

# Structural Variability from Noisy Tomographic Projections\*

Joakim Andén<sup>†</sup> and Amit Singer<sup>‡</sup>

**Abstract.** In cryo-electron microscopy, the 3D electric potentials of an ensemble of molecules are projected along arbitrary viewing directions to yield noisy 2D images. The volume maps representing these potentials typically exhibit a great deal of structural variability, which is described by their 3D covariance matrix. Typically, this covariance matrix is approximately low-rank and can be used to cluster the volumes or estimate the intrinsic geometry of the conformation space. We formulate the estimation of this covariance matrix as a linear inverse problem, yielding a consistent least-squares estimator. For  $n$  images of size  $N$ -by- $N$  pixels, we propose an algorithm for calculating this covariance estimator with computational complexity  $\mathcal{O}(nN^4 + \sqrt{\kappa}N^6 \log N)$ , where the condition number  $\kappa$  is empirically in the range 10–200. Its efficiency relies on the observation that the normal equations are equivalent to a deconvolution problem in 6D. This is then solved by the conjugate gradient method with an appropriate circulant preconditioner. The result is the first computationally efficient algorithm for consistent estimation of the 3D covariance from noisy projections. It also compares favorably in runtime with respect to previously proposed non-consistent estimators. Motivated by the recent success of eigenvalue shrinkage procedures for high-dimensional covariance matrix estimation, we incorporate a shrinkage procedure that improves accuracy at lower signal-to-noise ratios. We evaluate our methods on simulated datasets and achieve classification results comparable to state-of-the-art methods in shorter running time. We also present results on clustering volumes in an experimental dataset, illustrating the power of the proposed algorithm for practical determination of structural variability.

**Key words.** cryo-electron microscopy, heterogeneity, single-particle reconstruction, principal component analysis, deconvolution, Toeplitz matrices, shift invariance, conjugate gradient

**AMS subject classifications.** 92C55, 68U10, 44A12, 65R32, 62G05, 62H30, 62J10, 62J07

**1. Introduction.** A single biological macromolecule often exists in a variety of three-dimensional configurations. These can be due to deformations of the molecular structure, known as conformational variability, or smaller molecules being added or removed, known as compositional variability. Since molecular structure dictates biological function, properly resolving these different configurations is of great importance in structural biology. In some cases, it is possible to isolate the different structures experimentally and subsequently image each separately in order to reconstruct its three-dimensional structure. However, this is often not possible due to the similarity in shape and size of the various configurations. In this case, the particles are imaged in a single heterogeneous sample, and their structural variability must be taken into account at the reconstruction stage.

Traditional methods such as X-ray crystallography and nuclear magnetic resonance (NMR)

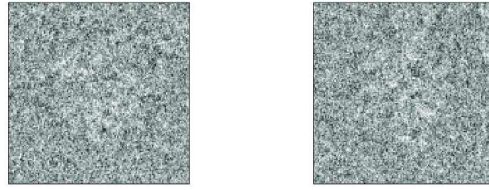
---

\*Submitted to the editors February 6th, 2018.

**Funding:** The authors were partially supported by Award Number R01GM090200 from the NIGMS, Simons Investigator Award, Simons Collaboration on Algorithms and Geometry from Simons Foundation, and the Moore Foundation Data-Driven Discovery Investigator Award.

<sup>†</sup>Center for Computational Biology, Flatiron Institute, New York, NY ([janden@flatironinstitute.org](mailto:janden@flatironinstitute.org)).

<sup>‡</sup>Department of Mathematics and Program in Applied and Computational Mathematics, Princeton University, NJ ([amits@math.princeton.edu](mailto:amits@math.princeton.edu)).



**Figure 1.** Two sample cryo-EM images from a 10000-image dataset depicting the 70S ribosome complex in *E. Coli* [46]. Each image measures 130-by-130 with a pixel size of 2.82 Å. The images depict two similar molecular structures projected in approximately the same viewing direction, but the high noise level makes it difficult to distinguish the difference in structure.

imaging are not well suited for this task, since both rely on aggregate measurements from the whole sample. In contrast, single-particle cryo-electron microscopy (cryo-EM) image each particle separately, and can thus potentially recover the structural variability in the sample. Unlike X-ray crystallography, cryo-EM does not require the crystallization of the sample and can handle larger molecules compared to NMR (as small as 64 kDa [39]), making it a more flexible method. New sample preparation techniques and better detectors have recently yielded reconstructions at near-atomic resolution and the method’s popularity has been steadily on the rise [42, 2, 48]. The 2017 Nobel Prize in Chemistry was awarded to three pioneers of cryo-EM and the technique was named Method of the Year in 2015 by Nature Methods.

To image a set of particles using single-particle cryo-EM, the sample is frozen in a thin layer of ice and exposed to an electron beam. The transmitted electrons are then recorded, forming a set of noisy projection images, one for each particle. Images are modeled as the integral of the particle’s electric potential along a particular viewing direction, followed by convolution with a point spread function and addition of noise [24]. Damaging ionization effects limit the allowable electron dose, so images are dominated by noise (see Figure 1).

The 3D reconstruction task in single-particle cryo-EM assumes that all particles have identical structure and attempts to reconstruct that structure. As mentioned above, however, this is not always the case. The task of reconstructing the structural variability of a heterogeneous population is known as the heterogeneity problem. This variability is typically assumed to be discrete or continuous. In the case of discrete variability, each particle takes on a finite number of possible molecular configurations. This is often referred to as the 3D classification problem in single-particle cryo-EM. For continuous variability, molecular structures vary continuously, forming a smooth manifold on which each point corresponds to a distinct configuration.

The single-particle reconstruction problem in cryo-EM has been approached from many directions, bringing together ideas from statistics and tomography [24, 12, 54, 83, 6]. The popular RELION software implements a regularized maximum-likelihood estimator using expectation-maximization [69]. For the 3D classification problem, it fits a parametric model of discrete variability, so the number of molecular structures needs to be specified in advance. Another problem is that a high-quality initialization is often required for successful reconstruction and there is no global convergence guarantee. In particular, the expectation-maximization algorithm suffers when the number of conformations is large since more populated components take over smaller ones. Lastly, such algorithms require a significant amount of computation

time, although this has recently been reduced using graphics processing units (GPUs) [40, 64].

Another approach has focused on the covariance of the 3D volumes as represented on an  $N$ -by- $N$ -by- $N$  voxel grid. Initial work by Liu & Frank [49] introduced the idea of estimating the variance of these vectors for the purpose of validating the accuracy of reconstructions. In addition, the authors discuss the possibility of quantifying conformational variability through variance estimates. Building on this, a bootstrap method for estimating the 3D variance was introduced by Penczek [60] which was later refined with applications to experimental data [63, 62]. In these methods, a single dataset is resampled multiple times, each yielding a reconstruction. The variance is then calculated from the set of reconstructions.

These works touch on estimation of the entire covariance matrix, but this approach was not fully explored until later by Penczek et al. [59]. Here, the bootstrap method is used to estimate the whole covariance of the volume vectors, known as the 3D covariance matrix. Typically, the covariance is approximately low-rank, since addition or removal of a substructure is captured by a single volume vector, while deformations are often limited in spatial extent and therefore well approximated by a small number of vectors. The top eigenvectors, or “eigenvolumes,” of the 3D covariance thus describe the dominant modes of variability in the volumes. Projecting these eigenvolumes in the viewing direction of each image and calculating the least-squares fit yields a set of coordinates for that image. Fitting a small number of coordinates significantly reduces the noise compared to the original images. Using the coordinates, the images are then clustered and each cluster is used to reconstruct a volume using standard tomographic inversion techniques. Another advantage is that the number of clusters need not be known in advance. For  $C$  volume states, the 3D covariance has rank at most  $C - 1$ , so one plus the number of dominant eigenvalues bounds the number of clusters.

Unfortunately, the heuristic bootstrap estimator used by Penczek et al. does not come with consistency guarantees. To remedy this, an alternate approach was proposed by Katsevich et al. [38], where the 3D covariance estimation problem is formulated as a linear inverse problem and a least-squares estimator is derived. While this estimator is consistent, its calculation involves solving a large-scale linear system, which is prohibitively expensive to invert directly for typical problem sizes. The authors therefore propose a block-diagonal approximation to the linear system in the large-sample limit which can be solved efficiently, but this is only valid for a uniform distribution of viewing angles and a fixed microscope point spread function. It is therefore of limited applicability in experimental datasets. A new approach was proposed by Andén et al. [3], where the exact linear system is solved using the conjugate gradient (CG) method [32]. As a result, the method is valid for non-uniform distributions of viewing directions and multiple point spread functions. However, it only converges after many iterations due to ill-conditioning and each iteration requires a separate pass through the entire dataset, resulting in long running times.

In this paper, we propose an improved version of the method of Andén et al. [3] for efficient and accurate estimation of the 3D covariance matrix. Our method exploits the fact that projection followed by its dual (backprojection) is a convolution operator [21], also known as a Toeplitz operator. This has already resulted in efficient reconstruction techniques in MRI [87, 23, 29] and cryo-EM [84, 88]. The 3D covariance least-squares estimator has a similar structure, letting us pose it as a deconvolution problem in six dimensions. As a result, only one pass through the dataset is required to calculate the convolution kernel, allowing each CG

iteration to be computed quickly. To reduce the number of iterations required for convergence, we employ a circulant preconditioner [81] to improve the conditioning of the system. Our method makes mild assumptions on the distribution of viewing angles and handles image-dependent point spread functions, providing a flexible method for covariance estimation. It is a consistent estimator of the 3D covariance, but unlike the methods of Katsevich et al. [38] and Andén et al. [3], it can be applied efficiently to a wide range of data.

The proposed algorithm has computational complexity  $\mathcal{O}(nN^4 + \sqrt{\kappa}N^6 \log N)$  for  $n$  images, where  $\kappa$  is the condition number of the preconditioned convolution operator and is typically in the range 1–200. This outperforms the algorithm of Katsevich et al. [38], which has computational complexity of  $\mathcal{O}(nN^6 + N^{9.5})$  [38]. It similarly outperforms the method of Andén et al. [3], which has a complexity of  $\mathcal{O}(\sqrt{\kappa'}nN^7)$ , whose condition number  $\kappa'$  is of the order of 5000. The computational complexity is also lower compared to the covariance matrix estimation method introduced by Liao et al. [47], which uses a block Kaczmarz method. Although the paper does not provide an explicit computational complexity of the algorithm, its complexity is at least  $\mathcal{O}(TN^{10})$ , where  $T$  is the number of iterations (typically around 20).

We also introduce a modified covariance estimator based on eigenvalue shrinkage, which lowers the estimation error in the high-dimensional regime. This technique is based on prior work for high-dimensional covariance estimation, where eigenvalue shrinkage methods have been shown to consistently outperform other approaches [19, 25, 18]. These ideas are most relevant when dealing with vectors of dimensionality comparable to the number of samples, which is often the case for cryo-EM. As a result, we can accurately estimate the 3D covariance at lower signal-to-noise ratios than is possible for the conventional least-squares estimator.

We evaluate the proposed algorithms on simulated datasets, showing their ability to deal with high noise levels, non-uniform distribution of viewing angles, and optical aberrations. In particular, we find that the number of images  $n$  necessary to obtain a given covariance estimation accuracy scales inversely with the square of the signal-to-noise ratio, with a phase transition occurring at a critical noise level for a fixed  $n$ . We also compare our algorithm to the state-of-the-art RELION software [69], obtaining superior accuracy in shorter computation time without using an initial reference structure. The eigenvalue shrinkage variant outperforms the standard estimator, achieving the same accuracy for signal-to-noise ratios up to a factor of 1.4 worse. Finally, we evaluate the algorithms on several experimental datasets, where we obtain state-of-the-art reconstruction results. GNU Octave/MATLAB code to reproduce the experiments and figures in this paper are provided by functions located in the heterogeneity folder of the ASPIRE toolbox available at <http://spr.math.princeton.edu/>.

The remainder of the paper is organized as follows. Section 2 presents the heterogeneity problem in cryo-EM, while some background and existing approaches are described in Section 3. Section 4 describes the least-squares estimators for the volume mean and covariance and proposes an eigenvalue shrinkage estimator to improve accuracy in the high-dimensional regime. The proposed algorithms for efficient calculation of these estimators are presented in Section 5. Once the mean and covariance have been estimated, we describe their use for image clustering in Section 6, while simulation and experimental results are provided in Sections 7 and 8, respectively. Possible directions for future work are discussed in Section 9.

**2. Image formation model with heterogeneity.** To model the cryo-EM imaging process, we equate molecular structure with its electric potentials in three dimensions. These potential maps, referred to as volumes, exist in a variety of states, characterizing the structural variability of the molecule. We consider the volume function  $\mathcal{X} : \mathbb{R}^3 \rightarrow \mathbb{R}$  to be a random field of unknown distribution such that  $\mathcal{X} \in L^1(\mathbb{R}^3)$ . In other words,  $\mathcal{X}$  is a random variable in the form of a function. Each realization corresponds to a particular structural configuration of the molecule. The distribution can be discrete, where each realization of  $\mathcal{X}$  takes on one of some finite number  $C$  of states, each of which is a function in  $L^1(\mathbb{R}^3)$ . Alternatively, the structures exist along some continuum, in which case the distribution is continuous, with each realization of  $\mathcal{X}$  being some function in  $L^1(\mathbb{R}^3)$  randomly selected from this continuum.

The electron microscope sends a stream of electrons through the particle represented by  $\mathcal{X}$ , which scatters the electrons. The result is a distorted tomographic projection of each volume, which can be modeled in the weak-phase approximation by an integral along a certain viewing angle followed by a convolution of the resulting image with a microscope-dependent point spread function [24]. The freezing process fixes each particle in a different orientation. We denote the rotation, or viewing direction, of the particle with respect to some reference frame by the 3-by-3 rotation matrix  $R$ , which we assume is drawn from some distribution over the rotation group  $\text{SO}(3)$ . We then define the projection of  $\mathcal{X}$  along  $R$  to be the image

$$(1) \quad \mathcal{Z}(\mathbf{u}) := \int_{\mathbb{R}} \mathcal{X}(R^T[\mathbf{u}; z]) dz,$$

where  $\mathbf{u} \in \mathbb{R}^2$  and  $[\mathbf{u}; z] \in \mathbb{R}^3$  is the concatenation of  $\mathbf{u}$  with  $z$ . This mapping is also known as the X-ray transform of  $\mathcal{X}$  along  $R$  [57]. In addition to the tomographic projection, the configuration of the microscope induces a certain amount of optical aberration, which is modeled by a convolution with some point spread function  $h \in L^1(\mathbb{R}^2)$  [86, 22]. Again, we assume that this is drawn from some (typically discrete) distribution over  $L^1(\mathbb{R}^2)$ . The convolution is defined by

$$(2) \quad \mathcal{Y}(\mathbf{u}) := \int_{\mathbb{R}^2} \mathcal{Z}(\mathbf{v} - \mathbf{u}) h(\mathbf{v}) d\mathbf{v},$$

for  $\mathbf{u} \in \mathbb{R}^2$ . Combining both operations, we have the projection mapping  $\mathcal{P} : L^1(\mathbb{R}^3) \rightarrow L^1(\mathbb{R}^2)$

$$(3) \quad \mathcal{P}\mathcal{X}(\mathbf{u}) := \int_{\mathbb{R}^2} \left( \int_{\mathbb{R}} \mathcal{X}(R^T[\mathbf{u} - \mathbf{v}; z]) dz \right) h(\mathbf{v}) d\mathbf{v}.$$

We can now state our forward model for the cryo-EM imaging process, which takes a volume  $\mathcal{X}$  and gives the image

$$(4) \quad \mathcal{Y} = \mathcal{P}\mathcal{X} + \mathcal{E},$$

where  $\mathcal{E} : \mathbb{R}^2 \rightarrow \mathbb{R}$  is a white Gaussian random field of variance  $\sigma^2$ . Since  $\mathcal{P}$ ,  $\mathcal{X}$  and  $\mathcal{E}$  are random variables,  $\mathcal{Y}$  is also a random variable. The noise  $\mathcal{E}$  represents error introduced into the image due to non-interacting electrons, inelastic scattering, and quantum noise [85, 7, 61]. While these error sources follow a Poisson distribution, the counts are typically high enough

for this to be well-approximated by a Gaussian distribution. Another potential problem is that the noise is rarely white but exhibits strong correlations which depend on the microscope configuration. Later in this section, we describe how to account for this discrepancy.

It is useful to consider the above mappings in the Fourier domain. Let us define the  $d$ -dimensional Fourier transform  $\mathcal{F}\mathcal{G}$  of some function  $\mathcal{G} \in L^1(\mathbb{R}^d)$  by

$$(5) \quad \mathcal{F}\mathcal{G}(\boldsymbol{\omega}) := \int_{\mathbb{R}^d} \mathcal{G}(\mathbf{u}) e^{-2\pi i \langle \boldsymbol{\omega}, \mathbf{u} \rangle} d\mathbf{u},$$

for any frequency  $\boldsymbol{\omega} \in \mathbb{R}^d$ . In this case, the tomographic projection mapping  $\mathcal{P}$  satisfies

$$(6) \quad \mathcal{F}\mathcal{P}\mathcal{X}(\boldsymbol{\omega}) = \mathcal{F}\mathcal{X}(R^T[\boldsymbol{\omega}; 0]) \cdot \mathcal{F}h(\boldsymbol{\omega}),$$

for any  $\boldsymbol{\omega} \in \mathbb{R}^2$ , which is known as the Fourier Slice Theorem [57]. In other words, projection in the spatial domain corresponds to restriction (or “slicing”) to a plane in the Fourier domain and multiplication by a transfer function  $\mathcal{F}h$ . The Fourier transform  $\mathcal{F}h$  of the point spread function  $h$  is known as the contrast transfer function (CTF). The CTF is an oscillatory function and is equal to zero for several frequencies. As a result, those frequencies are not available in that particular image. To mitigate this problem, cryo-EM datasets are collected for a number of different CTFs by varying certain microscope parameters.

The model presented above describes continuous images  $\mathcal{Y}$  obtained from continuous volume densities  $\mathcal{X}$ . While an accurate model of the physical process, it is not compatible with the output of an electron microscope, which is in the form of discrete images  $y$  with values on an  $N$ -by- $N$  pixel grid, where  $N$  typically ranges from 100 to 500. The images are therefore limited in resolution, imposing limit on the resolution of the reconstructed volumes.

To discretize the images, we define the  $N$ -point grid  $M_N$  as

$$(7) \quad M_N := \{-\lfloor N/2 \rfloor, \dots, \lfloor N/2 - 1 \rfloor\}.$$

An image  $\mathcal{P}$  is then represented by sampling evenly over the square  $[-1, +1]^2$  at points  $2M_N^2/N$ , yielding a function  $y : M_N^2 \rightarrow \mathbb{R}$ . We treat this function as a vector in  $\mathbb{R}^{N^2}$ .

There is more choice in representing the volumes. One popular approach is to consider the voxel samples of  $\mathcal{X}$  at points  $2M_N^3/N$  in the cube  $[-1, +1]^3$ , yielding a vector of dimension  $N^3$  [68, 84, 88]. While this basis has computational advantages, it is not always well suited to representing volumes of interest. As such, we will use a different basis described in Section 5.5 and convert between this and the voxel basis. Let us denote by  $Q$  the matrix whose columns corresponding to the  $p$  basis vectors of size  $N^3$ , each corresponding to an  $N$ -by- $N$ -by- $N$  volume. Here  $p = \mathcal{O}(N^3)$  and  $Q$  is of size  $N^3$ -by- $p$ . We will assume that expansion and evaluation in this basis is fast, that is both  $Q$  and its transpose  $Q^T$  can be applied in  $\mathcal{O}(N^3 \log N)$ . To simplify expressions, we introduce the notation  $v = Qx$ .

An important question is then how to properly discretize the forward model (4). One approach is embed the discretized volumes  $v \in \mathbb{R}^{N^3}$  into  $L^1(\mathbb{R}^3)$  using a sinc basis, apply  $\mathcal{P}$  and fit the result to a sinc basis expansion in  $L^1(\mathbb{R}^2)$  using least-squares. While this provides a matrix representation of  $\mathcal{P}$  that is accurate in a least-squares sense and converges to  $\mathcal{P}$  as  $N \rightarrow \infty$ , the mapping is computationally inefficient, requiring a full  $N^2$ -by- $N^3$  matrix multiplication to apply.



Another approach is to mimic the Fourier slice structure of  $\mathcal{P}$  described in (6), enabling speedups associated with fast Fourier transforms (FFTs) [14]. First, let us define the discrete Fourier transform of some function  $g : M_N^d \rightarrow \mathbb{R}$  in  $d$  dimensions

$$(8) \quad \mathcal{F}g(\mathbf{k}) := \sum_{\mathbf{i} \in M_N^d} g(\mathbf{i}) e^{-2\pi i \langle \mathbf{i}, \mathbf{k} \rangle / N}.$$

Here, we have abused notation slightly by having  $\mathcal{F}$  signify both the continuous and discrete Fourier transforms. The nature of the mapping should be clear from context. Although  $\mathcal{F}g(\mathbf{k})$  is defined for any frequency vector  $\mathbf{k} \in \mathbb{R}^d$ , it is traditionally restricted to the grid  $M_N^d$ .

We now define the mapping  $I$  transforming the voxel volume  $v$  into the image  $Iv$  through

$$(9) \quad Iv(\mathbf{i}) = \frac{1}{N^3} \sum_{\mathbf{j} \in M_N^3} v(\mathbf{j}) \sum_{\mathbf{k} \in M_{2\lceil N/2 \rceil - 1}^2} \mathcal{F}h(\mathbf{k}) e^{-2\pi i (\langle R^T[\mathbf{k}; 0], \mathbf{j} \rangle - \langle \mathbf{k}, \mathbf{i} \rangle) / N} \quad \text{for all } \mathbf{i} \in M_N^2.$$

Computing its discrete Fourier transform, we obtain

$$(10) \quad \mathcal{F}Iv(\mathbf{k}) = \begin{cases} \frac{1}{N} \mathcal{F}v(R^T[\mathbf{k}; 0]) \cdot \mathcal{F}h(\mathbf{k}), & \mathbf{k} \in M_{2\lceil N/2 \rceil - 1}^2 \\ 0, & \text{otherwise,} \end{cases}$$

for any  $\mathbf{k} \in M_N^2$ . The operator  $I$  therefore satisfies a discrete version of the Fourier Slice Theorem (6). Note that in the case of even  $N$ , the Nyquist frequencies at  $-N/2$  are set to zero to ensure a real image  $IQx$ . Enforcing this one-to-one mapping of frequencies allows us to derive the convolutional formulations in Sections 5.1 and 5.2. The entire mapping, from  $x$  to  $v = Qx$  to  $Iv = IQx$  is denoted by  $P = IQ$  and is called the volume imaging mapping. While  $P$  describes the whole imaging process, that is both projection and convolution with the point spread function, we shall often refer to it as projection for simplicity. Similarly, its adjoint  $P^T$  will be referred to as backprojection.

