4. *Eigenvalue histograms of $S$ for $N = 500$ and different levels of noise.*

when replacing randomly chosen entries of the low-rank matrix by zeros (see, for example, [21, page 3, Figure 1]).

Perhaps the most important consequence of our numerical simulations is that increasing the number of projections $N$ distinguishes the top three eigenvalues from the bulk of the spectrum (the semi-circle). For example, for $N = 100$ the top eigenvalues are clearly distinguished from the bulk for SNR $= 1/32$, while for $N = 500$ they can be distinguished for SNR $= 1/128$ (maybe even at SNR $= 1/256$), and for $N = 1000$ they are distinguished even at the most extreme noise level of SNR $= 1/512$. From the practical point of view, this means that 3D reconstruction is possible even at extreme levels of noise, provided that enough projections are available. Even if reconstruction is not possible from the raw noisy images, the eigenvector and SDP methods should allow to use class averages consisting of fewer images.

**6. The matrix $S$ as a convolution operator on SO(3).** Taking an even closer look into the numerical distribution of the eigenvalues of the "clean" $2N \times 2N$ matrix $S^{\text{clean}}$ corresponding to $p = 1$ (all common-lines detected correctly) reveals that its eigenvalues have the exact same multiplicities as the spherical harmonics, which are the eigenfunctions of the Laplacian on the unit sphere $S^2 \subset \mathbb{R}^3$. In particular, Figure 6.1(a) is a bar plot of the largest 50 eigenvalues of $S^{\text{clean}}$ with $N = 1000$, and is clearly showing numerical multiplicities of $3, 7, 11, \ldots$, corresponding to the multiplicity $2l + 1$ ($l = 1, 3, 5, \ldots$) of the odd spherical harmonics. Moreover, Figure 6.1(b) is a bar plot of the magnitude of the most negative eigenvalues of $S$. The multiplicities $5, 9, 13, \ldots$ corresponding to the multiplicity $2l + 1$ ($l = 2, 4, 6, \ldots$) of the even spherical harmonics are evident (the first even eigenvalue corresponding to $l = 0$ is missing).
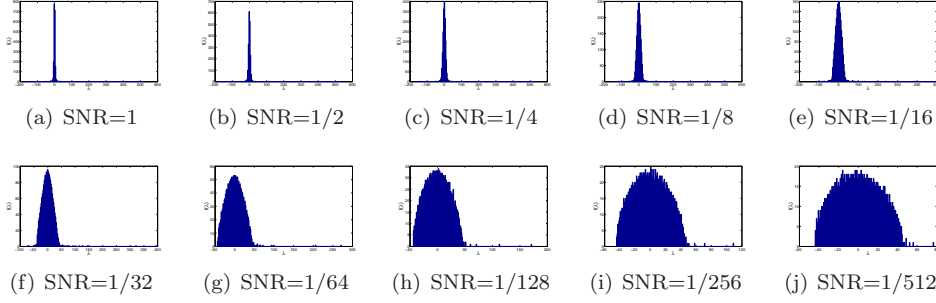
14

(a) SNR=1    (b) SNR=1/2    (c) SNR=1/4    (d) SNR=1/8    (e) SNR=1/16

(f) SNR=1/32    (g) SNR=1/64    (h) SNR=1/128    (i) SNR=1/256    (j) SNR=1/512

FIG. 5.5. *Eigenvalue histograms of $S$ for $N = 1000$ and different levels of noise.*