To project a volume  $x \in \mathbb{R}^{N^3}$ , first evaluate it on the  $2M_N^3/N$  voxel grid to obtain  $v = Qx$ , then calculate the discrete Fourier transform  $\mathcal{F}v$  using (8) on the grid defined by  $R^T[\mathbf{k}; 0]$ . We then multiply the Fourier transform pointwise by the contrast transfer function  $\mathcal{F}h(\mathbf{k})$ , set Nyquist frequencies to zero, and apply the inverse discrete Fourier transform, which gives  $Px \in \mathbb{R}^{N^2}$ . As mentioned earlier, the basis evaluation matrix  $Q$  can be applied in  $\mathcal{O}(N^3 \log N)$  time. The first discrete Fourier transform  $\mathcal{F}v$  is computed using a non-uniform fast Fourier transform (NUFFT), which has a computational complexity of  $\mathcal{O}(N^3 \log N)$  [21, 27], while pointwise multiplication and the 2D inverse FFT require  $\mathcal{O}(N^2)$  and  $\mathcal{O}(N^2 \log N)$ , respectively. The overall computational complexity is therefore  $\mathcal{O}(N^3 \log N)$ , which is a significant improvement over the direct matrix multiplication approach, which has a complexity of  $\mathcal{O}(N^5)$ . Another important advantage is that calculating multiple projections of the same volume  $x$ , the overall complexity scales as  $\mathcal{O}(N^3 \log N + nN^2)$ , where  $n$  is the number of projection images.

With the projection mapping  $P$ , we can now formulate our discrete forward model as

$$(11) \quad \mathbf{y} := P\mathbf{x} + \mathbf{e},$$

where  $\mathbf{e} \in \mathbb{R}^{N^2}$  is a standard white Gaussian noise image. Note that again  $P$  describes both projection and convolution with a point spread function  $h$ . To generate an image, a given

volume density  $x$  is rotated according to the viewing direction  $R$ , projected along the  $z$ -axis, convolved with a point spread function  $h$ , and finally a white Gaussian noise  $e$  is added. We do not include translation as part of  $P$  because we assume that translations have been previously estimated and subsequently removed from the images by translating them in the opposite direction. Apart from this and the white noise assumption, the above image formation model corresponds to those traditionally employed in single-particle cryo-EM [24, 69, 64, 78].

As in the continuous case, we note that  $P$ ,  $x$ , and  $e$  are all random variables, so the same is true of  $y$ . In particular,  $P$  depends on  $R$  and  $h$ , which are random variables with distributions over  $\text{SO}(3)$  and  $L^1(\mathbb{R}^2)$ , respectively (although in many cases we will condition  $y$  with respect to  $P$ , fixing that variable and by extension  $R$  and  $h$ ). The volume density  $x$  is a random variable defined over  $\mathbb{R}^p$  as described above, and finally  $e$  is a random variable with distribution over  $\mathbb{R}^{N^2}$ . The distribution of  $y$  depends on those of  $R$ ,  $h$ ,  $x$  and  $e$  and drawing realizations from this distribution provides us with the experimental images. It is not necessary for us to know what these distributions are, but this probabilistic framework will make it easier to derive and reason about the estimators introduced in the following.

In a single-particle cryo-EM experiment, we have more than one image, with multiple copies of the same molecule being imaged separately, each with a different structural configuration, projected at a different viewing angle, subjected to convolution by a different point spread function, and degraded by a different realization of noise. As a result, we consider identically distributed, independent copies  $x_1, \dots, x_n$  of  $x$ . Similarly, we have rotations  $R_1, \dots, R_n$ , point spread functions  $h_1, \dots, h_n$  and noise vectors  $e_1, \dots, e_n$  which are independent and identically distributed copies of  $R$ ,  $h$ , and  $e$ , respectively. These yield the images

$$(12) \quad y_s := P_s x_s + e_s$$

for  $s = 1, \dots, n$ , where  $P_s$  is the imaging operator corresponding to the viewing direction  $R_s$  and contrast transfer function  $\mathcal{F}h_s$ . Each  $y_s$  is a different random variable representing a different projection image from an experiment. The contrast transfer functions  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  are not all distinct, but are typically shared between particles picked from the same micrograph.

We note that the  $P_s$  mappings are not known for experimental data, but must be estimated. Several methods exist to estimate the CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  [55, 91]. Similarly, traditional methods for orientation estimation [69, 64, 28, 89] can be applied when  $x$  does not vary too much around its mean. Indeed, previous works have demonstrated that this is a feasible approach in a range of situations [59, 47]. In this work, we assume such estimation of  $R_1, \dots, R_n$  is possible and therefore  $P_1, \dots, P_n$  are known to a certain accuracy.<sup>1</sup>

In the process of estimating  $R_1, \dots, R_n$ , these methods also estimate the translations of the individual images. This is what allows us to cancel their effect in our forward model (11). Note that the effect of the viewing directions cannot be similarly removed since they define the tomographic projection from 3D to 2D. As such, it is necessary to include them in the forward model, unlike the translations.

As mentioned above, the noise component  $e$  is typically not white. One approach to handle this is to estimate the power spectrum of the noise and include a non-white noise component

---

<sup>1</sup>In the case of high variability in  $x$ , the orientations must be estimated simultaneously with the clustering of the images, a more challenging problem studied in recent works by Lederman and Singer [44, 45].



in our forward model. Most of the results in the remainder of this work follow, but with a certain loss of simplicity. Instead, we choose to prewhiten the images, rendering the data more compatible with our proposed forward model (11) which specifies white Gaussian noise.

To achieve this, the noise power spectrum is first estimated for each image. Due to changing experimental conditions, the intensities of the various noise sources (background, inelastic scattering, quantum noise, etc.) vary from image to image. To account for this, we employ a noise estimation algorithm which exploits such low-rank variability [4]. Alternatively, the software used to estimate the rotations  $R_1, \dots, R_n$  also estimates noise power spectra as part of the algorithm [68, 64].

Let us denote the estimated noise power spectrum of  $e_s$  to be  $\mathcal{F}m_s$ . To whiten  $y_s$ , we filter it with the transfer function  $\mathcal{F}m_s^{-1/2}$ . The noise component is now white, but the signal component has been altered. This is remedied by including the whitening filter into the CTF  $\mathcal{F}h_s$ , replacing it with a new “effective CTF”  $\mathcal{F}m_s^{-1/2} \cdot \mathcal{F}h_s$ . The resulting images have an approximately white noise component, with a signal still equal to  $P_s x_s$  since  $P_s$  now includes the effect of the prewhitening filter. Consequently, the prewhitening filter is inverted whenever we use the new  $P_1, \dots, P_n$  to reconstruct  $x$ .

The heterogeneity problem can now be stated more formally. From projection images  $y_1, \dots, y_n$ , we would like to characterize the distribution of  $x$ , which is equivalent to reconstructing the energy landscape inhabited by the molecule. This problem is unfortunately ill-posed. Indeed, third- and higher-order moments of  $x$  are impossible to estimate from its projections. The ill-posedness may be removed by only estimating partial information on  $x$ , such as its first- and second-order moments  $\mathbb{E}[x]$  and  $\text{Cov}[x]$ , or by restricting the class of distributions, such as those supported on a low-dimensional subspace. This latter restriction includes discrete variability and continuous variability on a smooth, low-dimensional manifold. In this work we shall make use of both restrictions to render the heterogeneity problem more tractable. First, given images  $y_1, \dots, y_n$ , we will first estimate  $\mathbb{E}[x]$  and  $\text{Cov}[x]$ . Second, we will use the fact that the distribution of  $x$  is supported on a low-dimensional subspace to improve our estimate of  $\text{Cov}[x]$ .

A further goal is to reconstruct the underlying volumes  $x_1, \dots, x_n$  from their projection images  $y_1, \dots, y_n$ . This problem is also ill-posed without any further assumptions. We therefore impose the same restriction on the class of distributions (namely, being supported on a low-dimensional subspace). Estimates of the mean and covariance lets us estimate the volumes themselves. However, this is only possible within limits dictated by the noise level. Indeed, if the noise is large enough, we may be able to estimate the mean and covariance accurately given enough images, but accurate reconstruction of individual volumes may not be feasible.

**3. Related work.** Due to the importance of determining structural variability from cryo-EM projections, much work has been focused on resolving this heterogeneity problem. Although various methods have been introduced that have some degree of experimental success, they do not possess any accuracy guarantees, so it is sometimes difficult to validate the reconstructions. In addition, they also often rely on good initializations, which can significantly bias the final result. Finally, the computational complexity of these methods is typically quite high, requiring a large amount of computational resources.

**3.1. Maximum likelihood.** One popular method for solving the heterogeneity problem has been to set up a probabilistic model for image formation and maximizing the likelihood function with respect to the model parameters given the data. This was first considered for class averaging in the space of images by Sigworth [74], where the probability density of the images  $y$  was modeled as a mixture of Gaussians, with each component center constituting a distinct image class. These centers, along with other parameters such as shifts were estimated by maximizing the likelihood using an expectation-maximization algorithm [17].

The maximum-likelihood method was subsequently extended by Scheres, who modeled the underlying volume vector  $x$  as a mixture of Gaussians and regularized the likelihood function using Bayesian priors on the parameters, which includes the viewing directions  $R_1, \dots, R_n$  [69]. The resulting algorithm, implemented in the RELION software package, has seen significant success and provides generally satisfactory volume estimates [68]. However, since the algorithm attempts to optimize a non-convex function, there is no guarantee that a globally optimal solution is obtained. The algorithm also needs to be initialized with single reference structure that is similar to the molecule being imaged, which can significantly bias the result if not chosen carefully. Similarly, the number of clusters is part of the model and needs to be specified in advance, limiting the method’s flexibility. The performance also degrades when a large number of classes is specified, as more populated classes tend to absorb smaller ones, making rare conformational changes hard to characterize. Finally, the algorithm has a long running time, although a recent GPU-based implementation has mitigated this problem [40].

**3.2. Common lines.** Another approach proposed by Shatsky et al. [71] clusters the projection images by defining a similarity measure between all pairs and applying a spectral clustering method. From the Fourier Slice Theorem (6), the Fourier transform  $\mathcal{F}y$  of a clean projection image  $y$  corresponds to the restriction of the volume Fourier transform  $\mathcal{F}x$  to a plane and multiplied by a contrast transfer function  $\mathcal{F}h$ . Two images  $y_s, y_t$  thus occupy two central planes of the volume Fourier transform and intersect along a common line. The Fourier transforms  $\mathcal{F}y_s, \mathcal{F}y_t$  of two noiseless projections of the same molecular structure should therefore coincide along this common line (up to differing contrast transfer functions), so their cross-correlation along this line provides a good similarity measure.

Using this common-lines affinity, the authors applied spectral clustering to group the images according to their underlying molecular structure. Unfortunately, the Fourier transforms of two projections of the same volume will not coincide exactly due to the image noise. In most cryo-EM experiments, the images are dominated by noise, making this particular approach unfeasible without some amount of denoising. Denoising images in cryo-EM is traditionally achieved by class averaging, where images that represent similar views are averaged together. However, for heterogeneous data this may break down since images belonging to different molecular structures could be averaged together.

**3.3. Covariance and low-rank approaches.** Instead of directly clustering the images, another line of work has focused on estimating the 3D covariance matrix  $\text{Cov}[x]$  of the volume  $x$  as a random vector. The first of these, Liu & Frank [49], introduced the notion of 3D covariance in the single-particle cryo-EM setting for the purposes of validation. The authors then proposed a method for estimating this covariance. Building on this, Penczek outlined a variant of the standard bootstrap algorithm for estimating the variance [60]. Here, multiple

subsets of the images  $y_1, \dots, y_n$  are drawn, each yielding a different 3D reconstruction. The sample covariance of these 3D reconstructions then yields an estimate for the 3D covariance of  $x$ . Since the distribution of molecular structures differs slightly between subsets, the idea is that this will capture the 3D variability of the volumes. Refinements of this method have successfully been applied to experimental data to estimate variance [63] and perform 3D classification [62]. Related work by Doerschuk and others estimate the covariance by fitting Gaussian mixture models of the volumes [94, 90].

A further refinement was proposed in Penczek et al. [59], where the bootstrap method was used to perform principal component analysis of the reconstructed volumes. In this work, the top eigenvectors, or eigenvolumes, of the estimated 3D covariance matrix are used to reconstruct the volumes in the sample. Indeed, as discussed above, this covariance matrix is typically approximated by a low-rank matrix, so its top eigenvolumes together with the mean volume form an affine space containing most of the volumes in the dataset. Projecting the mean volume and the eigenvolumes, the authors find the coordinates of each image in this affine space through a least-squares fit. Clustering the images using these coordinates, each cluster is then used to reconstruct a different molecular structure. Unfortunately, the heuristic nature of the covariance estimation does not provide any accuracy guarantees. A maximum-likelihood approach has also been proposed to estimate the top eigenvolumes [78], but this suffers from the same initialization problems and lack of guarantees as other non-convex approaches (see Section 3.1).

To address these problems, Katsevich et al. [38] formulated the 3D covariance estimation problem as a linear inverse problem and proposed a least-squares estimator with proven consistency results. Unfortunately, direct calculation of the estimator proved computationally intractable, so the authors introduced an approximation which relied on uniform distribution of viewing directions and a single fixed contrast transfer function. These conditions are rarely satisfied in experimental data, so practical use of this method was limited. A more flexible approach has been proposed, based on calculating a related estimator using the block Kaczmarz method, but unfortunately this suffers from slow convergence and reduced accuracy [47].

An improved version of the method of Katsevich et al. [38] was introduced in Andén et al. [3], where the exact linear system was solved iteratively using the CG method. This had the advantage of allowing for non-uniform distributions of viewing directions and varying contrast transfer functions. However, the method required a pass through the dataset at each iteration and a large number of iterations was needed to reach convergence. As a result, this method proved unfeasible for large datasets.

**3.4. Other methods.** Another successful approach has been to use an atomic reference structure to predict the possible motions of a molecule using normal mode analysis. These motions are then fit to the projection images to identify 3D variability in the sample [35]. A significant drawback of this approach is its requirement for an atomic-resolution reference structure, which may not always be available for the imaged molecule.

To capture continuous variability, Dashti et al. [15] group projection images by viewing direction and estimate the manifold structure in each group. The different manifolds are then assembled into a global manifold describing the variability of the entire molecule. Counting the number of projection images obtained from each point in that manifold, the authors derive an

energy landscape for that molecule. The authors apply this method to a dataset of ribosome projections and obtain impressive results. However, the heuristic nature of this method and its lack of accuracy guarantees make it problematic to apply in a general setting.

For a survey methods related to the heterogeneity problem, see Jonić [37].

**4. Mean and covariance estimators.** As discussed in the previous section, the 3D covariance is a powerful tool in characterizing variability for single particle cryo-EM [49]. In particular, applying it to perform a principal component analysis of the volumes is especially useful [59, 78]. Existing methods for covariance estimation, however, do not offer any accuracy guarantees [94, 63, 59, 78]. In the following, we describe the least-squares estimators for both volume mean and covariance previously introduced by Katsevich et al. [38]. The estimators take as input the projection images along with estimates of the viewing directions and CTFs and provide as output estimates of the mean and covariance of the volumes. These estimators have theoretical guarantees, ensuring that for a fixed noise level, the mean and covariance estimates converge to their population values as the number of images increases. We also introduce a modified variant of the original covariance estimator, where eigenvalue shrinkage is used to reduce error in the regime of high dimension.

**4.1. Mean estimator.** To estimate the mean  $\mathbb{E}[x]$  of the volume density  $x$ , we assume that the imaging operators  $P_1, \dots, P_n$  are known for all images (i.e., that the viewing directions  $R_1, \dots, R_n$  and CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  are known) and that the images are centered. As discussed in Section 2, the imaging operators can typically be estimated from the images with good accuracy. Likewise, their translations can be estimated and used to center the images.

We now consider the distribution of the images  $y_1, \dots, y_n$  with these mappings fixed. That is, we consider the distributions of  $y_s$  conditioned on  $P_s$ , where only the volume structure  $x_s$  and noise  $e_s$  are allowed to vary. Let us denote the expectation with respect to some variable  $a$  conditioned on some other variable  $b$  by  $\mathbb{E}_{a|b}$ . From the forward model (11), we obtain

$$(13) \quad \mathbb{E}_{x_s, e_s | P_s}[y_s] = P_s \mathbb{E}[x_s] = P_s \mathbb{E}[x],$$

for  $s = 1, \dots, n$ , since  $x_1, \dots, x_n$  are all identically distributed and  $\mathbb{E}[e_1] = \dots = \mathbb{E}[e_n] = 0$ . The above equation provides a constraint on  $\mathbb{E}[x]$  for each  $s = 1, \dots, n$ . We could therefore solve for  $\mathbb{E}[x]$  if we were given the left-hand side expectations  $\mathbb{E}_{x_s, e_s | P_s}[y_s]$ , but these are unavailable to us.

However, the image itself  $y_s$  is an unbiased estimator of  $\mathbb{E}_{x_s, e_s | P_s}[y_s]$ , albeit one with significant variance. Substituting  $y_s \approx \mathbb{E}_{x_s, e_s | P_s}[y_s]$  into (13) gives

$$(14) \quad y_s \approx P_s \mathbb{E}[x],$$

for all  $s = 1, \dots, n$ . Combining these approximate constraints into a regularized least-squares objective, we obtain the following estimator  $\mu_n$  for  $\mathbb{E}[x]$ :

$$(15) \quad \mu_n := \operatorname{argmin}_{\mu \in \mathbb{R}^p} \frac{1}{n} \sum_{s=1}^n \|y_s - P_s \mu\|^2 + \nu_n \|\mu\|^2,$$

where  $\nu_n \geq 0$  is a regularization parameter which ensures the problem remains well-posed (mitigating ill-posedness due to CTFs and distribution of viewing directions, as discussed

below). Since the problem is better conditioned for large  $n$ ,  $\nu_n$  typically decreases as  $n$  grows. This estimator minimizes the average square distance of the  $P_s\mu_n$  (that is,  $\mu_n$  projected along  $R_s$  and convolved with  $h_s$ ) to each of the images  $y_s$  subject to the regularization term  $\nu_n\|\mu\|^2$ . When there is no structural variability, that is, when  $\text{Cov}[x] = 0$ ,  $\mu_n$  is the regularized maximum-likelihood estimator for  $\mathbb{E}[x]$ .

In order to calculate  $\mu_n$ , we form its normal equations by differentiating the objective (15) and setting the derivative to zero. We thus have

$$(16) \quad A_n\mu_n = b_n,$$

where

$$(17) \quad A_n := \frac{1}{n} \sum_{s=1}^n P_s^T P_s + \nu_n I_p,$$

$$(18) \quad b_n := \frac{1}{n} \sum_{s=1}^n P_s^T y_s.$$

The right-hand side  $b_n$  is the average of the backprojected images  $P_s^T y_s$ , while  $A_n$  is the projection-backprojection operator corresponding to the set of viewing directions  $R_1, \dots, R_n$  and CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  plus the regularization term  $\nu_n I_p$ .

This least-squares estimator is a good estimator of molecular structure in cryo-EM samples when heterogeneity is not present. Indeed, if  $x$  has no variability,  $x_1 = \dots = x_n = \mathbb{E}[x]$  then  $\mu_n$  estimates the single volume present in the sample. When there is heterogeneity, it is a consistent estimate of the average volume  $\mathbb{E}[x]$ . As a result,  $\mu_n$  and closely related least-squares estimators have found widespread use in single-particle cryo-EM reconstruction [30, 84, 88, 6].

The standard Tikhonov regularization term in (15) can be replaced by more sophisticated regularizers to enforce smoothness and other properties. For example, RELION uses an adaptive weighting scheme where each radial frequency is assigned a different regularization parameter [68]. These are initialized using a reference structure and subsequently updated using the estimated structure at each iteration. The result is that higher frequencies are penalized more than lower frequencies, enforcing smoothness.

The constraints in (14) are by necessity loose due to the presence of noise in the data and the variability of  $x$ . Still, aggregating these over a large number of images  $n$  results in an estimator  $\mu_n$  that is consistent. Indeed, Katsevich et al. [38] showed that, for  $\nu_n = 0$ , we have

$$(19) \quad \|\mu_n - \mathbb{E}[x]\| = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right),$$

with high probability as  $n \rightarrow \infty$ , provided that  $x$  is bounded and  $A_n$  is invertible when  $n > n_0$  for some  $n_0$ . Informally, this invertibility condition is satisfied when the Fourier slices densely populate the Fourier domain and the zeros of the CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  do not overlap. We shall make this first requirement more specific below.

To understand the above result, it is helpful to consider a simpler toy example. Instead of a volume, we have a random vector  $x = [a, b, c]$  containing three values. We have observations for every  $s = 1, \dots, n$ , but we only observe two entries of  $x$  for each  $s$ . That is, our observations

$y_s$  are of the form  $[a_s, b_s, ?]$ ,  $[a_s, ?, c_s]$ , and  $[?, b_s, c_s]$ , where  $?$  denotes a missing value. From this data, we can estimate the mean of  $x = [a, b, c]$ .

First, to estimate  $\mathbb{E}[a]$ , we collect all observations in which the first entry is present, that is, observations of the form  $[a_s, b_s, ?]$  and  $[a_s, ?, c_s]$ . We then average all the first entries to obtain our estimate of  $\mathbb{E}[a]$ . We can similarly estimate the means of the other entries. All that is required is that we know which entry is missing in each observation (i.e., the location of the  $?$ ) and that all entries are represented among the observations.

The parallel with mean estimation in cryo-EM is established by considering the Fourier Slice Theorem (6). Let us first consider the case without CTF, that is for  $\mathcal{F}h = 1$ . In this case,  $P$  is pure projection and acts as the restriction to a central slice in the Fourier domain. Its adjoint, the backprojection mapping  $P^T$ , therefore inserts a two-dimensional Fourier transform into that plane, with remaining frequencies set to zero. Just like in the toy example, the Fourier transform  $\mathcal{F}y$  of each observation is a “subset” of the entries of  $\mathcal{F}x$ . If we can place each observed value at its appropriate frequency in the Fourier domain and average across all observations, we could estimate the mean Fourier transform  $\mathbb{E}[\mathcal{F}x]$ . Again, what is required is that we know the projection operators  $P_1, \dots, P_n$  (that is, we know the rotations  $R_1, \dots, R_n$ ) and that the collection of all central slices adequately covers the Fourier domain. The requirement on the viewing directions is not very strict. For example, a set of Fourier slices forming a fan-like pattern (where their normals are contained in a single plane) or a tilt series with no missing wedge are enough to guarantee accurate estimation of the mean.

This is exactly what it done by the least-squares estimator  $\mu_n$  as defined by (16)–(18) but in a more formal way. The right-hand side  $b_n$  takes the Fourier transform  $\mathcal{F}y_s$  of each image, places it onto the proper plane in three-dimensional Fourier space as defined by  $R_s$ , and averages across all images for  $s = 1, \dots, n$ . If the central slices cover enough of the three-dimensional Fourier domain, this will “reconstruct” the average volume  $\mathbb{E}[x]$  up to a frequency-dependent reweighting that depends on the distribution of viewing directions. Indeed, even in the case of uniform distribution of viewing directions  $R$  over  $\text{SO}(3)$ , frequencies close to the origin will be oversampled compared to frequencies farther away, and this must be compensated for. This reweighting is encoded in  $A_n$  through the average of the projection-backprojection operators  $P_s^T P_s$ . Again, (6) tells us that  $P_s^T P_s$  in the Fourier domain is restriction followed by insertion, which is equivalent to multiplication by the indicator function of the plane corresponding to  $R_s$ . Adding all of these indicator functions together yields the reweighting  $A_n$  relating  $\mu_n$  to  $b_n$ . Here, having an adequate coverage of the Fourier domain by the central slices implies that  $A_n$  is invertible.

While  $A_n$  may be invertible, it may still have a high condition number. This can happen, for example, if certain viewing directions are more common than others. In a given direction, higher frequencies are also undersampled with respect to lower frequencies. Reconstruction at a higher resolution  $N$  is therefore less well-conditioned compared to at lower  $N$ . These conditioning problems can be partially remedied by choosing an appropriate regularization parameter  $\nu_n$  at the cost of some bias in the estimation.

We note here that if the viewing directions  $R_1, \dots, R_n$  are sampled from the uniform distribution over  $\text{SO}(3)$ , the inverse of  $A = \lim_{n \rightarrow \infty} A_n$  is a ramp filter  $\|\omega\|$  so inverting it is particularly straightforward. When  $n$  is large, we can therefore calculate  $A^{-1}(b_n)$ , which is known as filtered backprojection [31, 57, 65], a popular estimator for reconstruction in



computerized tomography (CT) and related fields. As  $n \rightarrow \infty$ , this estimate will converge to  $\mathbb{E}[x]$ . However, in cryo-EM, the distribution of viewing directions is typically non-uniform, so this idealized ramp filter is not appropriate and the exact  $A_n$  must be used [30].

In the case of non-trivial CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$ , the same ideas hold, except that back-projection  $P_s^\top$  includes multiplication by the CTF  $\mathcal{F}h_s$  and projection-backprojection  $P_s^\top P_s$  involves multiplication by  $|\mathcal{F}h_s|^2$ . As a result, for each  $s$ , certain frequencies are zeroed out in the corresponding term of  $\sum_{s=1}^n P_s^\top P_s$  due to the CTF since  $y_s$  does not contain any information at those frequencies. However, this is mitigated by the fact that we have different CTFs for different images, and therefore different sets of zeros. As long as these do not all overlap, the matrix  $A_n$  is invertible. In addition, the CTF is small at low and high frequencies, acting as a bandpass filter. This would not make  $A_n$  invertible, but it does make it ill-conditioned. As before, increasing the regularization parameter  $\nu_n$  partly mitigates this.

**4.2. Covariance estimator.** To capture the variability of the volume density  $x$  as a random vector in  $\mathbb{R}^p$ , we consider its covariance  $\text{Cov}[x]$ . The same construction outlined in the previous section to estimate the mean can be used to estimate the covariance. Specifically, for each  $s = 1, \dots, n$ , computing the covariance of (11) conditioned on  $P_s$  gives

$$(20) \quad \text{Cov}_{x_s, e_s | P_s}[y_s] = P_s \text{Cov}[x] P_s^\top + \sigma^2 \mathbf{I}_{N^2},$$

where  $\text{Cov}[e] = \sigma^2 \mathbf{I}_{N^2}$ . The left-hand side is the covariance of the image  $y_s$  where  $P_s$  is fixed, but  $x_s$  and  $e_s$  are allowed to vary. To compute it, we would need an infinite number of realizations of  $y_s$  for a fixed viewing direction and CTF. However, we do not have an infinite number of images. We are only guaranteed to have one, but we can use this image to estimate the conditional covariance as in

$$(21) \quad (y_s - P_s \mu_n)(y_s - P_s \mu_n)^\top \approx \text{Cov}_{x_s, e_s | P_s}[y_s].$$

The left-hand side is available to us since we have already estimated  $\mu_n$  and  $P_s$  is known. In expectation it equals the right-hand side, so it forms an unbiased estimator, albeit again with high variance. Plugging this into (20), we obtain

$$(22) \quad (y_s - P_s \mu_n)(y_s - P_s \mu_n)^\top \approx P_s \text{Cov}[x] P_s^\top + \sigma^2 \mathbf{I}_{N^2}.$$

The high variance in the left-hand side makes this a loose constraint and we do not expect it to hold exactly. Instead, we compute the mean squared error between the two sides and attempt to minimize it over all images. This yields the regularized least-squares estimator  $\Sigma_n$  of  $\text{Cov}[x]$  defined by

$$(23) \quad \Sigma_n := \underset{\Sigma}{\text{argmin}} \frac{1}{n} \sum_{s=1}^n \left\| (y_s - P_s \mu_n)(y_s - P_s \mu_n)^\top - (P_s \Sigma P_s^\top + \sigma^2 \mathbf{I}_{N^2}) \right\|_{\text{F}}^2 + \xi_n \|\Sigma\|_{\text{F}}^2,$$

where  $\xi_n \geq 0$  is a regularization parameter. As with estimating the mean, any potential ill-posedness can be mitigated by the regularization term  $\xi_n \|\Sigma\|_{\text{F}}^2$ . Typically, the problem is less ill-posed at large  $n$ , so  $\xi_n$  decreases with growing  $n$ .

Similar to  $\mu_n$ , this estimator finds the covariance matrix that, when projected along  $R_s$  and convolved by  $h_s$  according to  $\Sigma \mapsto P_s \Sigma P_s^\top$ , minimizes the square Frobenius distance

to the outer product of the mean-subtracted images. Note that  $\Sigma_n$  is not the regularized maximum-likelihood estimator for  $\text{Cov}[x]$ . However, as we shall see later in this section, it converges to  $\text{Cov}[x]$  as  $n \rightarrow \infty$  with high probability under a wide range of conditions.

To solve the least-squares optimization problem in (23), we again differentiate and set the derivative to zero, obtaining

$$(24) \quad L_n(\Sigma_n) = B_n,$$

where

$$(25) \quad L_n(\Sigma) := \frac{1}{n} \sum_{s=1}^n P_s^\top P_s \Sigma P_s^\top P_s + \xi_n \mathbf{I}_{p^2},$$

$$(26) \quad B_n := \frac{1}{n} \sum_{s=1}^n P_s^\top (y_s - P_s \mu_n)(y_s - P_s \mu_n)^\top P_s - \sigma^2 P_s^\top P_s.$$

Calculating least-squares estimator defined in (23) is therefore equivalent to solving the linear system (24). In the same spirit as  $b_n$ , the right-hand side matrix  $B_n$  averages the backprojected outer product covariance estimators for each image with a noise term correction. The covariance projection-backprojection operator  $L_n$  plays the same role as  $A_n$  by describing the reweighting of the backprojected covariance matrix estimators.

As for the mean estimator, the Tikhonov regularization term in (23) may be replaced by other regularization terms. We also note that  $\Sigma_n$  is not constrained to be positive semi-definite and so may not qualify as a covariance matrix. Again, however, we are only interested in the dominant eigenvectors, so imposing this condition would not appreciably alter our results at the cost of significant computational expense.

Each of the constraints in (22) are loose because of the noise and the potentially large amount of variability in  $x$ . Consequently, it may appear that their derived least-squares estimator (23) would be a poor one. However, it has been shown that, given enough images,  $\Sigma_n$  does in fact provide a reasonable estimate of  $\text{Cov}[x]$ . Indeed, Katsevich et al. [38] have shown that, for  $\xi_n = 0$ ,

$$(27) \quad \|\Sigma_n - \text{Cov}[x]\|_F = \mathcal{O}\left(\frac{\log n}{\sqrt{n}}\right),$$

with high probability as  $n \rightarrow \infty$ , as long as that  $x$  is bounded and  $L_n$  is invertible when  $n > n_0$  for some  $n_0 > 0$ . This invertibility condition is satisfied when the central planes defined by the viewing directions  $R_1, \dots, R_n$  contain enough frequency pairs in the Fourier domain and the zeros of the CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$  do not overlap. Note that this requirement on the viewing directions is more strict than the corresponding requirement for the invertibility of  $A_n$ . Indeed, since the Fourier transform of the covariance matrix  $\text{Cov}[x]$  describes the correlation between any pairs of frequencies in the 3D Fourier domain, each pair must be represented in the data.

To clarify this, we return to the toy example introduced in the previous section. As before, we have a random vector  $x = [a, b, c]$  containing three values and observations  $y_s$  of the form  $[a_s, b_s, ?]$ ,  $[a_s, ?, c_s]$ , and  $[?, b_s, c_s]$  for  $s = 1, \dots, n$ . It is not possible to reconstruct the joint probability density of  $x = [a, b, c]$ , but we can estimate its covariance. Indeed, to estimate the

variance of  $a$ , we collect all observations of the form  $[a_s, b_s, ?]$  and  $[a_s, ?, c_s]$ , extract the first entry, and compute its variance. Similarly, for the covariance between  $a$  and  $b$ , we take all observations of the form  $[a_s, b_s, ?]$  and compute the covariance between the first two entries. Proceeding like this, we can “fill up” an entire matrix, which is a consistent estimator of the population covariance (that is, it converges to the latter as  $n \rightarrow \infty$ ). Note that at no point do we need to draw samples of the complete vector  $x = [a, b, c]$  or characterize its full distribution in order to estimate the covariance. All that is necessary is that we know which observation has what entries missing (which entry has the  $?$  in place of a value) and that all pairs of entries are observed. However, it is not possible to estimate all third-order moments in this scenario since this would require observing all three entries at once.

We can draw parallels to the cryo-EM covariance estimation problem by again making use of the Fourier Slice Theorem, which tells us that the imaging operator  $P$  in the Fourier domain corresponds to extracting a slice of the Fourier transform of a volume. The entries of  $\mathcal{F}x$  lying on the plane defined by  $R$  are kept, while others are discarded. The remaining entries are then multiplied by the CTF  $\mathcal{F}h$ . If we are to successfully estimate the covariance, we have to make sure that all pairs of 3D frequencies appear in our observed projections. For any given pair, we could then find the projection images whose Fourier slices contain that pairs and use these to compute the covariance. A given pair of frequencies, together with the origin, uniquely define a plane in 3D. Consequently, all such central planes must be present in our sample. In other words, the set of rotations  $R_1, \dots, R_n$  must cover all of  $\text{SO}(3)$ . This does not mean that they have to be uniformly distributed, but that any given rotation in  $\text{SO}(3)$  is sampled with non-zero probability.

Again, this is much stricter than the requirement for accurate reconstruction of the mean  $\mu_n$ , where a great circle of viewing directions is sufficient. This requirement was first observed by Liu & Frank [49], who argue that its stringency precludes accurate estimation of the covariance, which they refer to as “type-II variance.” In our work, this problem is mitigated by several factors. The first is that while we attempt to estimate the whole covariance matrix  $\text{Cov}[x]$ , our ultimate goal is its top eigenvectors. It follows from the Davis-Kahan  $\sin(\theta)$  theorem that it is possible to accurately estimate the leading eigenvectors of a matrix provided its eigenvalues are well-separated from the remaining eigenvalues (that is, there is a large eigengap) [16]. As will be described in the Section 4.4, random matrix theory suggests that such a separation in the eigenvalues of  $\Sigma_n$  does indeed appear as  $n$  grows. As a result, we can expect to obtain good estimates for the top eigenvectors of  $\Sigma_n$  even though  $\Sigma_n$  is not very accurate overall. The second point is that our proposed algorithm is typically applied at low resolution  $N$ . As a result, we only need the viewing directions to cover the sphere up to this low resolution (i.e., gaps smaller than  $1/N$  in the distribution of rotations are acceptable). Finally, adjusting the regularization parameter  $\xi_n$  allows us to regularize the entries of  $\Sigma_n$  whose corresponding pairs of frequencies are missing from the data.

To see how  $\Sigma_n$  performs this estimation, we apply the Fourier Slice Theorem to the continuous covariance matrix  $\mathcal{C} : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$  satisfying  $\mathcal{C} \in L^1(\mathbb{R}^3 \times \mathbb{R}^3)$ . This gives

$$(28) \quad (\mathcal{F} \times \mathcal{F})(\mathcal{P} \times \mathcal{P})\mathcal{C}(\omega_1, \omega_2) = (\mathcal{F} \times \mathcal{F})\mathcal{C}(R^T[\omega_1; 0], R^T[\omega_2; 0])\mathcal{F}h(\omega_1)\mathcal{F}h(\omega_2),$$

where  $(\mathcal{A} \times \mathcal{B})$  denotes the mapping that applies  $\mathcal{A}$  along the first variable and  $\mathcal{B}$  along the second variable [78]. As a result, projecting the Fourier transform of a covariance matrix

along its columns and rows therefore corresponds to restriction to frequency pairs belonging to a certain plane defined by  $R$  followed by multiplication by the CTF.

The dual formulation of (28) says that backprojecting a two-dimensional covariance matrix (i.e., applying  $(\mathcal{P}^T \times \mathcal{P}^T)$ ) corresponds to inserting its Fourier transform into a three-dimensional Fourier transform along pairs of frequencies both belonging to a certain plane. Each term in the sum (26) defining  $B_n$  therefore takes the two-dimensional matrix estimate  $(y_s - P_s \mu_s)(y_s - P_s \mu_s)^T - \sigma^2 I_{N^2}$ , places its Fourier transform along the correct plane in three-dimensional covariance Fourier space and multiplies by the appropriate CTF. These are then averaged across all images. Slice by slice, this provides a “reconstruction” of the three-dimensional covariance matrix.

Much like the case for mean estimation, however, this reconstruction by backprojection needs to be reweighted in order to obtain an accurate covariance estimate. The weighting is encoded by the covariance projection-backprojection operator  $L_n$  and depends on the distribution of viewing directions and CTFs. In the case of uniform distribution of viewing angles, the same consideration of  $A_n$  applies, where frequencies closer to the origin are weighted higher with respect to frequencies farther away. However, for the covariance we also need to take into account the relationship between pairs of frequencies. Indeed, for a given pair of frequencies, its weight depends on how many central planes, that is images, pass through both frequencies. Since frequencies that are nearly co-linear have more central planes passing through them, this results in higher weights compared to other pairs. By inverting  $L_n$ , we renormalize the backprojected covariance estimate  $B_n$  with respect to this density. For a more detailed discussion of this phenomenon, see Katsevich et al. [38].

As discussed above, the invertibility of  $L_n$  depends on the viewing directions  $R_1, \dots, R_n$  sufficiently covering  $\text{SO}(3)$  and the zeros of the CTFs not completely overlapping. However, as with  $A_n$ ,  $L_n$  may still be highly ill-conditioned if certain viewing directions are more common than others. Higher frequencies are also less well-conditioned, so a high resolution  $N$  gives a worse conditioning for  $L_n$ . Finally, the bandpass effect of the CTF increases its condition number. These can again be mitigated by an appropriate choice of the regularization parameter  $\xi_n$ , ensuring that  $L_n$  is well-conditioned without introducing too much bias.

**4.3. Resolution limits.** While their construction is similar, the least-squares mean  $\mu_n$  and covariance  $\Sigma_n$  estimators differ in their well-posedness and conditioning properties, as seen in the previous section. For a fixed number of images  $n$ , this results in the achievable resolution  $N$  for the covariance estimator being significantly lower compared to that of the mean.

To illustrate, we estimate the mean  $\mathbb{E}[x]$  at resolution  $N$ . Using the Fourier Slice Theorem, Section 4.1 shows that this is achieved by “filling up” the 3D Fourier domain with  $n$  central slices (each corresponding to a projection image) to obtain  $b_n$  and applying  $A_n^{-1}$ . Each central slice has  $\mathcal{O}(N^2)$  points, yielding a total of  $\mathcal{O}(nN^2)$  points. Since the 3D Fourier domain contains  $\mathcal{O}(N^3)$  points, we require that  $nN^2 \gg N^3$ , so  $n$  must be of order of at least  $N$ .

If we have clean data,  $n = \mathcal{O}(N)$  images would suffice. However, adding noise of variance  $\sigma^2$ , more images are needed for an accurate estimate. Specifically, for each of the  $\mathcal{O}(N^3)$  points in the 3D Fourier domain, we need  $\mathcal{O}(\sigma^2)$  samples to reduce the noise in  $b_n$  to order 1. The total number of required samples is therefore  $\mathcal{O}(N^3 \sigma^2)$ , so  $n$  has to be of order of at least  $N \sigma^2$ . This does not account for non-uniform distributions of viewing directions (which

increases the constant of proportionality) or the effect of the CTF, the power spectrum of the volumes, and non-white noise (which increase the number of samples necessary to estimate the higher frequencies). However, this provides a lower bound, requiring the number of images to be at least proportional to the desired resolution times the noise variance.

For the covariance, on the other hand, each image contributes  $\mathcal{O}(N^4)$  entries in the Fourier domain 3D covariance matrix as described in Section 4.2. The number of entries to estimate in the covariance matrix is  $\mathcal{O}(N^6)$ . Consequently, we need  $n$  to be at least of order  $N^2$  to adequately estimate  $\text{Cov}[x]$  from clean data. From another perspective, the viewing directions need to cover the unit sphere with separation  $1/N$ , so an order of  $N^2$  images is required.

Each term of the sum in  $B_n$  concerns the outer product of a (mean-centered) image with itself. As a result, adding noise of variance  $\sigma^2$  to the images results in noise of variance  $\sigma^4$  in each term. To reduce the total variance of the noise in  $B_n$  to  $\mathcal{O}(1)$ , we therefore need  $\mathcal{O}(N^6\sigma^4)$  samples. Consequently,  $n$  must be of order at least  $N^2\sigma^4$ . Again, this is an idealized setting, but this relationship provides a reasonable lower bound for  $n$ .

The difference in the required number of images for  $\mu_n$  and  $\Sigma_n$  is quite stark. Instead fixing the number of images  $n$ , we see that the best resolution  $N$  that we can achieve for  $\mu_n$  is  $\mathcal{O}(n\sigma^{-2})$ . The resolution limit for  $\Sigma_n$ , on the other hand, is  $\mathcal{O}(\sqrt{n}\sigma^{-2})$ . In other words, to increase the resolution by a factor of two, we need four times as many images. Estimating the covariance is therefore not only more computationally demanding (in terms of running time and memory usage), but is also more demanding in terms of data.

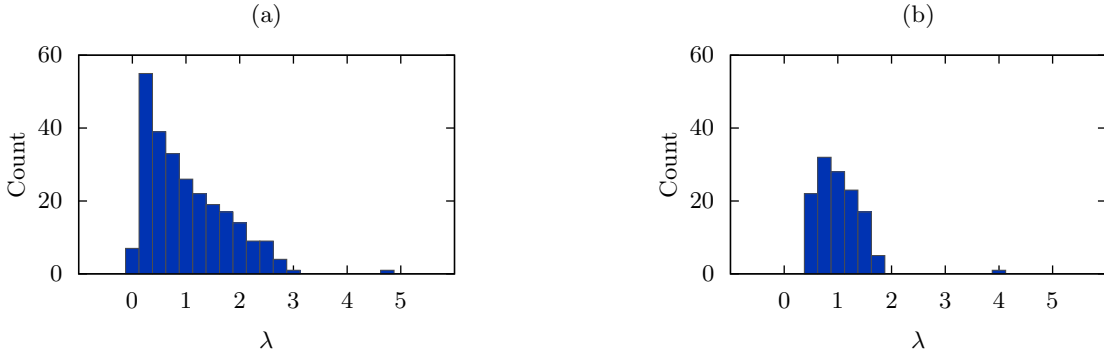
More concretely, let us suppose we have  $n = 10000$  images with  $\sigma^2 = 10$ . Assuming the clean projection images have a mean square intensity of order 1, this implies a signal-to-noise ratio of 0.1. The highest achievable resolution  $N$  for the  $\mu_n$  is then 1000. Again, this is an upper bound. The achieved resolution is much lower in practice due to the effects of pixel size, CTF, and lower signal-to-noise ratio at high frequencies. In contrast, the maximum achievable resolution  $N$  for  $\Sigma_n$  in our idealized setting is only 10.

One way to overcome these limitations is to impose strong assumptions on the covariance. For example, that it is low-rank, satisfies certain smoothness or sparsity constraints, or that it is generated by certain types of deformations. In this work, we make use of the low-rank property since we ultimately extract the leading eigenvectors of our covariance estimate  $\Sigma_n$ . This may not be optimal, however, as the low-rank constraint is not imposed when estimating the covariance matrix. We will explore these directions in future work.

**4.4. High-dimensional PCA.** The consistency result (27) shows that  $\Sigma_n$  converges to  $\text{Cov}[x]$  as  $n \rightarrow \infty$ . However, in many applications, while  $n$  may be large, it is not necessarily large with respect to the size  $N^3$  of the volume vectors.

In this case, a more appropriate setting is to consider the behavior of  $\Sigma_n$  as  $n$  and  $N$  both tend to infinity, but at potentially different rates. Indeed, Section 4.3 suggests that  $N$  should grow no faster than  $\sqrt{n}$  to ensure estimation is well-posed. To better understand the behavior of  $\Sigma_n$  in this high-dimensional regime, we first review related results from the literature on sample covariance. Let us consider the sample covariance of a set of independently sampled Gaussian noise vectors  $w_1, \dots, w_n$

$$(29) \quad \frac{1}{n} \sum_{s=1}^n w_s w_s^\top.$$



**Figure 2.** The eigenvalue distribution of the sample covariance matrix for (a)  $p = 256$ ,  $n = 512$  and (b)  $p = 128$ ,  $n = 1024$  with  $\sigma = 1$  and  $\ell = 3$  in both regimes. For (a), we have  $\gamma = 1/2$  and the spiked covariance model predicts a maximum noise eigenvalue at  $(1 + \sqrt{1/2})^2 \approx 2.91$  and a signal eigenvalue at  $\lambda(3, 1/2) \approx 4.67$ , while for (b),  $\gamma = 1/8$  gives  $(1 + \sqrt{1/8})^2 \approx 1.83$  and  $\lambda(3, 1/8) \approx 4.17$ .

where  $\text{Cov}[w_1] = \dots = \text{Cov}[w_n] = \sigma^2 \mathbf{I}_p$  for some dimension  $p > 1$ .

In the low-dimensional regime where  $n \gg p$ , all eigenvalues of this sample covariance are concentrated around the single population eigenvalue  $\sigma^2$ . However, for  $n, p \rightarrow \infty$  where  $p/n \rightarrow \gamma < 1$ , the spectrum will instead spread between  $\sigma^2(1 - \sqrt{\gamma})^2$  and  $\sigma^2(1 + \sqrt{\gamma})^2$ , following the Marčenko-Pastur distribution [52].