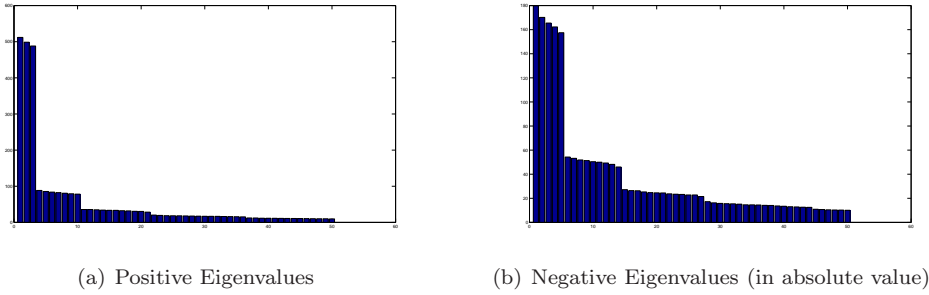(a) Positive Eigenvalues    (b) Negative Eigenvalues (in absolute value)

FIG. 6.1. *Bar plot of the positive (left) and the absolute values of the negative (right) eigenvalues of $S$ with $N = 1000$ and $p = 1$. The numerical multiplicities $2l + 1$ $(l = 1, 2, 3, \ldots)$ of the spherical harmonics are evident, with odd $l$ values corresponding to positive eigenvalues, and even $l$ values (except $l = 0$) corresponding to negative eigenvalues.*

The numerically observed multiplicities motivate us to examine $S^{\text{clean}}$ in more detail. To that end, it is more convenient to reshuffle the $2N \times 2N$ matrix $S$ defined in (3.1)-(3.2) into an $N \times N$ matrix $K$ whose entries are $2 \times 2$ rank-1 matrices given by

$$K_{ij} = \begin{pmatrix} x_{ij}x_{ji} & x_{ij}y_{ji} \\ y_{ij}x_{ji} & y_{ij}y_{ji} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} c_{ij}c_{ji}^T \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}^T, \quad i, j = 1, \ldots, N,$$

(6.1)

with $c_{ij}$ and $c_{ji}$ given in (2.4)-(2.5). From (2.6) it follows that the common line is given by the normalized cross product of $R_i^3$ and $R_j^3$, that is,

$$R_i c_{ij} = R_j c_{ji} = \pm \frac{R_i^3 \times R_j^3}{\|R_i^3 \times R_j^3\|},$$

(6.2)

because $R_i c_{ij}$ is a linear combination of $R_i^1$ and $R_i^2$ (perpendicular to $R_i^3$), while $R_j c_{ji}$ is a linear combination of $R_j^1$ and $R_j^2$ (perpendicular to $R_j^3$); a unit vector perpendicular to $R_i^3$ and $R_j^3$ must be given by either $\frac{R_i^3 \times R_j^3}{\|R_i^3 \times R_j^3\|}$ or $-\frac{R_i^3 \times R_j^3}{\|R_i^3 \times R_j^3\|}$. Equations (6.1)-(6.2) imply that $K_{ij}$ is a function of $R_i$ and $R_j$ given by

$$K_{ij} = K(R_i, R_j) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} R_i^{-1} \frac{(R_i^3 \times R_j^3)(R_i^3 \times R_j^3)^T}{\|R_i^3 \times R_j^3\|^2} R_j \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}^T,$$

(6.3)

15

for $i \neq j$ regardless of the choice of the sign in (6.2), and $K_{ii} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$.

The eigenvalues of $K$ and $S$ are the same, with the eigenvectors of $K$ being vectors of length $2N$ obtained from the eigenvectors of $S$ by reshuffling their entries. We therefore try to understand the operation of matrix-vector multiplication of $K$ applied to some arbitrary vector $f$ of length $2N$. It is convenient to view the vector $f$ as $N$ vectors in $\mathbb{R}^2$ obtained by sampling the function $f : \mathrm{SO}(3) \to \mathbb{R}^2$ at $R_1, \ldots, R_N$, that is,

$$f_i = f(R_i), \quad i = 1, \ldots, N. \tag{6.4}$$

The matrix-vector multiplication is thus given by

$$(Kf)_i = \sum_{j=1}^{N} K_{ij} f_j = \sum_{j=1}^{N} K(R_i, R_j) f(R_j), \quad i = 1, \ldots, N. \tag{6.5}$$