In the spiked covariance model [36], we have a clean signal  $a_s = \sqrt{\ell} v z_s$  for  $s = 1, \dots, n$ , where  $v$  is a unit vector,  $z_1, \dots, z_n$  are i.i.d., zero-mean and unit variance random variables, while  $\ell$  is the signal strength. The covariance of  $a_s$  is then equal to  $\ell v v^T$  and has rank one. Adding noise then gives the measurements

$$(30) \quad d_s = a_s + w_s,$$

for all  $s = 1, \dots, n$ . As before, the sample covariance is

$$(31) \quad \frac{1}{n} \sum_{s=1}^n d_s d_s^T.$$

When  $n \gg p$ , its spectrum converges to  $\{\ell + \sigma^2, \sigma^2, \dots, \sigma^2\}$ , with the dominant eigenvector equal to  $v$ . However, when  $n, p \rightarrow \infty$  and  $p/n \rightarrow \gamma < 1$ , there are two possible scenarios. If  $\ell/\sigma^2 < \sqrt{\gamma}$ , the spectrum will be the same as the pure noise case—the signal is lost in the noise. If instead  $\ell/\sigma^2 \geq \sqrt{\gamma}$ , all eigenvalues will follow the Marčenko-Pastur distribution except one [58]. This signal eigenvalue will “pop out” at

$$(32) \quad \lambda(\ell, \gamma) = (\sigma^2 + \ell)(1 + \gamma\sigma^2/\ell).$$

These distributions are illustrated for two regimes  $\gamma = 1/2$  and  $\gamma = 1/8$  in Figure 2. As  $\ell/(\sigma^2 \sqrt{\gamma})$  increases, the dominant eigenvector converges to  $v$  [58]. Specifically, the square correlation  $|\langle v, u \rangle|^2$  between  $v$  and the dominant eigenvector  $u$  of (31) tends to

$$(33) \quad c(\ell, \gamma) = \frac{1 - \gamma\sigma^4/\ell^2}{1 + \gamma\sigma^2/\ell}.$$



The spiked covariance model suggests a solution for estimating  $v$  from the noisy observations  $d_1, \dots, d_n$ . For signal covariance estimation, the eigenvalues below  $\sigma^2(1 + \sqrt{\gamma})^2$  are set to zero while those above are shrunk by an appropriate amount. A first approach would be to invert (32) to obtain

$$(34) \quad \ell(\lambda, \gamma) = \frac{1}{2} \left( \lambda + \sigma^2(1 - \gamma) + \sqrt{(\lambda + \sigma^2(1 - \gamma))^2 - 4\sigma^2\lambda} \right) - \sigma^2$$

and replace a given eigenvalue  $\lambda$  by  $\ell(\lambda, \gamma)$ . This “shrinks” the eigenvalues of the covariance matrix to provide better approximations of the population eigenvalues and consequently a better approximation of the covariance. Other functions can similarly be used to improve the estimation of the eigenvalues and are known as “shrinkers.” This leads to the question of which shrinker is optimal given a certain loss function on the covariance matrix. Such shrinkers have been derived by Donoho et al. [19] for 26 different loss functions.

This more general approach succeeds quite well in recovering covariance matrices for high-dimensional data. Given a sample covariance matrix, we calculate its eigendecomposition

$$(35) \quad \frac{1}{n} \sum_{s=1}^n d_s d_s^T = \sum_{m=1}^p \lambda_m u_m u_m^T,$$

where  $\{v_1, \dots, v_p\}$  form an orthonormal basis. The eigenvalues  $\{\lambda_1, \dots, \lambda_p\}$  are transformed using a shrinker function  $\rho$  into  $\{\rho(\lambda_1), \dots, \rho(\lambda_p)\}$ . Putting everything back together gives us

$$(36) \quad \sum_{m=1}^p \rho(\lambda_m) u_m u_m^T,$$

which is an estimator for the population covariance  $\text{Cov}[a]$ .

Donoho et al. showed how to choose the shrinker  $\rho$  to optimize the error with respect to some loss function on the covariance estimator [19]. The shrinker which achieves the lowest expected loss in the Frobenius norm as  $n, p \rightarrow \infty$  is given by

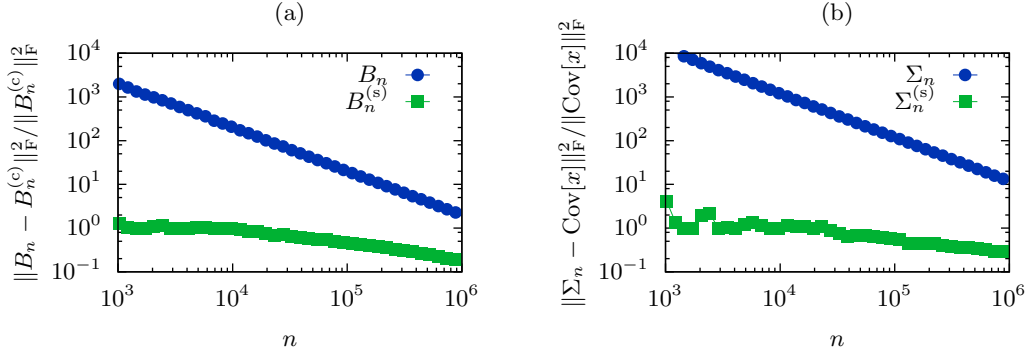
$$(37) \quad \rho(\lambda) = \ell(\lambda, \gamma) c(\ell(\lambda, \gamma), \gamma).$$

Among all shrinkage estimators of the form (36), the shrinker  $\rho$  given by (37) provides the lowest expected loss in the Frobenius norm [19] as  $n, p \rightarrow \infty$ . The expected loss with respect to the operator norm is minimized by  $\rho(\lambda) = \ell(\lambda, \gamma)$ . The authors consider a variety of norms, each of which is assigned a corresponding optimal shrinker [19]. Since our least-squares objective is formulated with respect to the Frobenius norm, we shall use the corresponding shrinker (37) in the following. Note that this estimator is not restricted to rank-one signals  $a$  but is optimal for arbitrary fixed finite rank.

By slight abuse of notation, we extend the action of  $\rho$  from scalars to symmetric matrices by its action on the eigenvalues, so that

$$(38) \quad \rho(A) := \rho \left( \sum_{m=1}^p \lambda_m v_m v_m^T \right) = \sum_{i=1}^p \rho(\lambda_m) v_m v_m^T,$$

provided  $\sum_{m=1}^p \lambda_m v_m v_m^T$  is an eigendecomposition of  $A$ .



**Figure 3.** Effect of shrinkage on the relative error of (a)  $B_n$  and (b)  $\Sigma_n$  for simulated data with  $N = 16$ ,  $C = 2$  and 7 distinct CTFs. The signal-to-noise ratio (see (91)) is 0.001.

**4.5. Shrinkage of  $B_n$ .** We can apply the above ideas to obtain a better estimate for the right-hand side  $B_n$ . This is the first major improvement over previous work on least-squares estimators for 3D covariance matrices in cryo-EM [38, 3]. Specifically, we will replace  $B_n$  by

$$(39) \quad B_n^{(s)} = A_n^{1/2} \rho(A_n^{-1/2} B_n A_n^{-1/2} + \sigma^2 I_p) A_n^{1/2},$$

where  $A_n$  is given by (17).

For clean data,  $\sigma = 0$  so  $y_s = P_s x_s$ . In this case, the definition of  $B_n$  (26) gives

$$(40) \quad B_n^{(c)} = \frac{1}{n} \sum_{s=1}^n P_s^\top (y_s - P_s \mu_n) (y_s - P_s \mu_n)^\top P_s,$$

Plugging in our forward model (12) for  $y_s$ , we get

$$(41) \quad B_n^{(c)} = \frac{1}{n} \sum_{s=1}^n P_s^\top P_s (x_s - \mu_n) (x_s - \mu_n)^\top P_s^\top P_s,$$

which is the sample covariance of the mean-subtracted and projected-backprojected volumes  $P_s^\top P_s x_s$  for  $s = 1, \dots, n$ .

We would now like to obtain something close to (41) also in the noisy case. Mean-subtracting and backprojecting noisy images  $y_s$ , we have

$$(42) \quad P_s^\top (y_s - P_s \mu_n) = P_s^\top P_s (x_s - \mu_n) + P_s^\top e_s.$$

This is similar to the spiked covariance model (30), except the noise term is not white but has covariance  $\mathbb{E}[P^\top P]$ . We can approximate the noise covariance using the regularized estimator

$$(43) \quad A_n = \frac{1}{n} \sum_{s=1}^n P_s^\top P_s + \nu_n I_p,$$

since  $\mathbb{E}[A_n] = \mathbb{E}[P^\top P] + \nu_n I_p$  and the law of large numbers guarantees that  $A_n$  converges to its expectation as  $n \rightarrow \infty$ . The regularization term  $\nu_n I_p$  ensures that the inverse of  $A_n$  will

be bounded. Note that the high-dimensional phenomena studied in the previous section do not appear here since  $P$  is sampled from a low-dimensional space of fixed dimension. Indeed, the viewing direction  $R$  is sampled from  $\text{SO}(3)$  (which has dimension three) while the CTF  $\mathcal{F}h$  depends on two defocus values and an azimuth angle, yielding a total of 6 dimensions.

Multiplying the backprojected images by  $A_n^{-1/2}$  whitens the noise, and we define

$$(44) \quad z_s = A_n^{-1/2} P_s^T (y_s - P_s \mu_n) = A_n^{-1/2} P_s^T P_s (x_s - \mu_n) + A_n^{-1/2} P_s^T e_s.$$

We now apply the standard shrinkage operator  $\rho$  defined in (37) and obtain

$$(45) \quad \rho \left( \frac{1}{n} \sum_{s=1}^n z_s z_s^T \right) = \rho \left( A_n^{-1/2} B_n A_n^{-1/2} + \sigma^2 \mathbf{I}_p \right).$$

Conjugating the shrunken covariance by  $A_n^{1/2}$ , we obtain a shrinkage variant  $B_n^{(s)}$  of  $B_n$  as

$$(46) \quad B_n^{(s)} := A_n^{1/2} \rho \left( A_n^{-1/2} B_n A_n^{-1/2} + \sigma^2 \mathbf{I}_p \right) A_n^{1/2},$$

providing an estimator of the clean right-hand side  $B_n^{(c)}$ . Note that the loss is minimized with respect to the Frobenius norm on  $A_n^{-1/2} B_n^{(s)} A_n^{-1/2}$ , not  $B_n^{(s)}$ , so there still might be some room for improvement. This is the subject of future work.

Replacing  $B_n$  by  $B_n^{(s)}$  in (24) yields a more accurate estimator  $\Sigma_n^{(s)} = L_n^{-1}(B_n^{(s)})$  since  $B_n^{(s)}$  is a more accurate estimate of  $B_n^{(c)}$  compared to  $B_n$ . We will refer to  $\Sigma_n^{(s)}$  as the shrinkage covariance estimator, in contrast to the standard  $\Sigma_n$  least-squares covariance estimator. To evaluate the effect of the shrinkage on estimation accuracy, we plot the error of  $B_n$  and  $B_n^{(s)}$  with respect to the clean  $B_n^{(c)}$  as a function of  $n$  in Figure 3(a). For the values of  $n$  considered, shrinkage introduces a significant reduction in error. The same effect is found for the resulting estimators  $\Sigma_n$  and  $\Sigma_n^{(s)}$  when compared with  $\text{Cov}[x]$  in Figure 3(b).

**5. Efficient computation.** The covariance estimators as formulated in the previous section are computed by solving the corresponding normal equations. However, direct matrix inversion is intractable for typical problem sizes. We therefore consider an iterative solution based on the conjugate gradient method applied to a convolutional formulation of the normal equations. We first illustrate this for the mean least-squares estimator  $\mu_n$  and then generalize this technique to the covariance estimator  $\Sigma_n$ . To speed up convergence, we employ circulant preconditioners for both  $A_n$  and  $L_n$ . The use of the conjugate gradient method to estimate the covariance was previously considered in Andén et al. [3], but this work lacked the convolutional formulation and appropriate preconditioners necessary for rapid convergence.

**5.1. Mean deconvolution.** The normal equations (16) for the mean estimator  $\mu_n$  can be solved by calculating the matrix  $A_n$ , the right-hand side  $b_n$ , and solving for  $\mu_n$  in  $A_n \mu_n = b_n$ . This is a linear system in  $N^3$  variables, so solving it directly has complexity  $\mathcal{O}(N^9)$ . A more sophisticated approach is therefore needed. Here, we shall exploit the convolutional structure of  $A_n$ . This approach has been successful in several related applications [87, 23, 29, 84], but we shall focus on its use in the Fourier-based iterative reconstruction method (FIRM) introduced

by Wang et al. [88]. This section will rederive that method with the goal of applying these ideas to the computation of  $\Sigma_n$ .

We first note that the projection-backprojection operator  $P_s^\top P_s$  is factored into  $Q^\top I_s^\top I_s Q$ , where  $I_s : \mathbb{R}^{N^3} \rightarrow \mathbb{R}^{N^2}$  is the voxel projection mapping corresponding to  $P_s$ . In the voxel domain,  $I_s^\top I_s$  is a convolution. Indeed, in the continuous case, projection integrates along a certain viewing direction and convolves with a point spread function, while backprojection “fills up” a volume along a certain viewing angle using an image convolved with that point spread function. The resulting volume is then constant along that viewing direction. The projection-backprojection operator is therefore a low-pass filter.

In the frequency domain, the Fourier slice theorem tells the same story. Projection is restriction of the volume Fourier transform to a plane followed by multiplication by the CTF. Backprojection multiplies a two-dimensional Fourier transform by a CTF and inserts it into a plane in a three-dimensional Fourier transform, filling the rest with zeros. The combined projection-backprojection operator is therefore a multiplication by a Dirac delta function along the projection direction times the squared CTF along the transverse direction.

Having engineered our voxel discretization  $I_s$  of the projection operator to satisfy a discrete Fourier slice theorem (10), these properties carry over to the discrete case. Consequently,

$$(47) \quad I_s^\top I_s x(\mathbf{i}) = x * k_s(\mathbf{i}) = \sum_{\mathbf{j} \in M_{2N-1}^3} x(\mathbf{i} - \mathbf{j}) k_s(\mathbf{j}),$$

where  $*$  denotes convolution and

$$(48) \quad k_s(\mathbf{i}) := \frac{1}{N^4} \sum_{\mathbf{l} \in M_{N-1}^2} |\mathcal{F}h_s(\mathbf{l})|^2 e^{\frac{2\pi i}{N} \langle \mathbf{i}, R_s^\top \mathbf{l}; 0 \rangle},$$

for all  $\mathbf{i} \in M_{2N-1}^3$ . This follows from calculating the dual  $I_s^\top$  of the voxel projection matrix (9) and applying it to the matrix  $I_s$  itself. In other words,  $I_s^\top I_s$  is a 3-Toeplitz matrix. The projection-backprojection operator  $P_s^\top P_s$  is thus factored into a basis evaluation  $Q$ , a convolution  $I_s^\top I_s$ , and a basis expansion  $Q^\top$ .

In the definition (17) of  $A_n$ , plugging in (47) now gives

$$(49) \quad \begin{aligned} A_n x &= \frac{1}{n} \sum_{s=1}^n P_s^\top P_s x + \nu_n x = \frac{1}{n} \sum_{s=1}^n Q^\top (Qx * k_s) + \nu_n x = Q^\top \left( Qx * \frac{1}{n} \sum_{s=1}^n k_s \right) + \nu_n x \\ &= Q^\top (Qx * f_n) + \nu_n x, \end{aligned}$$

where

$$(50) \quad f_n := \frac{1}{n} \sum_{s=1}^n k_s.$$

The sum over  $n$  is therefore factorized as basis evaluation  $Q$ , followed by application of a 3-Toeplitz matrix, then a basis expansion  $Q^\top$ .

The convolution kernel  $f_n$  can be calculated in one pass over the dataset using NUFFTs with complexity  $\mathcal{O}(N^3 \log N + nN^2)$ . Once this has been calculated, however, each application

of  $A_n$  using the convolutional formulation of (49) is achieved in  $\mathcal{O}(N^3 \log N)$  time using FFTs, independent of the number of images.

We can exploit this fast application of  $A_n$  to solve the system  $A_n \mu_n = b_n$ . Specifically, we apply the conjugate gradient (CG) method, which is an iterative algorithm for solving linear systems [32]. It computes an approximate solution at each iteration through a single application of  $A_n$ , so its efficiency depends on being able to apply this operator fast, which is the case for Toeplitz operators, as seen above [9]. To reach a given accuracy  $\mathcal{O}(\sqrt{\kappa(A_n)})$  iterations are needed, where  $\kappa(A_n)$  is the condition number of  $A_n$ . As we shall see,  $\kappa(A_n)$  can be reduced by the circulant preconditioner described in Section 5.3.

The above algorithm therefore consists of two steps. First, we precalculate the kernel  $f_n$  and the right-hand side  $b_n$ . These are both achieved in  $\mathcal{O}(N^3 \log N + nN^2)$  using NUFFTs. Second,  $\mathcal{O}(\sqrt{\kappa(A_n)})$  iterations CG are performed, each at a cost of  $\mathcal{O}(N^3 \log N)$ . The overall complexity is then  $\mathcal{O}(\sqrt{\kappa(A_n)} N^3 \log N + nN^2)$ . Note that this method is nearly optimal in the sense that simply reading the images requires  $\mathcal{O}(nN^2)$  operations, while the reconstructed volume takes up  $\mathcal{O}(N^3)$  in memory.

**5.2. Covariance deconvolution.** As discussed in the previous section, directly solving the normal equations for  $\mu_n$  can be computationally expensive. This is also the case for the covariance estimator  $\Sigma_n$ , which scales worse in  $N$ . Indeed, volume vectors  $x$  are of size  $N^3$  so the covariance estimate  $\Sigma_n$  is of size  $N^3$ -by- $N^3$  and thus contains  $N^6$  elements. Since  $L_n$  maps covariance matrices to covariance matrices, the matrix representation of  $L_n$  requires  $(N^6)^2 = N^{12}$  elements, which stored at single precision occupies 256 GB for  $N = 8$ . Direct inversion of this matrix would have computational complexity  $\mathcal{O}(N^{18})$ .

Previously, Katsevich et al. [38] defined a volume basis based on spherical harmonics in which  $L = \lim_{n \rightarrow \infty} L_n$  is a block diagonal matrix with sparse blocks under certain conditions. Unfortunately, the approximation is only valid if  $R$  is uniformly distributed on  $\text{SO}(3)$  and the CTF is fixed. These conditions typically do not hold for experimental data. In addition, the approximation of  $L_n$  by  $L$  only holds as  $n \rightarrow \infty$  and is not appropriate for smaller datasets. These problems are mitigated by solving the exact system  $L_n(\Sigma_n) = B_n$  using the CG method [3], but this approach passes through the entire dataset at each iteration and converges slowly.

A more practical approach is to apply the ideas from the mean estimation algorithm described in the previous section. According to (25),  $L_n$  is a sum of linear matrix operators

$$(51) \quad \Sigma \mapsto P_s^T P_s \Sigma P_s^T P_s$$

plus a regularization term  $\xi_n I_{p^2}$ . Since  $P_s^T P_s$  can be factored into basis evaluation/expansion and convolution in three dimensions, this mapping enjoys a similar factorization

$$(52) \quad P_s^T P_s \Sigma P_s^T P_s = Q^T (I_s^T I_s (Q \Sigma Q^T) I_s^T I_s) Q.$$

The conjugation by  $I_s^T I_s$  convolves both the rows and the columns of the matrix by  $k_s$ , which is a convolution in six dimensions by the outer product of  $k_s$  with itself. Specifically, we have

$$(53) \quad I_s^T I_s Z I_s^T I_s = Z * K_s,$$

for any matrix  $Z \in \mathbb{R}^{N^3 \times N^3}$ , where

$$(54) \quad K_s[\mathbf{i}_1, \mathbf{i}_2] := k_s[\mathbf{i}_1] k_s[\mathbf{i}_2],$$

for all  $(\mathbf{i}_1, \mathbf{i}_2) \in M_{2N-1}^6$ . Consequently

$$(55) \quad P_s^\top P_s \Sigma P_s^\top P_s = Q^\top (Q \Sigma Q^\top * K_s) Q.$$

One advantage of this formulation is that we average the convolutional kernels over the whole dataset to obtain a convolutional representation for  $L_n$ . This gives

$$(56) \quad L_n(\Sigma) = Q^\top (Q \Sigma Q^\top * F_n) Q + \xi_n \Sigma,$$

where

$$(57) \quad F_n := \frac{1}{n} \sum_{s=1}^n K_s.$$

Similar to  $A_n$ , the sum over  $s$  in  $L_n$  is factored into basis evaluations/expansions and a 6-Toeplitz matrix operator, allowing for rapid calculation of  $L_n(\Sigma)$ .

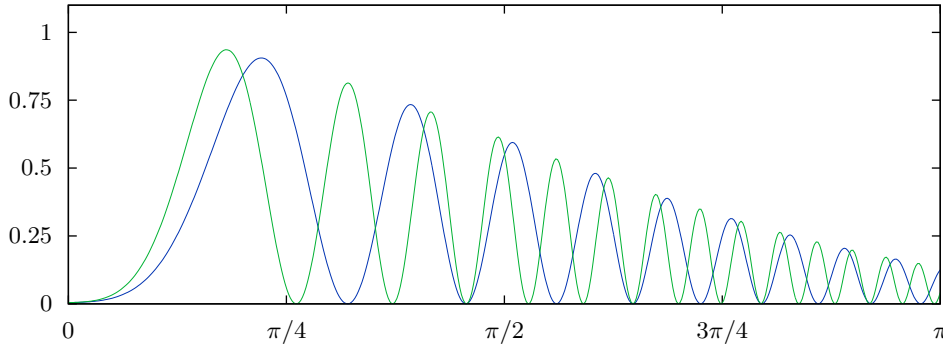
The kernel  $F_n$  is calculated using an NUFFT at a computational cost of  $\mathcal{O}(N^6 \log N + nN^4)$  by rewriting (57) as the six-dimensional non-uniform discrete Fourier transform of the  $N^2$ -by- $N^2$ -by- $n$  array formed by the outer products of  $\mathcal{F}h_s$  sampled on a two-dimensional  $M_{N-1}^2$  grid. Once this kernel is computed, applying  $L_n$  using the convolution formulation (56) costs  $\mathcal{O}(N^6 \log N)$ . This follows because the basis evaluations/expansions each have complexity  $\mathcal{O}(N^3 \log N)$  and there are  $\mathcal{O}(N^3)$  rows and columns in the matrix, while the six-dimensional convolution achieves complexity  $\mathcal{O}(N^6 \log N)$  using FFTs.

The right-hand side matrix  $B_n$  can also be computed as an NUFFT with complexity  $\mathcal{O}(N^6 \log N + nN^4)$ . As a result, the overall complexity for solving the normal equations (24) using CG is  $\mathcal{O}(\sqrt{\kappa(L_n)} N^6 \log N + nN^4)$ . We again note that this complexity is nearly optimal, since with respect to storing the  $\mathcal{O}(N^6)$  size covariance matrix, we only lose a condition number and logarithmic factor, while we only require a  $N^2$  factor increase with respect to storing the input images of size  $\mathcal{O}(nN^2)$ .

**5.3. Circulant preconditioners for  $A_n$  and  $L_n$ .** The number of iterations required for CG to converge scales with the square root of the condition number of the linear system [67, 5, 79]. As we shall see, both  $A_n$  and  $L_n$  are badly conditioned due to geometric considerations and the influence of the CTFs. One way to solve this is to increase the regularization parameters  $\nu_n$  and  $\xi_n$ . However, this increases the bias of the estimator and may not be desirable in all situations. Fortunately, the number of iterations in CG can be reduced without regularizing the original problem by introducing an appropriate preconditioner. In the following, we describe how this can be achieved for  $A_n$  and  $L_n$ .

For a uniform distribution of viewing directions  $R$  over  $\text{SO}(3)$  and with no microscope aberration, that is  $\mathcal{F}h(\omega) = 1$ , the unregularized (i.e.,  $\nu_n = 0$ ), continuous version of  $A_n$  is approximated in the Fourier domain by the filter  $2/\|\omega\|$  as  $n \rightarrow \infty$  (e.g., see [57, 38]). Qualitatively, this is the limit of the Fourier transform  $\mathcal{F}f_n$  of  $f_n$  as  $N, n \rightarrow \infty$ . As  $\mathcal{F}f_n$  approaches this limit, the  $1/\|\omega\|$  decay results in worse conditioning of  $A_n$ . The influence of the CTFs do not improve this situation. Indeed, the square magnitudes of the CTFs  $|\mathcal{F}h_s(\omega)|^2$  form a bandpass filter (see Figure 4), attenuating low and high frequencies, which is replicated in the Fourier transform of  $f_n$ , further worsening the conditioning of  $A_n$ .





**Figure 4.** The square magnitudes  $|\mathcal{F}h(\boldsymbol{\omega})|^2$  of two sample CTFs. Their sum forms a bandpass filter, worsening the conditioning of the least-squares estimators.

Similar results hold for  $L_n$ . Indeed, for uniform distribution of viewing directions and  $\mathcal{F}h(\boldsymbol{\omega}) = 1$  with no regularization (i.e.,  $\xi_n = 0$ ), it was shown by Katsevich et al. [38] that its continuous version acts as multiplication in the Fourier domain by  $2/\|\boldsymbol{\omega}_1 \times \boldsymbol{\omega}_2\|$  as  $n \rightarrow \infty$ . Not only does this decay as  $1/\|\boldsymbol{\omega}_1\|\|\boldsymbol{\omega}_2\|$ , but the kernel is singular for  $\boldsymbol{\omega}_1$  parallel to  $\boldsymbol{\omega}_2$ . Again, adding CTFs results in attenuation for low and high frequencies  $\|\boldsymbol{\omega}_1\|$  and  $\|\boldsymbol{\omega}_2\|$ . As a result, the Fourier transform  $\mathcal{F}F_n$  of the kernel  $F_n$  representing  $L_n$  is close to singular and decays rapidly, resulting in poor conditioning of  $L_n$ .

If the distribution of viewing directions is non-uniform, the condition number will be even larger. As a consequence of the above phenomena, a large number of iterations will be required in order to achieve convergence, which is not a desirable situation.

To remedy this, we precondition  $A_n$  and  $L_n$ . In other words, we define operators  $C_n$  and  $D_n$  that can be easily inverted and such that  $C_n^{-1}A_n$  and  $D_n^{-1}L_n$  are close to the identity. As a result,  $\kappa(C_n^{-1}A_n)$  and  $\kappa(D_n^{-1}L_n)$  are small, allowing us to reformulate the linear systems to achieve better conditioning. This is the idea of the preconditioned CG method, which converges in  $\mathcal{O}(\sqrt{\kappa(C_n^{-1}A_n)})$  and  $\mathcal{O}(\sqrt{\kappa(D_n^{-1}L_n)})$  steps for the mean and covariance estimators, respectively [67, 5, 79]. Note that the overall conditioning of the problem is still the same, it is only the CG convergence rate that changes.

A variety of preconditioners have been developed to improve convergence of the CG method. Popular alternatives include diagonal, or Jacobi, preconditioners and incomplete Cholesky or LU factorizations [67, 5]. As mentioned above, a good preconditioner is one whose inverse closely approximates the operator itself. This is commonly achieved by exploiting its structure. For example, a Toeplitz operator is well approximated by a circulant operator. Inverting a circulant operator is in turn achieved efficiently using FFTs [77].

In the current work, we therefore consider circulant approximations to the convolution factors of  $A_n$  and  $L_n$  as preconditioners. A circulant operator  $C$  operating on some  $d$ -dimensional vector  $w$  defined on  $M_N^d$  is given by a circular convolution

$$(58) \quad Cw(\mathbf{i}) = \sum_{\mathbf{j} \in M_N^d} w^{\text{per}}(\mathbf{i} - \mathbf{j})g(\mathbf{j}),$$

where  $w^{\text{per}}$  is the periodized extension of  $w$

$$(59) \quad w^{\text{per}}(\mathbf{i}) = w(\mathbf{i} + N\mathbf{m}) \quad \text{for } \mathbf{m} \text{ such that } \mathbf{i} + N\mathbf{m} \in M_N^d,$$

and  $g$  is a convolution kernel defined on  $M_N^d$ . The difference with the standard convolution encountered before in  $A_n$  and  $L_n$  (see (49) and (56)) is that the functions are periodized at the boundary instead of padded with zeros.

An important advantage of circulant operators is that they are diagonalized by the Fourier basis, and we can therefore write the action of  $C$  in the Fourier domain as

$$(60) \quad \mathcal{F}Cw(\mathbf{k}) = \mathcal{F}w(\mathbf{k})\mathcal{F}g(\mathbf{k}).$$

This makes circulant operators fast to apply, but also fast to invert, since

$$(61) \quad \mathcal{F}C^{-1}w(\mathbf{k}) = \mathcal{F}w(\mathbf{k})\mathcal{F}g(\mathbf{k})^{-1},$$

provided  $\mathcal{F}g(\mathbf{k}) \neq 0$  for all  $\mathbf{k} \in M_N^d$ . Since standard and circular convolutions differ principally at the boundary, they are similar when both  $w$  and  $g$  concentrate around the origin.

Circulant operators are therefore good preconditioners to standard convolutions, provided we can calculate them efficiently. Let  $C_n$  be a circulant operator with kernel  $g_n$  such that

$$(62) \quad \tilde{C}_n := \underset{\tilde{C}_n \text{ circ.}}{\operatorname{argmin}} \|\tilde{A}_n - \tilde{C}_n\|_{\text{F}},$$

and let

$$(63) \quad \tilde{A}_n = \frac{1}{n} \sum_{s=1}^n I_s^{\text{T}} I_s + \nu_n \mathbf{I}_{N^3},$$

be the voxel version of  $A_n$  such that  $\tilde{A}_n v = v * (f_n + \nu_n \delta_0)$ , where  $\delta_0$  is the three-dimensional Dirac delta function with value 1 at  $\mathbf{i} = 0$  and zero elsewhere. Such approximations have been previously studied by Tyrtshnikov [81], from whose results we derive the formula

$$(64) \quad g_n(\mathbf{i}) := \frac{1}{N^3} \sum_{\substack{\mathbf{m} \in M_{2N-1}^3 \\ \mathbf{m} = \mathbf{i} \pmod{N}}} (N - |m_1|)(N - |m_2|)(N - |m_3|) f_n(\mathbf{m}) + \nu_n \delta_0(\mathbf{i}),$$

for all  $\mathbf{i} \in M_N^3$ . This periodizes the original kernel  $f_n + \nu_n \delta_0$  with periodicity  $N$  and weights by a multiplier that attenuates points far from the origin.

This particular circulant approximation  $\tilde{C}_n$  of the Toeplitz operator  $\tilde{A}_n$  has the advantage of being computed with low computational complexity. Indeed, calculating  $g_n$  using (64) has complexity  $\mathcal{O}(N^3)$ . Furthermore, it minimizes the distance in Frobenius norm to  $\tilde{A}_n$  and preserves the positive semidefiniteness and invertibility of  $\tilde{A}_n$  [81], which is not true for other circulant preconditioners, such as those proposed by Strang and Chan [77, 10].

Since the circulant approximation  $\tilde{C}_n$  can be inverted easily using three-dimensional FFTs, we use  $C_n^{-1} := Q^{\text{T}} \tilde{C}_n^{-1} Q$  as a preconditioner in the CG method when solving  $A_n \mu_n = b_n$ . In this case, the effective condition number of the preconditioned linear system is equal to

---

**Algorithm 1** The least-squares mean estimator  $\mu_n$ 


---

**function** MEANESTIMATION( $\{R_s\}_{s=1}^n, \{h_s\}_{s=1}^n, \{y_s\}_{s=1}^n, Q, \nu_n$ )  
 Set  $f_n[\mathbf{i}] \leftarrow \frac{1}{nN^4} \sum_{s=1}^n \sum_{\mathbf{l} \in M_{2N-1}^2} |\mathcal{F}h_s(\mathbf{l})|^2 e^{\frac{2\pi i}{N} \langle \mathbf{i}, R_s^T \mathbf{l}; 0 \rangle}$   
 Set  $g_n[\mathbf{i}] \leftarrow \frac{1}{N^3} \sum_{\substack{\mathbf{m} \in M_{2N-1}^3 \\ \mathbf{m} = \mathbf{i} \pmod{N}}} (N - |m_1|)(N - |m_2|)(N - |m_3|) f_n(\mathbf{m}) + \nu_n \delta_0(\mathbf{i}) \quad \forall \mathbf{i} \in M_N^3$   
 Calculate  $\bar{g}_n$  such that  $\mathcal{F}\bar{g}_n = (\mathcal{F}g_n)^{-1}$   
 Set  $b_n \leftarrow \frac{1}{n} \sum_{s=1}^n P_s^T y_s$   
 Apply CG to  $Q^T(Q\mu_n * f_n) + \nu_n \mu_n = b_n$  with preconditioner  $x \mapsto Q^T(Qx * \bar{g}_n)$   
**return**  $\mu_n$   
**end function**

---

$\kappa(C_n^{-1}A_n)$ , which is small if the approximation is accurate. Numerical experiments in Section 7.4 indicate that this preconditioner brings the condition number down to 1–50.

The same type of circulant approximation can also be found for the convolution factor in  $L_n$ , with the circulant operator

$$(65) \quad \tilde{D}_n := \operatorname{argmin}_{\tilde{D}_n \text{ circ.}} \|\tilde{L}_n - \tilde{D}_n\|_{\mathbb{F}},$$

where

$$(66) \quad \tilde{L}_n(Z) = \frac{1}{n} \sum_{s=1}^n I_s^T I_s Z I_s^T I_s + \xi_n Z = Z * (F_n + \xi_n \delta_0)$$

is the voxel version of  $L_n$  where  $\delta_0$  is now the six-dimensional Dirac delta function. The kernel  $G_n$  of  $\tilde{D}_n$  can be found to equal

$$(67) \quad G_n(\mathbf{i}) := \frac{1}{N^6} \sum_{\substack{\mathbf{m} \in M_{2N-1}^6 \\ \mathbf{m} = \mathbf{i} \pmod{N}}} (N - |m_1|) \cdots (N - |m_6|) F_n(\mathbf{m}) + \xi_n \delta_0(\mathbf{i}),$$

for all  $\mathbf{i} \in M_N^6$ . Again,  $G_n$  is a weighted periodization of  $F_n + \xi_n \delta_0$ . The computational complexity of calculating  $G_n$  is  $\mathcal{O}(N^6)$ .

The circulant operator  $\tilde{D}_n$  can also be inverted quickly using six-dimensional FFTs, so we use  $\Sigma \mapsto D_n^{-1}(\Sigma) := Q^T \tilde{D}_n^{-1}(Q \Sigma Q^T) Q$  to precondition the normal equations  $L_n(\Sigma_n) = B_n$  of the least-squares covariance estimator or  $L_n(\Sigma_n) = B_n^{(s)}$  for the shrinkage variant. For the same reasons as in the mean estimation case, the condition number  $\kappa(D_n^{-1}L_n)$  is expected to be small. Again, numerical simulations in Section 7.4 indicate that the condition number of this operator stays in the regime 1–200.

**5.4. Conjugate gradient & thresholding.** We are now ready to formulate the algorithms for estimating  $\mu_n$  and  $\Sigma_n$  given the input images  $y_1, \dots, y_n$  and projection mappings  $P_1, \dots, P_n$  (or equivalently, the viewing directions  $R_1, \dots, R_n$  and CTFs  $\mathcal{F}h_1, \dots, \mathcal{F}h_n$ ).

The mean estimation algorithm is given in Algorithm 1. First, the convolutional kernel  $f_n$  associated with  $A_n$  and the right-hand side  $b_n$  are computed at a cost of  $\mathcal{O}(N^3 \log N + nN^2)$ .

---

**Algorithm 2** The covariance estimators  $\Sigma_n$  ( $do\_shrink = false$ ) and  $\Sigma_n^{(s)}$  ( $do\_shrink = true$ )

---

**function** COVARIANCEESTIMATION( $\{R_s\}_{s=1}^n, \{h_s\}_{s=1}^n, \{y_s\}_{s=1}^n, \mu_n, \sigma, Q, \xi_n, do\_shrink$ )  
 Set  $f_n[\mathbf{i}_1, \mathbf{i}_2] \leftarrow \frac{1}{nN^8} \sum_{s=1}^n \sum_{\mathbf{l}_1, \mathbf{l}_2 \in M_{N-1}^2} |\mathcal{F}h_s(\mathbf{l}_1)|^2 |\mathcal{F}h_s(\mathbf{l}_2)|^2 e^{\frac{2\pi i}{N} (\langle \mathbf{i}_1, R_s^T[\mathbf{l}_1; 0] \rangle - \langle \mathbf{i}_2, R_s^T[\mathbf{l}_2; 0] \rangle)}$   
 Set  $G_n[\mathbf{i}] \leftarrow \frac{1}{N^6} \sum_{\substack{\mathbf{m} \in M_{2N-1}^6 \\ \mathbf{m} = \mathbf{i} \pmod{N}}} (N - |m_1|) \cdots (N - |m_6|) F_n(\mathbf{m}) + \xi_n \delta_0(\mathbf{i}) \quad \forall \mathbf{i} \in M_N^6$   
 Calculate  $\bar{G}_n$  such that  $\mathcal{F}\bar{G}_n = (\mathcal{F}G_n)^{-1}$   
 Set  $B_n \leftarrow \frac{1}{n} \sum_{s=1}^n P_s^T (y_s - P_s \mu_n) (y_s - P_s \mu_n)^T P_s - \sigma^2 P_s^T P_s$   
**if**  $do\_shrink$  **then**  
   Set  $B_n \leftarrow B_n^{(s)} = A_n^{1/2} \rho(A_n^{-1/2} B_n A_n^{-1/2} + \sigma^2 I_p) A_n^{1/2}$   
**end if**  
 Apply CG to  $Q^T(Q\Sigma_n Q^T * F_n) + \xi_n \Sigma_n = B_n$  with preconditioner  $X \mapsto Q^T(QXQ^T * \bar{G}_n)Q^T$   
**return**  $\Sigma_n$   
**end function**

---

The circulant approximation kernel  $g_n$  is then calculated from  $f_n$ , which takes  $\mathcal{O}(N^3)$ . We then apply CG to (16), with each iteration requiring application of  $A_n$  and  $C_n^{-1}$  which are obtained by multiplications by  $Q$ ,  $Q^T$  and convolutions by  $f_n$  and  $\bar{g}_n$ , all of which have computational complexity of  $\mathcal{O}(N^3 \log N)$ . After  $\mathcal{O}(\sqrt{\kappa(C_n^{-1} A_n)})$  iterations, we have  $\mu_n$ . The overall computational complexity of Algorithm 1 is then

$$(68) \quad \mathcal{O}\left(\sqrt{\kappa(C_n^{-1} A_n)} N^3 \log N + nN^2\right),$$

where  $\kappa(C_n^{-1} A_n)$  is typically in the range 1–50.

The covariance estimation method listed in Algorithm 2 is qualitatively similar. The convolutional kernel  $F_n$  associated with  $L_n$  and the right-hand side matrix  $B_n$  are computed with complexity  $\mathcal{O}(N^6 \log N + nN^4)$ . The circulant kernel  $G_n$  is obtained at cost  $\mathcal{O}(N^6)$ . Applying  $L_n$  and  $D_n$  now involves multiplication by  $Q$  and  $Q^T$  as well as convolution with  $F_n$  and  $G_n^{-1}$ , each of which has computational complexity  $\mathcal{O}(N^6 \log N)$ . Now  $\mathcal{O}(\sqrt{\kappa(D_n^{-1} L_n)})$  iterations are needed to obtain  $\Sigma_n$ . The overall complexity of Algorithm 2 is then

$$(69) \quad \mathcal{O}\left(\sqrt{\kappa(D_n^{-1} L_n)} N^6 \log N + nN^4\right),$$

where  $\kappa(D_n^{-1} L_n)$  is in the range 1–200.

To obtain the shrinkage variant of the estimator,  $\Sigma_n^{(s)}$ , the additional step of calculating  $B_n^{(s)}$  from  $B_n$  is added before the CG step in Algorithm 2. This is done using (46), where  $B_n$  is whitened by conjugation with  $A_n^{-1/2}$ , the whitened matrix is shrunk using  $\rho$ , and the result is unwhitened by conjugation with  $A_n^{1/2}$ . The number of top eigenvalues of  $A_n^{-1/2} B_n A_n^{-1/2}$  which exceed the Marčenko-Pastur threshold is typically small, so we can exploit Lanczos method for finding the top eigenvalues and eigenvectors. For this, we need to apply  $A_n^{-1/2} B_n A_n^{-1/2}$  fast. Since  $A_n$  can be applied fast using its convolution kernel, applying its inverse square root  $A_n^{-1/2}$

to a volume can be approximated iteratively using Krylov subspace methods [20, 66, 33]. The number of iterations needed is typically small, so we take its complexity to be  $\mathcal{O}(N^3 \log N)$ . As a result, applying  $A_n^{-1/2} B_n A_n^{-1/2}$  has complexity  $\mathcal{O}(N^6)$ , since matrix multiplication by  $B_n$  takes  $\mathcal{O}(N^6)$ . The overall eigendecomposition calculation therefore has complexity  $\mathcal{O}(N^6)$ , where we have assumed that the number of non-trivial eigenvectors is taken to be constant.

We note that an alternative to Krylov subspace methods for approximating  $A_n^{-1/2}$  is to exploit the Toeplitz structure in  $A_n$  and use this to calculate its Cholesky factors, which have similar properties to the matrix square root but can be inverted efficiently. In one dimension, this can be done in  $\mathcal{O}(N^2)$  using the Schur algorithm [70, 56, 1]. While this has been generalized to matrices of block Toeplitz structure [1], these do not take into account 2-Toeplitz structure, also known as Toeplitz-block-Toeplitz, and so have complexity  $\mathcal{O}(N^5)$  instead of the desired  $\mathcal{O}(N^4)$ . Designing an appropriate generalization of the Schur algorithm for  $d$ -Toeplitz operators where  $d \geq 2$  is the subject of future work.

Both in the case of the standard estimator and the shrinkage variant, the estimated covariance matrix  $\Sigma_n$  will contain a considerable amount of error in the form of a bulk noise distribution similar to that observed in the spiked covariance model. A final step of selecting the dominant eigenvectors is therefore necessary to extract the relevant part of the covariance matrix structure. Since we expect the population covariance matrix to be of low rank, it must have a small number of non-zero eigenvectors. This number can be estimated by looking for a ‘‘knee’’ in the spectrum of  $\Sigma_n$  where the dominant eigenvalues separate from the bulk noise distribution. In the case of a discrete distribution of molecular structures, this is at most one minus the number of resolved structures in the dataset. However, this determination has to be done manually. A heuristic method for validating this choice could be to inspect the corresponding eigenvectors and determine how ‘‘noise-like’’ they appear, using a suitable prior. Future work will focus on enabling the algorithm to perform this selection automatically. The computational complexity of calculating the leading  $r = \mathcal{O}(1)$  eigenvectors of  $\Sigma_n$  is  $\mathcal{O}(rN^6)$ .

An important feature of (69) is that the algorithm scales as  $N^6$  in the resolution  $N$  of the images. Since we are estimating the entire covariance matrix  $\text{Cov}[x]$ , this is unavoidable since that matrix has  $N^6$  entries. However, it has the unfortunate consequence of limiting the attainable resolution of covariance estimation using the proposed algorithm. For example, at  $N = 16$ , the covariance estimate  $\Sigma_n$  requires 128 MB to store in double precision. The kernel  $f_n$  is of size  $2N$  and therefore requires  $2^6 = 64$  times as much space, or 8 GB. Increasing  $N$  beyond this becomes impractical for a typical workstation.

That being said, a large amount of useful information can be obtained at these resolutions. Indeed, since our goal is to classify rather than reconstruct, all we need is for the features that discriminate between various conformations to be present at low resolution. This is not an unreasonable assumption. Indeed, if one subunit of a molecule moves with respect to another, we can capture that movement at low resolution as long as that subunit is large enough. Similarly, the binding of external complexes to larger molecule is visible provided that those complexes are large enough. Once we can distinguish such differences, the dataset can be partitioned and higher-resolution reconstructions can then be produced from each subset. We have observed this in experimental data, suggesting that the restriction to  $N = 16$  is not as debilitating as it may first seem to be. In our experiments in this paper, we shall therefore

restrict ourselves to  $N = 16$ .

**5.5. Choice of basis.** To represent a volume  $x$ , we can store its values on the voxel grid  $M_N^3/N$ . We will call this a decomposition in the voxel basis. The problem with this basis is that the electric potential of a molecule is supported in the central ball  $\{\|\mathbf{u}\| < 1\}$ , with no energy in the ‘‘corners’’ of the cube  $[-1, +1]^3$ . Indeed, any energy in this region will be captured by projections along a subset of viewing angles and will not be well reconstructed. We can therefore safely assume that the support is contained in the central ball. The same holds in the frequency domain, where frequency samples outside the Nyquist ball  $\{\|\mathbf{k}\| < N/4\}$  are expected to be negligible. In addition, the low sampling density of these frequencies leads to ill-conditioning of the reconstruction problem, which we would like to avoid.

To solve this, we will use different bases which are concentrated on  $\{\|\mathbf{u}\| < 1\}$  in space and within  $\{\|\mathbf{k}\| < N/4\}$  in frequency. One solution to this spectral concentration is given by the 3D Slepian functions [76], but their implementation is quite complicated. Instead, we will focus on an alternative basis with similar properties, the spherical Fourier-Bessel basis. It consists of functions given in spherical coordinates  $(r, \theta, \phi)$  by

$$(70) \quad \phi_{\ell,k,m}(r, \theta, \phi) = \begin{cases} C_{\ell,k} j_\ell(r z_{\ell,k}) Y_{\ell,m}(\theta, \phi) & 0 \leq r < 1 \\ 0 & 1 \leq r \end{cases}$$

where  $j_\ell$  is the spherical Bessel function of order  $\ell$ ,  $z_{\ell,k}$  is the  $k$ th zero of  $j_\ell$ , and  $Y_{\ell,m}$  is the spherical harmonic function of order  $\ell$  and degree  $m$ , and  $C_{\ell,k} = \sqrt{2} |j_{\ell+1}(z_{\ell,k})|^{-1}$ . The indices  $m$  and  $k$  satisfy  $|m| \leq \ell$  and  $k \leq k_{\max}(\ell)$ , where  $k_{\max}(\ell)$  is the largest integer such that  $z_{\ell, k_{\max}(\ell)+1} < N\pi/4$ . This is the same sampling criterion used in Bhamre et al. [8] and Cheng [11], which generalizes similar criteria for the 2D Fourier-Bessel basis [41, 92]. This condition on  $k$  ensures that  $\mathcal{F}\phi$  is concentrated within the Nyquist ball since this function is concentrated around a ring centered at  $\|\mathbf{k}\| = z_{\ell,k}/\pi$ . Finally, the constant  $C_{\ell,k}$  ensures that the basis functions have unit norm. For  $\ell$  up to some  $\ell_{\max}$ , we therefore have the basis  $\{\phi_{\ell,k,m}\}_{\ell \leq \ell_{\max}, k \leq k_{\max}(\ell), |m| \leq \ell}$  which we use to decompose  $x$ .