If the rotations $R_1, \ldots, R_N$ are i.i.d random variables uniformly distributed over $\mathrm{SO}(3)$, then the expected value of $(Kf)_i$ conditioned on $R_i$ is

$$\mathbb{E}\left[(Kf)_i \mid R_i\right] = (N-1) \int_{\mathrm{SO}(3)} K(R_i, R) f(R) \, dR, \tag{6.6}$$

where $dR$ is the Haar measure (recall that by being a zero matrix, $K(R_i, R_i)$ does not contribute to the sum in (6.5)). The eigenvectors of $K$ are therefore discrete approximations to the eigenfunctions of the integral operator $\mathcal{K}$ given by

$$(\mathcal{K}f)(R) = \int_{\mathrm{SO}(3)} K(R_1, R_2) f(R_2) \, dR_2, \tag{6.7}$$

due to the law of large numbers, with the kernel $K : \mathrm{SO}(3) \times \mathrm{SO}(3) \to \mathbb{R}^{2 \times 2}$ given by (6.3). We are thus interested in the eigenfunctions of the integral operator $\mathcal{K} : L_2\left(\mathrm{SO}(3), dR\right)^2 \to L_2\left(\mathrm{SO}(3), dR\right)^2$ given by (6.7).

The integral operator $\mathcal{K}$ is a convolution operator over $\mathrm{SO}(3)$. Indeed, note that $K$ given in (6.3) satisfies

$$K(gR_1, gR_2) = K(R_1, R_2), \quad \text{for all } g \in \mathrm{SO}(3), \tag{6.8}$$

because $(gR_1^3) \times (gR_2^3) = g(R_1^3 \times R_2^3)$ and $gg^{-1} = g^{-1}g = I$. It follows that the kernel $K$ depends only upon the "ratio" $R_1^{-1} R_2$, because we can choose $g = R_1^{-1}$ so that

$$K(R_1, R_2) = K(I, R_1^{-1} R_2),$$

and the integral operator $\mathcal{K}$ of (6.7) becomes

$$(\mathcal{K}f)(R_1) = \int_{\mathrm{SO}(3)} K(I, R_1^{-1} R_2) f(R_2) \, dR_2. \tag{6.9}$$

We will therefore define the convolution kernel $\tilde{K} : \mathrm{SO}(3) \to \mathbb{R}^{2 \times 2}$ as

$$\tilde{K}(U^{-1}) \equiv K(I, U) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \frac{(I^3 \times U^3)(I^3 \times U^3)^T}{\|I^3 \times U^3\|^2} U \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}^T, \tag{6.10}$$

16

where $I^3 = (0\ 0\ 1)^T$ is the third column of the identity matrix $I$. We rewrite the integral operator $\mathcal{K}$ from (6.7) in terms of $\tilde{K}$ as

$$(\mathcal{K}f)(R_1) = \int_{SO(3)} \tilde{K}(R_2^{-1}R_1)f(R_2)\,dR_2 = \int_{SO(3)} \tilde{K}(U)f(R_1U^{-1})\,dU, \qquad (6.11)$$

where we used the change of variables $U = R_2^{-1}R_1$. Equation (6.11) implies that $\mathcal{K}$ is a convolution operator over $SO(3)$ given by [22, page 158]

$$\mathcal{K}f = \tilde{K} * f. \qquad (6.12)$$

Similar to the convolution theorem for functions over the real line, the Fourier transform of a convolution over $SO(3)$ is the product of their Fourier transforms, where the Fourier transform is defined by a complete system of irreducible matrix valued representations of $SO(3)$ (see, e.g., [22, Theorem (4.14), page 159]).