However, as discussed in Section 2, the standard voxel basis allows for fast projection through using NUFFTs. To take advantage of this, we need a fast change-of-basis mapping between the voxel basis and the spherical Fourier-Bessel basis. For this, we can use NUFFTs and separation of variables to evaluate the basis at voxel grid points in  $\mathcal{O}(N^4)$  complexity. Using fast spherical harmonic transforms [80, 43] and fast Fourier-Bessel transforms [53], we can reduce this further to  $\mathcal{O}(N^3 \log N)$ .

A simpler alternative is provided by the truncated Fourier basis

$$(71) \quad \phi_{\mathbf{k}}(x) = \begin{cases} C e^{2\pi i \langle x, \mathbf{k} \rangle} & x \in M_{N-2}^3/N \\ 0 & \text{otherwise} \end{cases}$$

for  $\mathbf{k} \in M_{N-2}^3 \cap \{\|\boldsymbol{\omega}\| < (N-2)/2\}$ . Again,  $C$  is chosen so that  $\phi_{\mathbf{k}}$  has unit norm for all  $\mathbf{k}$ . The functions are zero outside a central box  $M_{N-2}$ , providing a padding of one voxel in each direction, and only consists of frequencies inside the Nyquist ball. While less concentrated compared to the spherical Fourier-Bessel basis, it has the advantage of providing efficient change-of-basis mappings through standard FFTs.



In the following, we will use the spherical Fourier-Bessel basis since it enjoys better concentration properties. We note, however, that for large values of  $N$ , it may be more computationally efficient to use the truncated Fourier basis since the constant associated with FFTs is much smaller than those of the fast spherical harmonic and Fourier-Bessel transforms.

**6. Reconstruction of states.** Having estimates  $\mu_n$  and  $\Sigma_n$  of the mean  $\mathbb{E}[x]$  and covariance  $\text{Cov}[x]$  provides us with partial information on the distribution of the volumes  $x_1, \dots, x_n$ . However, unless the distribution of  $x$  is Gaussian, this is not enough to fully characterize it. To do so, more information has to be extracted from the images  $y_1, \dots, y_n$ . We shall consider two types of singular distributions: those supported on a finite number of points and those supported on a continuous low-dimensional manifold.

**6.1. Wiener filter.** For a fixed viewing direction  $R$ , the variability in the random density  $x$  encoded by the 3D covariance  $\text{Cov}[x]$  induces variability in the clean images  $Px$  through the 2D covariance  $PCov[x]P^T$ . Classical Wiener filtering leverages this covariance to denoise images or estimate the underlying volume corresponding to each image.

Recall that we have the image formation model (12), restated here as

$$(72) \quad y_s = P_s x_s + e_s,$$

where, as before,  $P_s$  defines projection along viewing direction  $R_s$  and convolution with  $h_s$ . As we saw earlier, this induces the relation

$$(73) \quad \text{Cov}_{x_s|P_s}[P_s x_s] = P_s \text{Cov}[x] P_s^T,$$

for the signal term, which allows us to estimate the image covariance as

$$(74) \quad \text{Cov}_{x_s|P_s}[P_s x_s] \approx P_s \Sigma_n P_s^T.$$

The noise covariance  $\text{Cov}[e_s]$  is assumed to be  $\sigma^2 \mathbf{I}_{N^2}$ .

We now use the estimated mean and covariance to define a Wiener filter estimator [51]

$$(75) \quad \hat{x}_s := H_s (y_s - P_s \mu_n) + \mu_n,$$

of  $x_s$  where

$$(76) \quad H_s := \Sigma_{n,r} P_s^T (P_s \Sigma_{n,r} P_s^T + \sigma^2 \mathbf{I}_{N^2})^{-1}.$$

If  $\Sigma_{n,r} = \text{Cov}[x]$  and  $\mu_n = \mathbb{E}[x]$ , this linear filter minimizes the expected mean-squared error

$$(77) \quad \mathbb{E}_{x_s, e_s|P_s} \|\hat{x}_s - x_s\|^2.$$

In the finite-sample case, this no longer holds, but we can expect the Wiener filter to perform better than the pseudo-inverse for reasonably accurate mean and covariance estimates.

As discussed in the previous section,  $\Sigma_n$  is often of low rank  $r$  following the final thresholding step, giving  $\Sigma_{n,r}$ . We can use this to significantly reduce the complexity of calculating  $\hat{x}_s$ , which through direct evaluation of (76) takes  $\mathcal{O}(N^6 \log N)$  operations since it involves calculating  $P_s Q \Sigma_n Q^T P_s^T$ . Let

$$(78) \quad O_s U_s = P_s V_{n,r}$$

be the “thin” QR decomposition of  $P_s V_{n,r}$ , where  $O_s$  is an  $N^2$ -by- $r$  orthonormal matrix and  $U_s$  is an  $r$ -by- $r$  upper triangular matrix [79, 26]. Using these matrices, we rewrite  $\hat{x}_s$  as

$$(79) \quad \hat{x}_s = V_{n,r} \Lambda_{n,r} U_s^T (U_s \Lambda_{n,r} U_s^T + \sigma^2 \mathbf{I}_r)^{-1} O_s^T (y_s - P_s \mu_n) + \mu_n,$$

for  $s = 1, \dots, n$ . This shows that  $\hat{x}_s$  equals  $\mu_n$  plus a linear combination of the vectors  $V_{n,r}$ . A more compact, but isometrically equivalent, representation is therefore given by

$$(80) \quad \hat{\alpha}_s := V_{n,r}^T (\hat{x}_s - \mu_n) = \Lambda_{n,r} U_s^T (U_s \Lambda_{n,r} U_s^T + \sigma^2 \mathbf{I}_r)^{-1} O_s^T (y_s - P_s \mu_n).$$

These coordinates can be calculated in  $\mathcal{O}(rN^3 \log N + nr^2N^2)$ , since the images  $P_s V_{n,r}$  and  $P_s \mu_n$  require  $\mathcal{O}(rN^3 \log N + nrN^2)$ , the QR decompositions have computational complexity  $\mathcal{O}(nr^2N^2)$  and inverting  $n$   $r$ -by- $r$  matrices takes  $\mathcal{O}(nr^3)$ , where we assume that  $N^2 \gg r$ . We note that this is close to optimal, since  $r$  eigenvolumes require  $\mathcal{O}(rN^3)$  in storage, while the images are stored in  $\mathcal{O}(nN^2)$ , so we only lose a factor of  $\log N$  and  $r^2$ , respectively.

The Wiener filter estimate of the volumes is now

$$(81) \quad \hat{x}_s = V_{n,r} \hat{\alpha}_s + \mu_n.$$

The traditional denoising Wiener filter of the 2D images is obtained by projecting these volume estimates. Specifically, we define

$$(82) \quad \hat{y}_s := P_s \hat{x}_s = P_s (V_{n,r} \hat{\alpha}_s + \mu_n).$$

This is the same estimator obtained by minimizing the expected loss  $\mathbb{E}_{x_s, e_s | P_s} [\|\hat{y}_s - y_s\|^2]$  of a linear estimator  $\hat{y}_s$  and substituting our estimates for the volume mean and covariance.

**6.2. Volume distance measures.** Given the images  $y_1, \dots, y_n$  together with our mean and covariance estimates  $\mu_n$  and  $\Sigma_n$ , we can also define distance measures on the underlying volumes. This will allow us to cluster them using methods described in Section 6.3 or to describe their manifold structure using the manifold learning techniques in Section 6.4.

The simplest distance is the Euclidean norm on the volume estimates  $\hat{x}_1, \dots, \hat{x}_n$  given by

$$(83) \quad d_{st}^{(\text{eucl})} = \|\hat{x}_s - \hat{x}_t\|,$$

for  $s, t = 1, \dots, n$ .

Unfortunately, this distance measure weights all directions equally regardless of their accuracy for a given pair. To see why this is a problem, consider the fact that the columns of  $P_s V_{n,r}$  have different norms which depend on  $s$ . For example, a volume which is highly oscillatory along one axis will project to almost zero for viewing directions along that axis. Since these vectors are used to estimate  $\hat{x}_s$ , this means that the power of the noise is different for each coordinate. A distance measure that takes this into account would therefore be more robust than  $d_{st}^{(\text{eucl})}$ .

One way to do this is to instead consider distances on the denoised images  $\hat{y}_1, \dots, \hat{y}_n$ . While we still have the problem of low-energy basis vectors, these do not have a large energy once reprojected, so the situation is better. We then use the common-lines distance between the images [71]. From the Fourier Slice Theorem (6), we see that the Fourier transforms of

two images  $y_s$  and  $y_t$  occupy the planes orthogonal to  $R_s^{(3)}$  and  $R_t^{(3)}$ , respectively, where  $R_s^{(i)}$  denotes the  $i$ th row of  $R_s$ . As such, they intersect along the line defined by the unit vector

$$(84) \quad \frac{R_s^{(3)} \times R_t^{(3)}}{\|R_s^{(3)} \times R_t^{(3)}\|}.$$

We therefore define the common-lines vector  $\mathbf{c}_{st} \in \mathbb{R}^2$  for image  $s$  with respect to image  $t$  by rotating this vector into the image coordinates

$$(85) \quad \mathbf{c}_{st} = \begin{bmatrix} R_s^{(1)}; R_s^{(2)} \end{bmatrix} \left( \frac{R_s^{(3)} \times R_t^{(3)}}{\|R_s^{(3)} \times R_t^{(3)}\|} \right).$$

The common lines of  $\hat{y}_s$  and  $\hat{y}_t$  with respect to one another are then  $\mathcal{F}\hat{y}_s(k\mathbf{c}_{st})$  and  $\mathcal{F}\hat{y}_t(k\mathbf{c}_{ts})$ , respectively, where  $k \in M_N$ . If  $y_s$  and  $y_t$  are projections of the same molecular structure, that is  $x_s = x_t$ , we expect that these common lines should be close since they are restrictions of the same volume Fourier transform along the same line. A useful distance is therefore

$$(86) \quad d_{st}^{(\text{cl})} = \sum_{k \in M_N} |\mathcal{F}\hat{y}_s(k\mathbf{c}_{st}) - \mathcal{F}\hat{y}_t(k\mathbf{c}_{ts})|^2,$$

for  $s, t = 1, \dots, n$ , which we call the common-lines distance. Note that this does not take into account the different CTFs of  $\hat{y}_s$  and  $\hat{y}_t$ , which puts it at a disadvantage compared to the Euclidean distance  $d_{st}^{(\text{eucl})}$ .

**6.3. Clustering.** In the previous sections, we estimated the 3D covariance matrix and used it to calculate estimates  $\hat{x}_1, \dots, \hat{x}_n$  of the volumes  $x_1, \dots, x_n$ , or more specifically, their low-dimensional coordinate vectors  $\hat{\alpha}_1, \dots, \hat{\alpha}_n$ . These were then used to define distances  $d^{(\text{eucl})}$  and  $d^{(\text{cl})}$  between the volumes. Without any additional assumptions, we cannot extract more information. For this, we need prior information on the distribution of the volume  $x$ .

For example, if  $x$  is a discrete random variable, we can fit a discrete distribution by clustering the volume estimates  $\hat{x}_1, \dots, \hat{x}_n$ . Such a model is realistic for many molecules, where the majority of their time is spent in a small number of states.

Instead of clustering the volume estimates  $\hat{x}_1, \dots, \hat{x}_n$  themselves, we work on the coordinates  $\hat{\alpha}_1, \dots, \hat{\alpha}_n$ , since these are lower-dimensional but isometrically parametrize the volumes. One clustering approach is the  $k$ -means vector quantization algorithm [50]. While widely used for clustering,  $k$ -means has several problems, one of which is that it favors partitions that distribute points uniformly between clusters. This can be partly mitigated by modeling  $x$  using a Gaussian mixture model (GMM) and fitting its parameters using the expectation-maximization algorithm [17].

Given the distances  $d^{(\text{eucl})}$  or  $d^{(\text{cl})}$  instead of coordinates, we can use standard graph clustering algorithms such as normalized cut [72]. These algorithms partition the points into subsets that optimize certain criteria, the goal being to minimize the distances within clusters while maximizing distances between clusters.

Regardless of the clustering mechanism, we obtain a cluster assignment associated with each image. We then average the corresponding volume estimates to obtain a reconstruction

of that class. However, as we shall see, the algorithm will typically be applied to downsampled images, so this reconstruction is by necessity of low resolution. A more accurate reconstruction is obtained by partitioning the dataset according to the cluster assignments and reconstructing each subset separately at full resolution using tools such as RELION [68], cryoSPARC [64], FREALIGN [28], or ASPIRE [93, 73, 75].

We note that performing full-resolution reconstruction for each subset provides refined estimates of the viewing directions  $R_1, \dots, R_n$ . While these are assumed given for our algorithm, they are not necessarily very accurate, since they must be estimated from the average molecular structure as discussed in Section 2. Since these subsets of particles given by clustering should be more homogeneous, we expect the estimates to be more accurate. The covariance estimation and clustering steps can then be repeated using the refined estimates to achieve better results. This approach is known as iterative refinement, and has proved useful for other cryo-EM problems [82, 24].

**6.4. Continuous variability and diffusion maps.** Certain molecules do not primarily exist in a discrete set of states, but exhibit continuous variability. In this case, the clustering approach outlined above fails. However, due to physical constraints on the molecular dynamics, this continuum of states can often be described by a small number of dominant flexible motions. In this section, we describe a method for analyzing this low-dimensional manifold using diffusion maps [13].

These tools have previously been applied to study the continuous variability of molecular structure by calculating diffusion maps for images in each viewing direction and “patching” these together to yield a diffusion map for the whole volume [15]. Unfortunately, this is a heuristic method which is not guaranteed to yield an accurate description of the low-dimensional conformation manifold. We propose to use the mean and covariance estimates, together with the derived distance measures  $d^{(\text{eucl})}$  or  $d^{(\text{cl})}$  to calculate a diffusion map embedding that more closely captures the underlying structural variability.

The first step is to form a similarity matrix  $W$  whose entries are

$$(87) \quad W_{st} := \exp\left(-\frac{d_{st}^2}{\epsilon}\right),$$

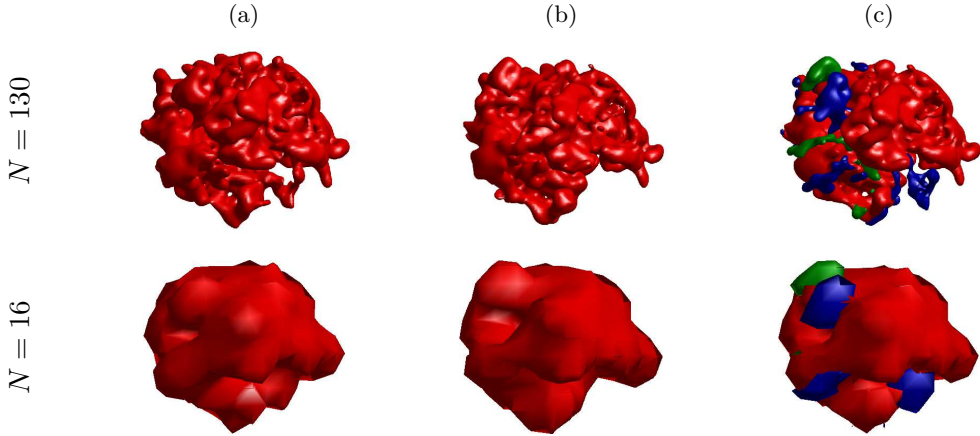
for  $s, t = 1, \dots, n$ , where  $d_{st}$  is a distance between volumes  $s$  and  $t$  while  $\epsilon$  is a scale parameter that depends on the smoothness of the manifold. The distance  $d_{st}$  can be one of  $d^{(\text{eucl})}$  or  $d^{(\text{cl})}$ . In the following, we shall use  $d^{(\text{eucl})}$ .

We now sum the entries along each row and define the diagonal matrix  $D$  with entries

$$(88) \quad D_{ss} := \sum_{t=1}^n W_{st},$$

for  $s = 1, \dots, n$ . Renormalizing  $W$  by  $D$ , we obtain the row-stochastic Markov transition matrix  $A := D^{-1}W$ . This defines a random walk on the graph of volume estimates, with the transition probability between two points proportional to their similarity. Let us calculate the eigenvalues and eigenvectors of  $A$ . We then have

$$(89) \quad A\phi_i = \lambda_i\phi_i,$$



**Figure 5.** The simulation ground truth at  $N = 130$  (top) and  $N = 16$  (bottom). (a, b) Two conformations of the 70S ribosome. (c) Their mean volume (red) and difference map (positive in blue, negative in green).

for  $i = 1, \dots, n$ . These eigenvalues satisfy  $|\lambda_i| \leq 1$  for all  $i$  and there is at least one eigenvalue equal to 1 with the corresponding eigenvector parallel to the all-ones vector (this follows from the fact that  $A$  is row-stochastic). We thus order the eigenvalues as  $1 = \lambda_1 \geq |\lambda_2| \geq \dots \geq |\lambda_n|$ .

The diffusion coordinates at diffusion time  $\tau$  for the  $s$ th volume are now

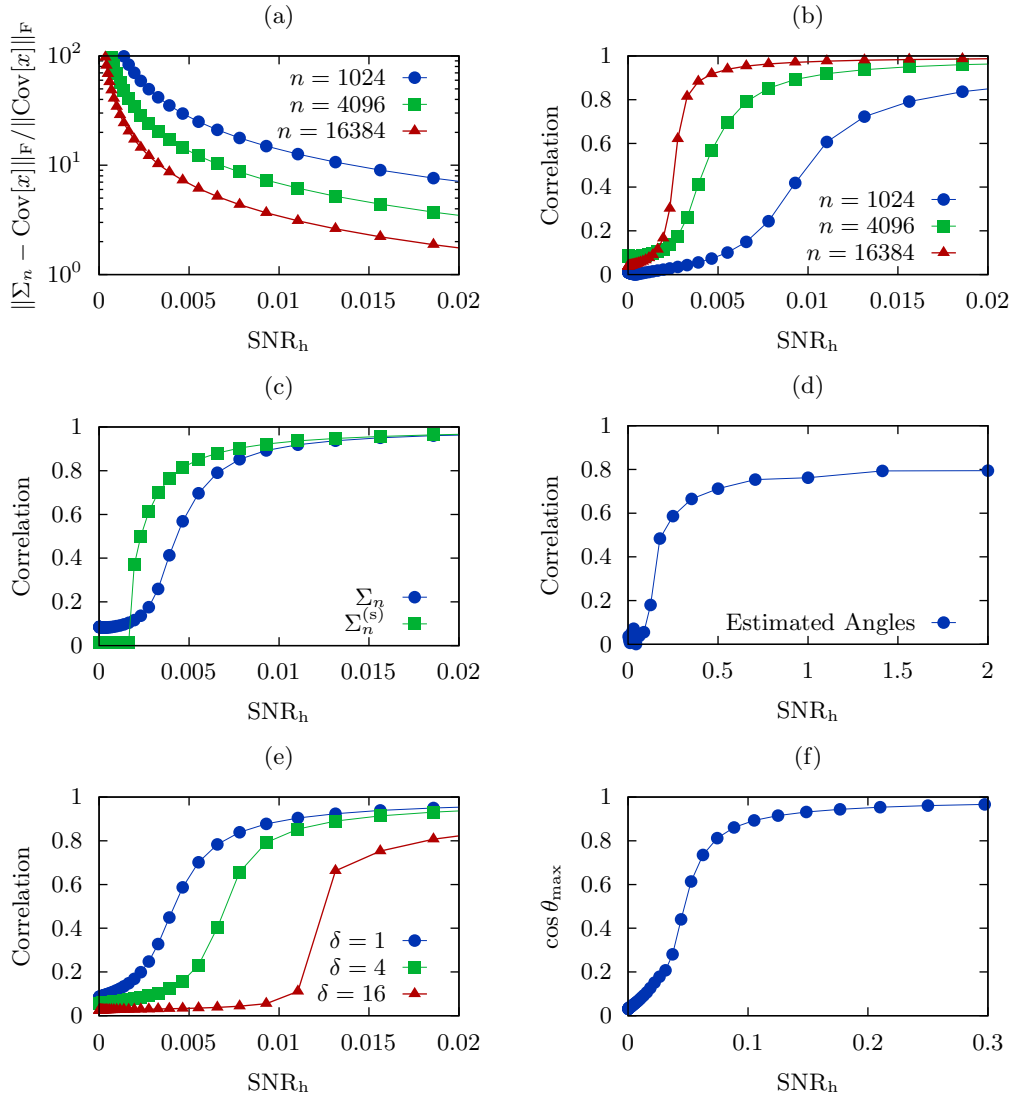
$$(90) \quad \widehat{\beta}_s^{(\tau)} := [\lambda_2^\tau \phi_{2,s}, \dots, \lambda_n^\tau \phi_{n,s}],$$

where  $\phi_{i,s}$  is the  $s$ th element of the  $i$ th eigenvector  $\phi_i$  [13]. The diffusion time  $\tau$  specifies the scale of the diffusion embedding. Specifically, the distance  $\|\widehat{\beta}_s^{(\tau)} - \widehat{\beta}_t^{(\tau)}\|$  between two diffusion map coordinates  $\widehat{\beta}_s^{(\tau)}$  and  $\widehat{\beta}_t^{(\tau)}$  approximates the distance between the probability distributions obtained from random walks on the volume graph starting at  $\widehat{x}_s$  and  $\widehat{x}_t$  after  $\tau$  steps. As  $\tau$  increases, these distributions start to overlap for points close together on the manifold. By increasing  $\tau$ , we obtain a set of coordinates  $\widehat{\beta}_1^{(\tau)}, \dots, \widehat{\beta}_n^{(\tau)}$  that are more robust to noise compared to  $\widehat{\alpha}_1, \dots, \widehat{\alpha}_n$  at the expense of smoothing out the fine-scale manifold structure.

Restricting  $\widehat{\beta}_s$  to the first two or three coordinates, we obtain a two- or three-dimensional embedding of the estimated volumes, which provides a helpful visualization to determine the global geometric structure of the continuous manifold of conformations. We shall see some examples of this in Section 7.3.