Let $\rho_\theta \in SO(3)$ be a rotation by the angle $\theta$ around the $z$ axis, and $\tilde{\rho}_\theta \in SO(2)$ be a planar rotation by the same angle

$$\rho_\theta = \begin{pmatrix} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \tilde{\rho}_\theta = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}.$$

The kernel $\tilde{K}$ satisfies the invariance property

$$\tilde{K}((\rho_\theta U \rho_\alpha)^{-1}) = \tilde{\rho}_\theta \tilde{K}(U^{-1})\tilde{\rho}_\alpha, \quad \text{for all } \theta, \alpha \in [0, 2\pi). \qquad (6.13)$$

To that end, we first observe that $\rho_\theta I^3 = I^3$ and $(U\rho_\alpha)^3 = U^3$ so

$$I^3 \times (\rho_\theta U \rho_\alpha)^3 = (\rho_\theta I^3) \times (\rho_\theta U \rho_\alpha)^3 = \rho_\theta(I^3 \times (U\rho_\alpha)^3) = \rho_\theta(I^3 \times U^3), \qquad (6.14)$$

from which it follows that

$$\|I^3 \times (\rho_\theta U \rho_\alpha)^3\| = \|\rho_\theta(I^3 \times U^3)\| = \|I^3 \times U^3\|, \qquad (6.15)$$

because $\rho_\theta$ preserves length, and it also follows that

$$(I^3 \times (\rho_\theta U \rho_\alpha)^3)(I^3 \times (\rho_\theta U \rho_\alpha)^3)^T = \rho_\theta(I^3 \times U^3)(I^3 \times U^3)^T \rho_\theta^{-1}. \qquad (6.16)$$

Combining (6.15) and (6.16) yields

$$\frac{(I^3 \times (\rho_\theta U \rho_\alpha)^3)(I^3 \times (\rho_\theta U \rho_\alpha)^3)^T}{\|I^3 \times (\rho_\theta U \rho_\alpha)^3\|^2}\rho_\theta U \rho_\alpha = \rho_\theta \frac{(I^3 \times U^3)(I^3 \times U^3)^T}{\|I^3 \times U^3\|^2}U\rho_\alpha, \qquad (6.17)$$

which together with the definition of $\tilde{K}$ in (6.3) and (6.10) demonstrate the invariance property (6.13).

The fact that $\mathcal{K}$ is a convolution satisfying the invariance property (6.13) implies that the eigenfunctions of $\mathcal{K}$ are related to the spherical harmonics $Y_{lm}$ ($l = 0, 1, 2, \ldots$, $m = -l, \ldots, l$) whose eigenvalues have multiplicities $2l + 1$. This relation, as well as the exact computation of the eigenvalues will be established in a separate publication [23]. We note that the spectrum of $\mathcal{K}$ would have been much easier to compute if the normalization factor $\|I^3 \times U^3\|^2$ did not appear in the kernel function $\tilde{K}$ of (6.10). Indeed, in such a case, $\tilde{K}$ would have been a third-order polynomial, and all eigenvalues corresponding to higher order representations must have vanished.

We note that (6.6) implies that the top eigenvalue of $S^{\text{clean}}$, denoted $\lambda_1(S^{\text{clean}})$ scales linearly with $N$, that is, with high probability,

$$\lambda_1(S^{\text{clean}}) = N\lambda_1(\mathcal{K}) + O(\sqrt{N}), \tag{6.18}$$

where the $O(\sqrt{N})$ term is the standard deviation of the sum in (6.5). Moreover, from the top eigenvalues observed in Figures 5.1(a), 5.1(e), 5.3(a), 5.4(a), and 5.5(a) corresponding to $p = 1$ and $p$ values close to 1, it is safe to speculate that

$$\lambda_1(\mathcal{K}) = \frac{1}{2}, \tag{6.19}$$

as the top eigenvalues are approximately 50, 250 and 500 for $N = 100, 500$, and 1000, respectively.