**7. Simulation results.** We evaluate the performance of the covariance estimation algorithm using simulated data. These are obtained by applying the forward model (11) to different configurations of volumes, projection mappings, and noise sources. The resulting images are then given as input to the covariance estimation algorithms outlined in Sections 4–6 to study their computational efficiency and accuracy.

**7.1. Covariance estimation results.** We first evaluate the performance of the covariance matrix in the absence of noise. We create a synthetic dataset from two 70S ribosome maps downsampled to  $N = 16$  shown in the bottom row of Figure 5. These are projected in viewing directions drawn from a uniform distribution over  $\text{SO}(3)$  and convolved with one of seven



**Figure 6.** Covariance estimation results for different simulations with  $N = 16$ ,  $C = 2$ , seven distinct CTFs, and unless otherwise noted,  $n = 4096$  and uniform distribution of orientations. (a) The relative error in  $\Sigma_n$  as a function of  $\text{SNR}_h$  for  $n = 1024$ ,  $n = 4096$ , and  $n = 16384$ . (b) The correlation of the top eigenvector of  $\Sigma_n$  with that of  $\text{Cov}[x]$ . (c) Top eigenvector correlations for  $\Sigma_n$  and  $\Sigma_n^{(s)}$ . (d) Top eigenvector correlations for  $\Sigma_n$  with orientations estimated using the ASPIRE toolbox. (e) Top eigenvector correlations for  $\Sigma_n$  with different orientation distributions over  $\text{SO}(3)$  described by (92). (f) The cosine of the maximum principal angle between the top three eigenvectors of  $\Sigma_n$  and those of  $\text{Cov}[x]$  for a simulation with  $C = 4$  classes.

distinct CTFs. From this we obtain  $n = 1024$  simulated images of size  $N = 16$ . Applying the mean estimation algorithm (Algorithm 1) followed by the covariance estimation algorithm (Algorithm 2), we obtain a covariance matrix estimate  $\Sigma_n$ . Here and throughout this section,



we use the unregularized variants of  $\mu_n$  and  $\Sigma_n$ , setting  $\nu_n = 0$  and  $\xi_n = 0$ .

The relative error of  $\Sigma_n$  compared to  $\text{Cov}[x]$  is  $\|\Sigma_n - \text{Cov}[x]\|_F / \|\text{Cov}[x]\|_F \approx 4.4 \cdot 10^{-2}$ . Note that the error is not zero due to the finite size of the dataset. Since we have two configurations,  $C = 2$  and  $\text{Cov}[x]$  has rank one. Consequently, we are only interested in the top eigenvector of  $\Sigma_n$ . If it is well-correlated with the single non-trivial eigenvector  $\text{Cov}[x]$ , this is another important performance measure. In this case, that correlation is  $1 - 8 \cdot 10^{-5}$ . If we redo the experiment for  $n = 16384$ , the relative error of  $\Sigma_n$  drops to  $2.3 \cdot 10^{-3}$  while the correlation of the top eigenvector increases to  $1 - 2 \cdot 10^{-7}$ . For clean data, the proposed method accurately estimates the covariance matrix as  $n$  increases.

We now add noise to the above simulation and consider the performance of our algorithm with respect to the signal-to-noise ratio (SNR) of the images. Since the task is to extract the heterogeneity structure of the data, the standard SNR comparing the average signal power to that of the noise is insufficient. We instead consider the heterogeneous signal-to-noise ratio

$$(91) \quad \text{SNR}_h = \frac{\sum_{s=1}^n \|P_s(x_s - \mathbb{E}[x])\|^2}{nN^2\sigma^2}.$$

That is, we center the clean images  $P_s x_s$  by subtracting the projection of the mean volume  $\mathbb{E}[x]$  and compute the square norm of these coefficients before dividing by the noise power  $\sigma^2$ . As a result, for a fixed  $\sigma^2$ , a dataset with low variability will yield a lower  $\text{SNR}_h$  compared to a dataset with higher variability. This is the same definition used by Katsevich et al. [38].

We now consider the simulation described previously at different noise levels with  $n = 1024$  images and  $N = 16$ . The relative error in the Frobenius norm  $\|\Sigma_n - \text{Cov}[x]\|_F / \|\text{Cov}[x]\|_F$  is shown in Figure 6(a) as a function of  $\text{SNR}_h$  for  $n = 1024$ ,  $n = 4096$ , and  $n = 16384$ . This agrees with the guarantee provided by (27) in that for larger  $n$ , the error goes to zero. Note that this is independent of the noise level. Higher noise simply requires a larger number of images  $n$  to achieve a given accuracy.

Although decreasing, the error in  $\Sigma_n$  is still high for  $\text{SNR}_h$  range shown in Figure 6(a). As we saw above, however, we are only concerned with the top eigenvector. Figure 6(b) therefore plots the correlation of the top eigenvector of  $\Sigma_n$  with that of  $\text{Cov}[x]$  for different values of  $\text{SNR}_h$  and  $n$ . Below a certain  $\text{SNR}_h$ , the correlation drops to zero while above a critical threshold the eigenvector correlation approaches one. This behavior is typical of the high-dimensional PCA model discussed in Section 4.4 and the critical  $\text{SNR}_h$  threshold can similarly be observed to vary proportionally to the inverse square root of the number of images  $n^{-1/2}$ . Indeed, this behavior is present in Figure 6 as the threshold  $\text{SNR}_h$  decreases with increasing  $n$ . Note that this is consistent with the derivations of Section 4.3, where we found that the achievable resolution  $N$  was proportional to  $n^{1/2}\sigma^{-2}$  which is proportional to  $n^{1/2}\text{SNR}_h$ . From this, we expect the critical  $\text{SNR}_h$  to grow faster than  $n^{-1/2}N$ .

Having established a baseline performance, we study the effect of replacing  $B_n$  with its shrinkage variant  $B_n^{(s)}$  in the estimator. Figure 6(c) plots the eigenvector correlation as a function of  $\text{SNR}_h$  for both  $\Sigma_n$  and  $\Sigma_n^{(s)}$ . The shrinkage makes a difference when  $\text{SNR}_h$  is between 0.001 and 0.01. Above 0.01, while the shrinkage provides a good estimate of the covariance, the top eigenvector is already well-correlated even without shrinkage, so there is little difference in performance. Below 0.001, the signal eigenvalues are absorbed by the noise bulk, so there is no possibility of extracting accurate eigenvectors and both variants perform

badly. Between these values, however, shrinkage makes a difference, obtaining an eigenvector correlation of 0.8 for an  $\text{SNR}_h$  about 1.4 times lower than the standard estimator.

We also study the robustness of the estimation algorithm with respect to errors in viewing angle estimation. Running the standard orientation estimation algorithms in the ASPIRE toolbox, which relies on class averaging [93] followed by a common-lines based synchronization [73, 75], we apply the covariance estimation method using the estimated viewing angles. The results are shown in Figure 6(d). We see that for this particular set of molecular structures, we recover the viewing angles with enough accuracy to allow accurate covariance estimation.

So far, the distribution of viewing angles has been uniform over  $\text{SO}(3)$ , but this is not necessary. To demonstrate robustness of  $\Sigma_n$  to non-uniform distributions, we draw the rotations  $R_s$  from a family of distributions on  $\text{SO}(3)$  indexed by a parameter  $\delta$  which determines the skew of the distribution towards the identity rotation  $I_3$ . Representing the rotation matrices using Euler angles  $(\alpha, \beta, \gamma)$  in the relative z-y-z convention, we consider the distributions

$$(92) \quad \begin{aligned} \alpha &\sim U[0, 2\pi] \\ \beta &\sim \cos^{-1}(2U[0, 1]^\delta - 1) \text{ ,} \\ \gamma &\sim U[0, 2\pi] \end{aligned}$$

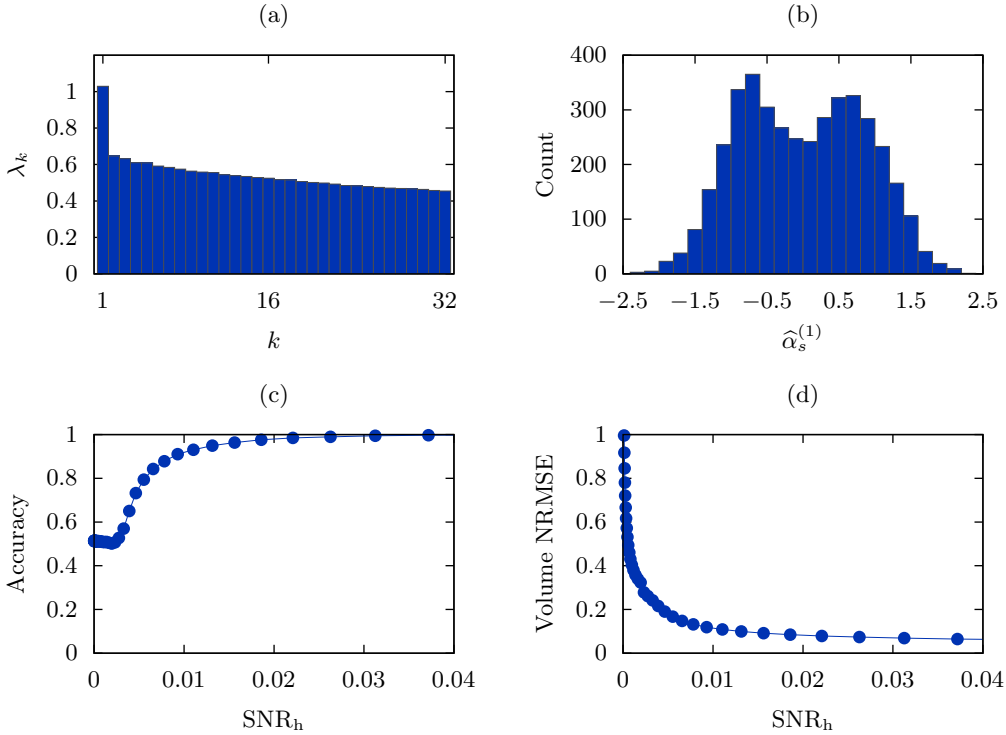
where  $\delta = 1$  is a uniform distribution over  $\text{SO}(3)$  and higher values of  $\delta$  concentrate the distribution closer to  $I_3$ . The resulting eigenvector correlations are shown in Figure 6(e) as a function of  $\text{SNR}_h$  for different values of  $\delta$ . As long as the distribution is not too skewed, we are able to recover the covariance structure accurately at low signal-to-noise ratios.

Increasing the number of classes to  $C = 4$ , we see that the method handles this type of variability just as well as for two classes. Since  $C = 4$ , the population covariance is of rank 3, so instead of evaluating the top eigenvector of  $\Sigma_n$ , we need to consider the top 3 eigenvectors. A simple correlation will not do, and we instead consider the maximum principal angles between the subspace spanned by the top three eigenvectors of  $\Sigma_n$  and those of  $\text{Cov}[x]$ . The cosine  $\cos \theta_{\max}$  of the maximum principal angle  $\theta_{\max}$  between two subspaces  $U$  and  $V$  is given by

$$(93) \quad \cos \theta_{\max} = \min_{u \in U, v \in V} \frac{|\langle u, v \rangle|}{\|u\| \|v\|}$$

If  $U$  and  $V$  are the ranges of two orthogonal matrices  $Q_U$  and  $Q_V$ ,  $\cos \theta_{\max}$  is smallest non-zero singular value of  $Q_U^T Q_V$ . This value is plotted as a function of  $\text{SNR}_h$  in Figure 6(f).

**7.2. Clustering results.** We now consider the clustering and reconstruction steps for the baseline estimator at  $\text{SNR}_h = 0.01$ , where the top eigenvector correlation equals 0.91. Figure 7(a) shows the spectrum of the covariance estimate  $\Sigma_n$ . Since  $C = 2$ , we expect there to be one dominant eigenvalue since the population covariance is of rank one. Indeed, there is one eigenvalue that stands out from the bulk noise distribution, so we form  $\Sigma_{n,1}$  by extracting the dominant eigenvector and eigenvalue. The Wiener filter described in Section 6.1 gives us a set of scalar coordinate estimates  $\{\hat{\alpha}_1^{(1)}, \dots, \hat{\alpha}_n^{(1)}\}$ . Their histogram is shown in Figure 7(b). A clear bimodal distribution suggests that we do indeed have two molecular structures present in the data. Clustering the coordinates using  $k$ -means as described in Section 6.3 and comparing with the ground truth assignments, we achieve 91.8% accuracy.

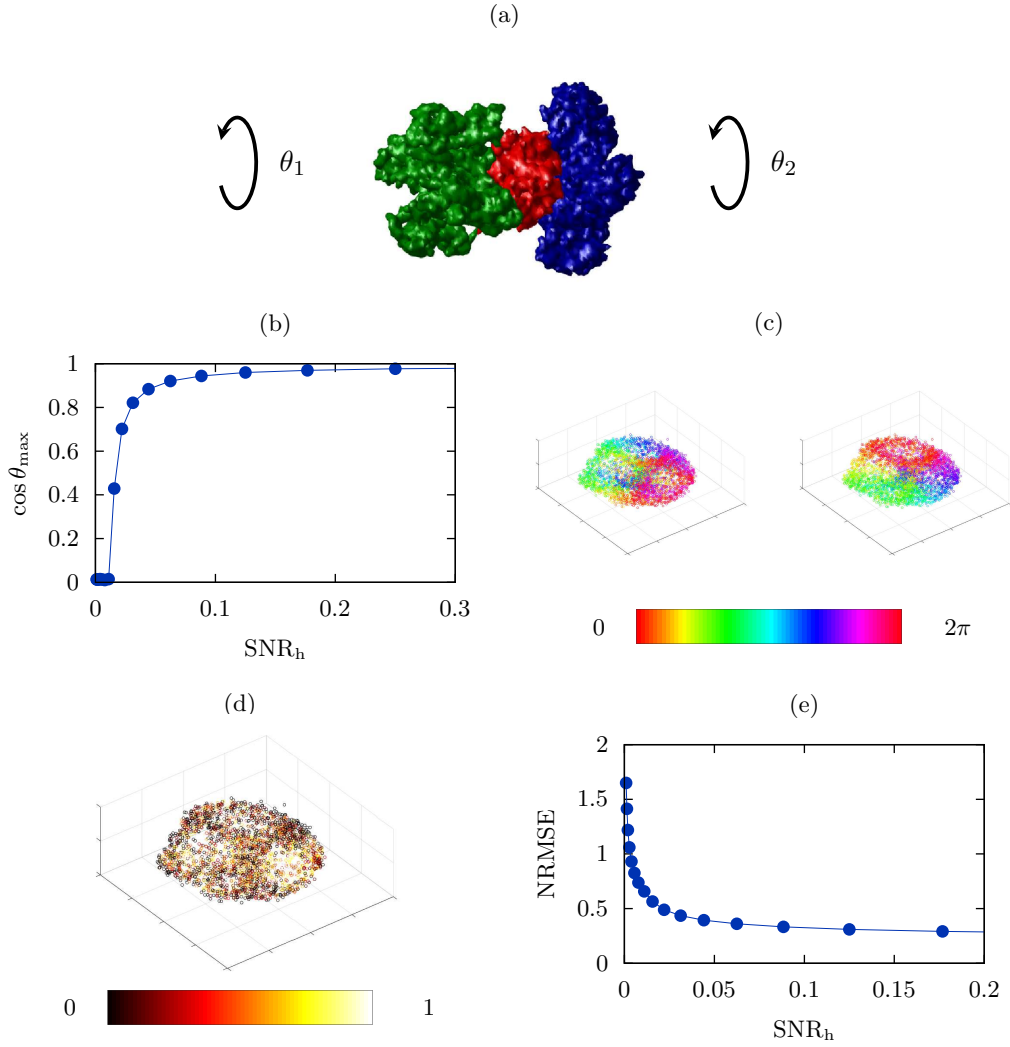


**Figure 7.** Clustering results for discrete variability with  $C = 2$  classes imaged using  $n = 4096$  images with resolution  $N = 16$  for uniform distribution of viewing angles and seven distinct CTFs. (a) The top 32 eigenvalues of  $\Sigma_n$  obtained at  $\text{SNR}_h = 0.01$ . (b) A histogram of the coordinates  $\hat{\alpha}_s^{(1)}$  corresponding to the images  $y_s$  for  $s = 1, \dots, n$  subject to the same  $\text{SNR}_h$ . (c) The fraction of images classified correctly as a function of  $\text{SNR}_h$ . (d) The normalized root mean square error (NRMSE) of the reconstructed volumes.

We now plot the clustering accuracy with respect to  $\text{SNR}_h$  in Figure 7(c). Similarly, the reconstruction error with respect to  $\text{SNR}_h$  is plotted in Figure 7(d). As expected, we observe a phase transition phenomenon similar to that of the top eigenvector correlations. Once a certain threshold is passed, we classify well and obtain high-quality reconstructions. Below the threshold, however, the estimated eigenvectors correlate badly with the population eigenvectors so we do not identify the important directions of variability in the molecules. As a result, the subsequent clustering and reconstruction steps fail.

**7.3. Manifold learning results.** To simulate continuous variability, we deform a potassium channel molecule by independently rotating two parts by angles  $\theta_1$  and  $\theta_2$  as shown in Figure 8(a). This yields a two-parameter family of molecular structures. The manifold described by these volumes is the two-dimensional torus, which can be embedded in three dimensions.

The population covariance of the volumes is not strictly low-rank, but 83% of the variance is concentrated in the leading 4 eigenvectors. If we can recover this eigenspace using our covariance estimation method, we should be able to estimate the manifold structure of the continuous variability. The cosine of the maximum principal angle between the top four population eigenvectors and those of the estimated covariance  $\Sigma_n$  (again with  $\nu_n = \xi_n = 0$ ) is



**Figure 8.** Manifold learning results for continuous variability. (a) Volumes are generated by independently rotating two parts (green and blue) by angles  $\theta_1$  and  $\theta_2$  while keeping the remainder (red) fixed. (b) The cosine of the maximum principal angle between the top four population eigenvectors and those of  $\Sigma_n$ . (c) Three-dimensional diffusion map embedding coordinates of the volume coordinates  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_n\}$ , colored according to the first and second rotation angles. (d) The NRMSE of each volume estimate as a function of its diffusion map coordinate. (e) The NRMSE of the reconstructed volumes as a function of  $\text{SNR}_h$ .

plotted in Figure 8(b) as a function of  $\text{SNR}_h$ . For an  $\text{SNR}_h$  above 0.05 these top eigenvectors are well-estimated, with the cosine of the maximum principal angle in excess of 0.90.

Fixing the  $\text{SNR}_h$  at 0.125, we calculate the coordinates  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_n\}$  of the images using the mean estimate  $\mu_n$  and the top four eigenvectors of  $\Sigma_n$ . From these we compute a diffusion map embedding. Figure 8(c) plots first three embedding coordinates: first colored according to one rotation angle, second colored according to the other. The embedding successfully reproduces these angles, indicating that the procedure captures the two-parameter structure quite well.

**Table 1**

The condition numbers obtained for  $A_n$ ,  $C_n^{-1}A_n$ ,  $L_n$ , and  $D_n^{-1}L_n$  defined from projection mappings  $P_1, \dots, P_n$  where  $n = 16384$  or  $n = 1024$  and  $N = 16$ . The projection mappings are obtained from uniform or non-uniform viewing direction distributions with or without CTFs.

$n$	Distribution of $R$	CTF	$\kappa(A_n)$	$\kappa(C_n^{-1}A_n)$	$\kappa(L_n)$	$\kappa(D_n^{-1}L_n)$
16384	uniform	no	23	6.6	650	170
	uniform	yes	16	12	1400	230
	non-uniform	yes	45	17	4800	720
1024	uniform	no	24	6.5	800	180
	uniform	yes	18	12	4600	400
	non-uniform	yes	53	17	24000	1800

Comparing the estimated volumes  $\hat{x}_s$  with the ground truth volumes  $x_s$  for  $s = 1, \dots, 4096$ , we obtain an NRMSE of 0.31. From this, the recovered range of molecular structure appears to be quite accurate. Plotting the NRMSE as a function of the diffusion map coordinates in Figure 8(d), we see that in areas with high sampling density, the reconstruction is more accurate compared to areas with more sparse sampling. Since the neighborhoods in the sparsely sampled regions have fewer images, this loss of accuracy is expected.

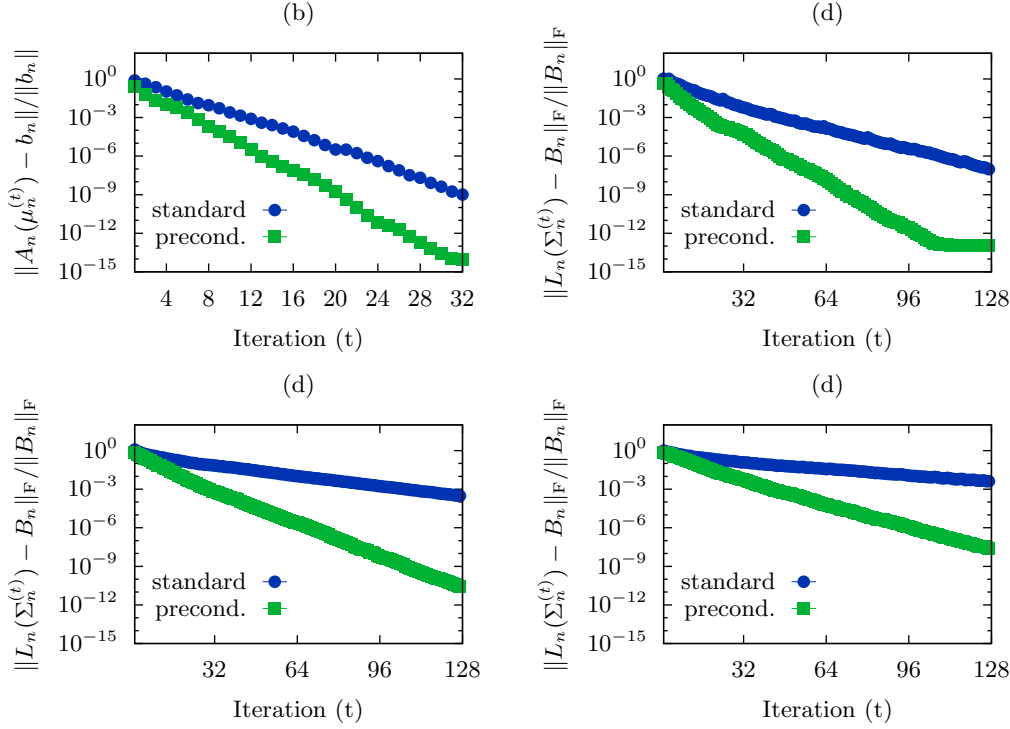
The accuracy of manifold learning reconstruction as a function of  $\text{SNR}_h$  is shown in Figure 8(e). As in the discrete case, higher  $\text{SNR}_h$  yields more accurate reconstructions for fixed  $n$ .