We calculate $\lambda_1(\mathcal{K})$ analytically by showing that the three columns of

$$f(U) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} U^{-1} \tag{6.20}$$

are eigenfunctions of $\mathcal{K}$. Plugging (6.20) in (6.11) and employing (6.10) give

$$(\mathcal{K}f)(R) = \int_{\text{SO}(3)} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \frac{(I^3 \times U^3)(I^3 \times U^3)^T}{\|I^3 \times U^3\|^2} U \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} U^{-1} R^{-1} \, dU. \tag{6.21}$$

From $UU^{-1} = I$ it follows that

$$U \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} U^{-1} = UIU^{-1} - U \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} U^{-1} = I - U^3 U^{3T}. \tag{6.22}$$

Combining (6.22) with the fact that $(I^3 \times U^3)^T U^3 = 0$, we obtain

$$\frac{(I^3 \times U^3)(I^3 \times U^3)^T}{\|I^3 \times U^3\|^2} U \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} U^{-1} = \frac{(I^3 \times U^3)(I^3 \times U^3)^T}{\|I^3 \times U^3\|^2}. \tag{6.23}$$

Letting $U^3 = (x \ y \ z)^T$, the cross product $I^3 \times U^3$ is given by

$$I^3 \times U^3 = (-y \ x \ 0)^T, \tag{6.24}$$

whose squared norm is

$$\|I^3 \times U^3\|^2 = x^2 + y^2 = 1 - z^2, \tag{6.25}$$

and

$$(I^3 \times U^3)(I^3 \times U^3)^T = \begin{pmatrix} y^2 & -xy & 0 \\ -xy & x^2 & 0 \\ 0 & 0 & 0 \end{pmatrix}. \tag{6.26}$$

It follows from (6.21) and identities (6.23)-(6.26) that

$$(\mathcal{K}f)(R) = \int_{\text{SO}(3)} \frac{1}{1 - z^2} \begin{pmatrix} y^2 & -xy & 0 \\ -xy & x^2 & 0 \end{pmatrix} dU R^{-1}. \tag{6.27}$$

The integrand in (6.27) is only a function of the axis of rotation $U^3$. The integral over SO(3) therefore collapses to an integral over the unit sphere $S^2$ with the uniform measure $d\mu$ (satisfying $\int_{S^2} d\mu = 1$) given by

$$(\mathcal{K}f)(R) = \int_{S^2} \frac{1}{1-z^2} \begin{pmatrix} y^2 & -xy & 0 \\ -xy & x^2 & 0 \end{pmatrix} d\mu R^{-1}. \tag{6.28}$$

From symmetry it follows that $\int_{S^2} \frac{xy}{1-z^2} d\mu = 0$, and that $\int_{S^2} \frac{x^2}{1-z^2} d\mu = \int_{S^2} \frac{y^2}{1-z^2} d\mu$. As $\frac{x^2}{1-z^2} + \frac{y^2}{1-z^2} = 1$ on the sphere, we conclude that $\int_{S^2} \frac{x^2}{1-z^2} d\mu = \int_{S^2} \frac{y^2}{1-z^2} d\mu = \frac{1}{2}$, and

$$(\mathcal{K}f)(R) = \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} R^{-1} = \frac{1}{2} f(R). \tag{6.29}$$

This shows that the three functions defined by (6.20) are three eigenfunctions of $\mathcal{K}$ with the corresponding eigenvalue $\lambda_1(\mathcal{K}) = \frac{1}{2}$, as was speculated before in (6.19) based on the numerical evidence.

The remaining spectrum is analyzed in [23], where it is shown that the eigenvalues of $\mathcal{K}$ are

$$\lambda_l(\mathcal{K}) = \frac{(-1)^{l+1}}{l(l+1)}, \tag{6.30}$$

with multiplicities $2l+1$ for $l = 1, 2, 3, \dots$. An explicit expression for all eigenfunctions is also given in [23]. In particular, the spectral gap between the top eigenvalue $\lambda_1(\mathcal{K}) = \frac{1}{2}$ and next largest eigenvalue $\lambda_3(\mathcal{K}) = \frac{1}{12}$ is

$$\Delta(\mathcal{K}) = \lambda_1(\mathcal{K}) - \lambda_3(\mathcal{K}) = \frac{5}{12}. \tag{6.31}$$

**7. Wigner's semi-circle law and the threshold probability.** As indicated by the numerical experiments of Section 5, false detections of common-lines due to noise lead to the emergence of what seems to be Wigner's semi-circle for the distribution of the eigenvalues of $S$. In this section we provide a simple mathematical explanation for this phenomenon.

Consider the simplified model of Section 5.1 that assumes that every common line is detected correctly with probability $p$, independently of all other common-lines, and with probability $1 - p$ the common lines are falsely detected and are uniformly distributed over the unit circle. The expected value of the noisy matrix $S$, whose entries are correct with probability $p$, is given by

$$\mathbb{E}S = pS^{\text{clean}}, \tag{7.1}$$

because the contribution of the falsely detected common lines to the expected value vanishes by the assumption that their directions are distributed uniformly on the unit circle. From (7.1) it follows that $S$ can be decomposed as

$$S = pS_{\text{clean}} + W, \tag{7.2}$$

where $W$ is a $2N \times 2N$ zero-mean random matrix whose entries are given by

$$W_{ij} = \begin{cases} (1-p)S_{ij}^{\text{clean}} & \text{with probability } p, \\ -pS_{ij}^{\text{clean}} + X_{ij}X_{ji} & \text{w.p. } 1-p, \end{cases} \tag{7.3}$$

19

where $X_{ij}$ and $X_{ji}$ are two independent random variables obtained by projecting two independent random vectors uniformly distributed on the unit circle onto the $x$-axis. For small values of $p$, the variance of $W_{ij}$ is dominated by the variance of the term $X_{ij}X_{ji}$. Symmetry implies that $\mathbb{E}X_{ij}^2 = \mathbb{E}X_{ji}^2 = \frac{1}{2}$, from which we have that

$$\mathbb{E}W_{ij}^2 = \mathbb{E}X_{ij}^2 X_{ji}^2 + O(p) = \frac{1}{4} + O(p). \tag{7.4}$$

The eigenvalues of $W$ are therefore distributed according to Wigner's semi-circle law whose support, up to small $O(p)$ terms, is $[-\sqrt{2N}, \sqrt{2N}]$. In other words, the spectral norm, which is also the top eigenvalue of $W$, denoted $\lambda_1(W)$, is given by

$$\lambda_1(W) = \sqrt{2N(1 + O(p))}. \tag{7.5}$$

This prediction is in full agreement with the numerically observed supports in Figure 5.1 and in Figures 5.3-5.5, noting that for $N = 100$ the right edge of the support is located near $\sqrt{200} = 14.14\ldots$, for $N = 500$ near $\sqrt{1000} = 31.62\ldots$, and for $N = 1000$ near $\sqrt{2000} = 44.72\ldots$. The agreement is striking especially for Figures 5.3-5.5 that were obtained from simulated noisy projections without imposing the artificial probabilistic model of Section 5.1 that was used here to actually derive (7.5).

The threshold probability $p_c$ depends on the spectral gap of $S^{\mathrm{clean}}$, denoted $\Delta(S^{\mathrm{clean}})$, and on the top eigenvalue $\lambda_1(W)$ of $W$. From (6.31) it follows that

$$\Delta(S^{\mathrm{clean}}) = \frac{5}{12}N + O(\sqrt{N}). \tag{7.6}$$

In [24, 25, 26] it is proven that the top eigenvalue of the matrix $A + W$, composed of a rank-one matrix $A$ and a random matrix $W$, will be pushed away from the semi-circle with high probability if the condition

$$\lambda_1(A) > \frac{1}{2}\lambda_1(W) \tag{7.7}$$

is satisfied. Clearly, for matrices $A$ that are not necessarily of rank one, the condition (7.7) can be replaced by

$$\Delta(A) > \frac{1}{2}\lambda_1(W), \tag{7.8}$$

where $\Delta(A)$ is the spectral gap. Therefore, the condition

$$p\Delta(S^{\mathrm{clean}}) > \frac{1}{2}\lambda_1(W) \tag{7.9}$$

guarantees that the top three eigenvalues of $S$ will reside away from the semi-circle. Substituting (7.5) and (7.6) in (7.9) results in

$$p\left(\frac{5}{12}N + O(\sqrt{N})\right) > \frac{1}{2}\sqrt{2N(1 + O(p))}, \tag{7.10}$$

from which it follows that the threshold probability $p_c$ is given by

$$p_c = \frac{6\sqrt{2}}{5\sqrt{N}} + O(N^{-1}). \tag{7.11}$$

For example, the threshold probabilities predicted for $N = 100$, $N = 500$, and $N = 1000$, are $p_c \approx 0.17$, $p_c \approx 0.076$, and $p_c \approx 0.054$, respectively. These values are in correspondence with the numerical results of Section 5.1. As for the numerical experiments with the noisy projections presented in subsection 5.2, the prediction (7.11) is in the right ballpark, though it seems to be a bit too pessimistic, as we observe spectral gaps for $p$ values below this threshold. This is perhaps a result of the fact that the falsely detected common lines are not random when considering noisy images.

**8. Summary and discussion.** In this paper we presented efficient methods for computing the rotations of all cryo-EM particles from common lines information in a globally consistent way. Our algorithms, one based on a spectral method (computation of eigenvectors) and the other based on semidefinite programming (a version of max-cut) are able to find the correct set of rotations even at very low common lines detection rates. Using random matrix theory and harmonic analysis on SO(3) we showed that rotations obtained by the eigenvector method can lead to a meaningful ab-initio model as long as the proportion of correctly detected common lines exceeds $\frac{6\sqrt{2}}{5\sqrt{N}}$. It remains to be seen how these algorithms will perform on real raw projection images or on their class averages, and to compare their performance to the recently proposed voting recovery algorithm [12], whose usefulness has already been demonstrated on real datasets.

We note that the techniques and analysis applied here to solve the cryo-EM problem can be translated to the computer vision problem of structure from orthographic images at unknown directions, where lines perpendicular to the epipolar lines play the role of the common-lines. This particular application will be the subject of a separate publication [27].

REFERENCES

[1] Frank, J. (2006) *Three-Dimensional Electron Microscopy of Macromolecular Assemblies: Visualization of Biological Molecules in Their Native State*, Oxford.

[2] Wang, L. and Sigworth, F. J. (2006) Cryo-EM and single particles, *Physiology* (Bethesda). **21**:13-8. Review. PMID: 16443818 [PubMed - indexed for MEDLINE]

[3] Henderson, R. (2004) Realizing the potential of electron cryo-microscopy. *Q Rev Biophys.* **37**(1):3-13. Review. PMID: 17390603 [PubMed - indexed for MEDLINE]

[4] Ludtke, S. J., Baker, M. L., Chen, D. H., Song, J. L., Chuang, D. T., and Chiu, W. (2008) De novo backbone trace of GroEL from single particle electron cryomicroscopy. *Structure*, **16** (3):441–448.

[5] Zhang, X., Settembre, E., Xu, C., Dormitzer, P. R., Bellamy, R., Harrison, S. C., and Grigorieff, N. (2008) Near-atomic resolution using electron cryomicroscopy and single-particle reconstruction. *Proceedings of the National Academy of Sciences*, **105** (6):1867-1872

[6] Chiu, W., Baker M. L., Jiang W., Dougherty M. and Schmid M .F. (2005) Electron cryomicroscopy of biological machines at subnanometer resolution. *Structure.* **13**(3):363-72. Review. PMID: 15766537 [PubMed - indexed for MEDLINE]

[7] Van Heel, M. (1987) Angular reconstitution: a posteriori assignment of projection directions for 3D reconstruction. *Ultramicroscopy.* **21**(2):111-23. PMID: 12425301 [PubMed - indexed for MEDLINE]

[8] Vainshtein, B. and Goncharov, A. (1986) Determination of the spatial orientation of arbitrarily arranged identical particles of an unknown structure from their projections. *Proc. llth Intern. Congr. on Elec. Mirco.*, 459-460.

[9] Farrow, M. and Ottensmeyer, P. (1992) A Posteriori Determination Of Relative Projection Directions Of Arbitrarily Oriented Macrmolecules. *JOSA-A*, **9** (10):1749–1760.

[10] Penczek, P. A., Zhu, J. and Frank, J. (1996) A common-lines based method for determining orientations for $N > 3$ particle projections simultaneously. *Ultramicroscopy* **63**, pp. 205-218.

[11] Mallick, S. P., Agarwal, S., Kriegman, D. J., Belongie, S. J., Carragher, B., and Potter, C. S. (2006) Structure and View Estimation for Tomographic Reconstruction: A Bayesian Approach. *Computer Vision and Pattern Recongnition (CVPR)*, Volume II, 2253–2260.

[12] Singer, A., Coifman, R. R., Sigworth, F. J., Chester, D. W., and Shkolnisky, Y. (2009) Detecting Consistent Common Lines in Cryo-EM by Voting, *submitted*.

[13] Vandenberghe, L., and Boyd, S. (1996) Semidefinite programming. *SIAM Review*, **38**(1):49–95.

[14] Goemans, M. X. and Williamson, D. P. (1995) Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM (JACM)* **42** (6), pp. 1115–1145.

[15] Natterer, F. (2001) *The Mathematics of Computerized Tomography*, SIAM: Society for Industrial and Applied Mathematics, Classics in Applied Mathematics.

[16] Arun, K., Huang, T., and Bolstein, S. (1987). Least-Squares Fitting of Two 3-D Point Sets. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **9** (5), pp. 698–700.

[17] Wigner, E. P. (1955) Characteristic vectors of bordered matrices with infinite dimensions, *Annals of Mathematics* **62** pp. 548–564.

[18] Wigner, E. P. (1958) On the distribution of tile roots of certain symmetric matrices, *Annals of Mathematics* **67** pp. 325–328

[19] Toh, K.C., Todd, M.J., and Tutuncu, R.H. (1999) SDPT3 — a Matlab software package for semidefinite programming, *Optimization Methods and Software,* **11**, pp. 545–581.

[20] Tutuncu, R.H., Toh, K.C., and Todd, M.J. (2003) Solving semidefinite-quadratic-linear programs using SDPT3. *Mathematical Programming Ser. B,* **95**, pp. 189–217.

[21] Keshavan, R. H., Montanari, A., and Oh, S. (2009) Matrix Completion from a Few Entries. *submitted*.

[22] Coifman, R. R., and Weiss, G. (1968) Representations of compact groups and spherical harmonics. *L'Enseignement Mathématique,* **14**, pp. 121–173.

[23] Hadani, R., and Singer, A. (2009) *in preparation*.

[24] Péché, S. (2006) The largest eigenvalues of small rank perturbations of Hermitian random matrices, *Prob. Theo. Rel. Fields* **134** (1): 127–174 .

[25] Féral, D. and Péché, S. (2007) The Largest Eigenvalue of Rank One Deformation of Large Wigner Matrices, *Communications in Mathematical Physics* **272** (1): 185–228.

[26] Füredi, Z., and Komlós, J. (1981) The eigenvalues of random symmetric matrices. *Combinatorica*, **1**, pp. 233–241.

[27] Basri, R., and Singer, A. (2009) *in preparation*.