**7.4. Conditioning and convergence results.** The conjugate gradient method requires  $O(\sqrt{\kappa(Z)})$  iterations to invert the operator  $Z$  up to a fixed accuracy, where  $\kappa(Z)$  is its condition number [67, 5, 79]. As such, we would like to compare this quantity for the operators  $A_n$  and  $C_n^{-1}A_n$ , as well as  $L_n$  and  $D_n^{-1}L_n$  in order to evaluate the effect of the strategy outlined in Section 5.3. In addition, we examine its effect on the convergence rate of the CG method.

To estimate the condition numbers, we generate a dataset with  $n = 16384$  projection mappings  $P_1, \dots, P_n$  with  $N = 16$ . Setting the regularization parameters  $\nu_n$  and  $\xi_n$  to zero, we estimate condition numbers in three scenarios. First, we have no CTF (that is  $\mathcal{F}h_s = 1$  for all  $s = 1, \dots, n$ ) and uniform distribution of viewing angles over  $\text{SO}(3)$ . The second scenario includes three distinct CTFs, but with uniform distribution of viewing directions. Finally, we generate viewing directions using a non-uniform distribution and combine with CTFs. The resulting condition numbers for  $A_n$ ,  $C_n^{-1}A_n$ ,  $L_n$ , and  $D_n^{-1}L_n$  are shown in Table 1.

For all operators, adding CTFs and making the viewing direction distribution non-uniform generally worsens the condition number. The exception is adding the CTF for  $A_n$ , which improves its conditioning slightly. Since the CTF may boost certain high frequencies, such an improvement is not unreasonable. Another feature of these results is that covariance estimation is inherently more ill-conditioned compared to mean estimation for a given resolution  $N$  and sample size  $n$ . This confirms our analysis of Section 4.3 showing that the larger number of unknowns in covariance estimation renders it fundamentally harder than mean estimation.

From the same analysis, we see that  $A_n$  requires  $n$  to scale as  $N$  to maintain well-posedness, while  $L_n$  requires  $n$  to scale as  $N^2$ . Consequently, reducing  $n$  for fixed  $N$  should have a greater impact on the conditioning of  $L_n$  compared to  $A_n$ . This is indeed what we see in Table 1



**Figure 9.** (a) The relative residuals for each iteration of CG applied to  $A_n\mu_n = b_n$ , denoted by  $\mu_n^{(t)}$ , with no CTF and uniform distribution of viewing angles. (b-d) The relative residuals for each CG iterate  $\Sigma_n^{(t)}$  of  $L_n(\Sigma_n) = B_n$  with (b) no CTF and uniform distribution of viewing angles, (c) three distinct CTFs and uniform distribution of viewing angles, and (d) three distinct CTFs and non-uniform distribution of viewing angles. For all plots, the residuals of the standard (unpreconditioned) CG method is compared with using a circulant preconditioner. All methods were applied to  $n = 16384$  images with size  $N = 16$  and  $\sigma^2 = 1$ .

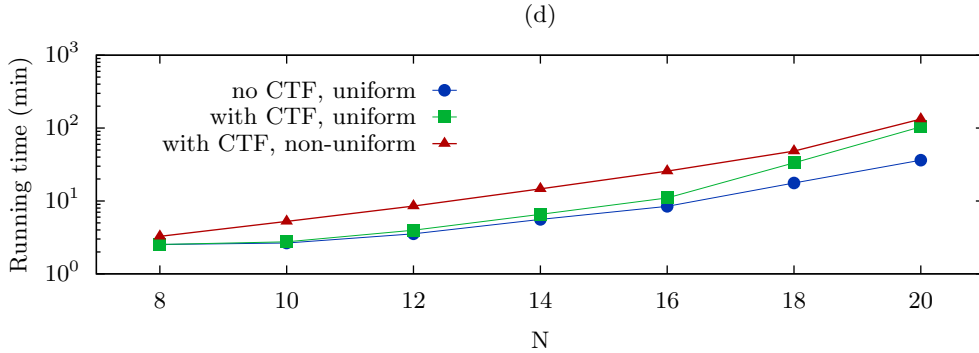
when we redo the experiments for  $n = 1024$ . While the condition numbers of  $A_n$  and  $C_n^{-1}A_n$  remain unchanged, those of  $L_n$  and  $D_n^{-1}L_n$  deteriorate significantly.

Finally, we see that the preconditioning gives a decent improvement in the condition number for all operators, even for small  $n$ . Notably, it reduces  $\kappa$  by a factor of 6 for  $L_n$  in the case of non-uniform distribution of viewing directions and with CTF included when  $n = 16384$ . Since the number of iterations required for convergence of CG scales with  $\sqrt{\kappa}$ , this should reduce the number of iterations by a factor of at least 2.5.

We can observe the effects of preconditioning on the convergence rate of the CG method. Here we generate a dataset with  $C = 2$  classes by applying the forward model specified by  $P_1, \dots, P_n$  as generated previously and adding a noise of variance  $\sigma^2 = 1$ .

In Figure 9(a) we have the CG residuals at different iterations for  $A_n\mu_n = b_n$ , with and without preconditioning, for the baseline case of no CTF and uniform viewing angles. Preconditioning allows CG to converge much faster, reaching a relative residual of  $10^{-3}$  in 8 iterations instead of 13, a reduction by a factor of 1.6. This is close to the factor 1.9 predicted by examining the condition numbers  $\kappa(A_n)$  and  $\kappa(C_n^{-1}A_n)$ .





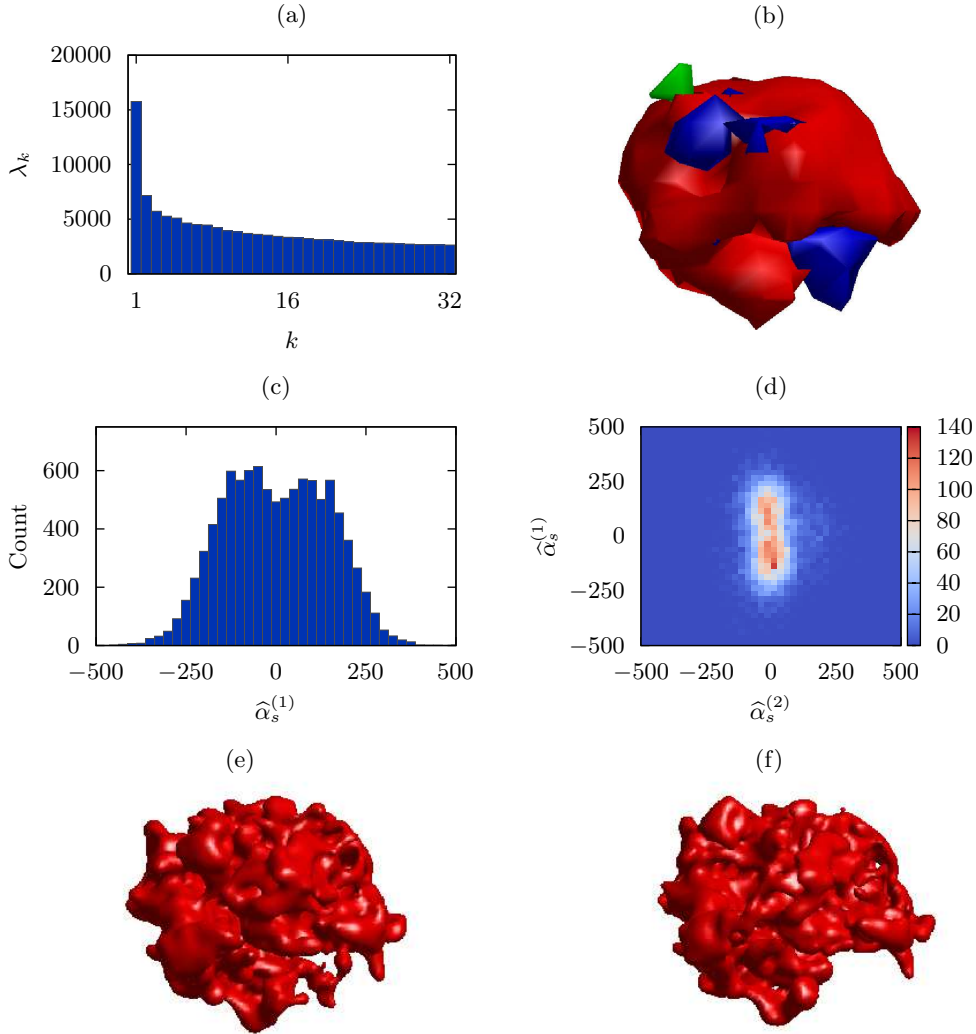
**Figure 10.** Running times for the whole covariance estimation algorithm applied to a dataset of size  $n = 16384$  with varying image size  $N$ . Three scenarios are considered: no CTF with uniform distribution of viewing angles, three distinct CTFs with uniform distribution of viewing angles, and three CTFs with non-uniform distribution of viewing angles.

Similar plots are obtained for CG applied to  $L_n(\Sigma_n) = B_n$  in Figure 9(b), which similarly has no CTF and uniform distribution of viewing angles. Reaching relative residual  $10^{-3}$  takes 20 iterations compared to 47 for the non-preconditioned algorithm. In Figure 9(c), we see that adding CTFs slows convergence for both methods but with the preconditioned algorithm showing a similar improvement. It reduces the number of iterations required to attain a relative residual of  $10^{-3}$  from 107 to 31, resulting in a speedup of at least 3 times. Finally, making the distribution of viewing angles non-uniform gives an even worse convergence rate as shown in Figure 9(d). However, the preconditioned algorithm only requires 44 iterations to converge while without preconditioning, more than 128 iterations are required. These reduction in the number of iterations are consistent with the condition numbers in Table 1.

We note that applying  $D_n^{-1}L_n$  requires more time than  $L_n$ . However,  $D_n^{-1}$  only involves 6D FFTs of size  $N$ , while the 6D FFTs in  $L_n$  have size  $2N$ , that is 64 times larger. The additional time required is therefore expected to be small. Indeed, one preconditioned iteration of CG requires, on average, 15% more time compared to the non-preconditioned iteration.

**7.5. Running time results.** To evaluate the efficiency of our method, we measure its running time on the datasets outlined in the previous section. That is, we have  $n = 16384$  images in three different configurations. The first has no CTF and uniform distribution of viewing angles, whereas the second introduces three distinct CTFs into the imaging model. Finally, the third set of experiments adds a non-uniform distribution of viewing directions. The running times on dual 14-core 2.4 GHz Intel Xeon CPUs are reported in Figure 10.

We see that in the relatively well-conditioned case of uniform distribution of viewing angles for  $N = 16$ , the algorithm terminates quite quickly. It takes around 5 minutes when no CTFs are present and 6 minutes with CTFs. When we no longer have a uniform distribution of viewing angles, we have a running time of 8 minutes. In all scenarios, the dominant step is the conjugate gradient method, since each iteration requires a 6D FFT of size  $2N$  and a large number of iterations can be required when the system is ill-conditioned. That being said, the results compare favorably with other methods for CPU-based heterogeneity, such as RELION,



**Figure 11.** Covariance estimation on the 70S ribosome dataset. (a) Largest eigenvalues of the estimated covariance matrix  $\Sigma_n$ . (b) The estimated mean volume (red), together with the positive (blue) and negative (green) components of the top eigenvector. (c) Histogram of coordinates  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_n\}$  from the Wiener filter estimator. (d) Two-dimensional histogram of coordinates  $\{(\hat{\alpha}_1^{(1)}, \hat{\alpha}_1^{(2)}), \dots, (\hat{\alpha}_n^{(1)}, \hat{\alpha}_n^{(2)})\}$  from the Wiener filter estimator. (e, f) Full-resolution reconstructions obtained using RELION applied to the clusters identified in (c).

which can take up to 800 minutes to run on the same dataset.

**8. Experimental results.** Although simulations are useful to understand the workings of an algorithm, they do not suffice to demonstrate the practical usefulness of a tool. We therefore investigate our covariance estimation approach on an experimental dataset obtained from real cryo-EM samples. This is a standard dataset published by Joachim Frank's lab [46] consisting of projections of a 70S ribosome in two distinct conformations. The dataset comes with a labeling of the images as coming from one of the two states. This labeling was obtained

using supervised classification.

For the 70S ribosome dataset, we have 10000 images of size 130-by-130. We first run the RELION software [68] with the number of classes set to one. This provides an estimate of the viewing angles and translations associated with each projection image which are then given as input to our algorithms. To speed up computations, we downsample the images to 16-by-16. Since our goal is to cluster the images, it is not necessary to do this at full resolution if the discriminant features are already present at low resolution. The images are then whitened so that the noise is approximately white with variance  $\sigma^2 = 1$ . Given the images and the estimates of viewing angles and translations, we apply our mean estimation algorithm (see Algorithm 1). Using the mean estimate  $\mu_n$  (for  $\nu_n = 0$ ), we then apply our covariance estimation method (see Algorithm 2) with  $\xi_n = 2^{-10}$ .

The top eigenvalues of the covariance matrix estimate  $\Sigma_n$  are shown in Figure 11(a). There is a large eigenvalue well-separated from the rest, suggesting that the variability in molecular structure is at least one-dimensional, which is consistent with a two-class configuration. However, there is also a second eigenvalue of significant amplitude, possibly indicating the presence of a small third class. Figure 11(b) shows the estimated mean volume (red) with the leading eigenvector superimposed (positive part in blue, negative part in green). In the top part of the molecule, we observe the rotation of a small subunit indicated by the negative values on one side and positive values on the other side of that subunit. In the bottom, there is a subunit that attaches and detaches depending on the class.

To investigate this further, we plot the one-dimensional histogram obtained from the first coordinates  $\hat{\alpha}_s^{(1)}$  for  $s = 1, \dots, n$ , and a two-dimensional heatmap obtained from  $(\hat{\alpha}_s^{(1)}, \hat{\alpha}_s^{(2)})$  for  $s = 1, \dots, n$  in Figures 11(c) and 11(d), respectively. The heatmap is obtained by dividing the plane into square boxes and counting the number of points in each box. There does seem to be at least two structures while a third structure is hard to discern. It is therefore likely that this second dimension is due to some continuous variability between the two states.

Clustering the coordinates  $\hat{\alpha}_1^{(1)}, \dots, \hat{\alpha}_n^{(1)}$  into two classes and reconstructing from these subsets at full resolution, we obtain the two molecular structures shown in Figures 11(e,f). The two are very similar, with two differences consisting of the rotating subunit at the top and the appearance of the bottom subunit, which agrees the changes observed in the leading eigenvector (see Figure 11(b)). This is in agreement with the presumed structures of the dataset, which consists of a ribosome with and without EF-G. If we compare our clustering with the given annotation of the dataset, we obtain 88.7% accuracy. To compare, the accuracy achieved by RELION on the same dataset with number of classes set to two is 84.6%.

The total running time of Algorithms 1 and 2 was 9 minutes on a quad-core 3.4 GHz Intel Core i7 CPU. The initial estimation of viewing angles and translations using RELION took 356 minutes. This compares with running RELION configured with two classes, which had a running time of 520 minutes. Enabling support for GPU with a GeForce GTX 980 Ti, these running times dropped to 22 and 28 minutes for one and two classes, respectively.

**9. Future Work.** The proposed algorithm scales as  $N^6 \log N$  in the image size  $N$ . As a result, doubling the size of the image results in more than a 64-fold increase in running time and memory usage. This makes estimation for  $N$  larger than 20 prohibitive in most cases. In addition, the large number of unknowns severely limits the achievable resolution as described

in Section 4.3. To resolve these problems, we need to incorporate further structure into the estimation problem. One approach is to better leverage the approximate low rank of the population covariance. Since the number of unknowns drops from  $\mathcal{O}(N^6)$  to  $\mathcal{O}(N^3)$ , this is a more well-posed problem which could be solved faster and using less memory. This could be done by explicitly fixing the rank of  $\Sigma_n$  during estimation. We will also explore related approaches for low-rank matrix recovery such as low-rank matrix sensing via alternating minimization [34] and direct shrinkage of singular values [18].

Another drawback of the above approach is that it requires the viewing directions and translations to be known in advance. While this may be possible for small, localized heterogeneity, where homogeneous reconstruction methods can provide reasonable estimates, this is not always feasible. In that case, methods which combine heterogeneous reconstruction and parameter estimation enjoy a significant advantage [69, 64, 28, 44, 45]. Extending the proposed method for this more general setting provides another avenue for future research.

Finally, more work is needed to process the coordinate estimates  $\hat{\alpha}_1, \dots, \hat{\alpha}_n$ . While standard clustering and manifold learning approaches outlined in this paper provide reasonable results, they do not perform well at high noise levels. A subject of further investigation is therefore how to incorporate informative priors on the space of volumes, such as variability caused by deformation, which should prove useful in this regime.

**10. Conclusion.** We have introduced a computationally efficient method for least-squares estimation of the 3D covariance matrix of the molecular structure from noisy 2D projection images. Given  $n$  images of size  $N$ -by- $N$ , it has computational complexity  $\mathcal{O}(\sqrt{\kappa}N^6 \log N + nN^4)$ , where  $\kappa$  is in the range 1–200 for typical problems. This is achieved by reformulating the normal equations as a deconvolution problem solved using the preconditioned conjugate gradient method. We also introduced a shrinkage variant which improves accuracy at low signal-to-noise ratios, decreasing by a factor of 1.4 the necessary signal power for accurate estimation. The estimated covariance matrices are then used to reconstruct the three-dimensional structures a Wiener filter, which are then clustered into a number of discrete states or fitted to a continuous manifold structure. The accurate performance of both methods is confirmed through experiments on simulated and experimental datasets.

**11. Acknowledgments.** The authors would like to thank Fred Sigworth and Joachim Frank for invaluable discussions regarding the heterogeneity problem and single-particle reconstruction more generally. They are also very grateful to the reviewers for their valuable comments. Initial results on manifold learning for continuous heterogeneity were obtained by Hugh Wilson. This work was performed while the first author was a postdoctoral research associate in the Program in Applied and Computational Mathematics at Princeton University.

## REFERENCES

- [1] G. S. AMMAR AND W. B. GRAGG, *The generalized Schur algorithm for the superfast solution of Toeplitz systems*, in Rational approximation and its applications in mathematics and physics, J. Gilewicz, M. Pindor, and W. Siemaszko, eds., vol. 1237 of Lecture Notes in Mathematics, Springer, 1987, pp. 315–330, <https://doi.org/10.1007/BFb0072474>.
- [2] A. AMUNTS, A. BROWN, X.-C. BAI, J. L. LLÁCER, T. HUSSAIN, P. EMSLEY, F. LONG, G. MURSHUDOV, S. H. W. SCHERES, AND V. RAMAKRISHNAN, *Structure of the yeast mitochondrial large ribosomal*

- subunit*, Science, 343 (2014), pp. 1485–1489, <https://doi.org/10.1126/science.1249410>.
- [3] J. ANDÉN, E. KATSEVICH, AND A. SINGER, *Covariance estimation using conjugate gradient for 3D classification in CRYO-EM*, in Proc. ISBI, April 2015, pp. 200–204, <https://doi.org/10.1109/ISBI.2015.7163849>.
- [4] J. ANDÉN AND A. SINGER, *Factor analysis for spectral estimation*, in Proc. SampTA, July 2017, pp. 169–173, <https://doi.org/10.1109/SAMP.2017.8024447>.
- [5] O. AXELSSON, *Iterative solution methods*, Cambridge university press, 1996.
- [6] A. BARNETT, L. GREENGARD, A. PATAKI, AND M. SPIVAK, *Rapid solution of the cryo-EM reconstruction problem by frequency marching*, SIAM J. Imaging Sci., 10 (2017), pp. 1170–1195, <https://doi.org/10.1137/16M1097171>.
- [7] W. T. BAXTER, R. A. GRASSUCCI, H. GAO, AND J. FRANK, *Determination of signal-to-noise ratios and spectral SNRs in cryo-EM low-dose imaging of molecules*, J. Struct. Biol., 166 (2009), pp. 126–132, <https://doi.org/10.1016/j.jsb.2009.02.012>.
- [8] T. BHAMRE, T. ZHANG, AND A. SINGER, *Anisotropic twicing for single particle reconstruction using autocorrelation analysis*. Submitted, 2017, <https://arxiv.org/abs/1704.07969>.
- [9] R. H. CHAN AND M. K. NG, *Conjugate gradient methods for Toeplitz systems*, SIAM Rev., 38 (1996), pp. 427–482, <https://doi.org/10.1137/S0036144594276474>.
- [10] T. F. CHAN, *An optimal circulant preconditioner for Toeplitz systems*, SIAM J. Sci. and Stat. Comput., 9 (1988), pp. 766–771, <https://doi.org/10.1137/0909051>.
- [11] X. CHENG, *Random Matrices in High-dimensional Data Analysis*, PhD thesis, Princeton University, 2013, <http://arks.princeton.edu/ark:/88435/dsp01wh246s26t>.
- [12] Y. CHENG, N. GRIGORIEFF, P. PENCZEK, AND T. WALZ, *A primer to single-particle cryo-electron microscopy*, Cell, 161 (2015), pp. 438–449, <https://doi.org/https://doi.org/10.1016/j.cell.2015.03.050>.
- [13] R. R. COIFMAN AND S. LAFON, *Diffusion maps*, Appl. Comput. Harmon. Anal., 21 (2006), pp. 5–30, <https://doi.org/https://doi.org/10.1016/j.acha.2006.04.006>.
- [14] J. W. COOLEY AND J. W. TUKEY, *An algorithm for the machine calculation of complex Fourier series*, Math. Comp., 19 (1965), pp. 297–301, <https://doi.org/10.2307/2003354>.
- [15] A. DASHTI, P. SCHWANDER, R. LANGLOIS, R. FUNG, W. LI, A. HOSSEINIZADEH, H. Y. LIAO, J. PALLESEN, G. SHARMA, V. A. STUPINA, A. E. SIMON, J. D. DINMAN, J. FRANK, AND A. OURMAZD, *Trajectories of the ribosome as a Brownian nanomachine*, Proc. Natl. Acad. Sci. U.S.A., 111 (2014), pp. 17492–17497, <https://doi.org/10.1073/pnas.1419276111>.
- [16] C. DAVIS AND W. M. KAHAN, *The rotation of eigenvectors by a perturbation. III*, SIAM J. Numer. Anal., 7 (1970), pp. 1–46, <https://doi.org/10.1137/0707001>.
- [17] A. P. DEMPSTER, N. M. LAIRD, AND D. B. RUBIN, *Maximum likelihood from incomplete data via the EM algorithm*, J. Royal Stat. Soc. B Stat. Methol., 39 (1977), pp. 1–38, <http://www.jstor.org/stable/2984875>.
- [18] E. DOBRIBAN, W. LEEB, AND A. SINGER, *Optimal prediction in the linearly transformed spiked model*. Submitted, 2017, <https://arxiv.org/abs/1709.03393>.
- [19] D. L. DONOHO, M. GAVISH, AND I. M. JOHNSTONE, *Optimal shrinkage of eigenvalues in the spiked covariance model*. 2013, <https://arxiv.org/abs/1311.0851>.
- [20] V. DRUSKIN AND L. KNIZHNERMAN, *Two polynomial methods of calculating functions of symmetric matrices*, USSR Comput. Math. Math. Phys., 29 (1989), pp. 112–121, [https://doi.org/10.1016/S0041-5553\(89\)80020-5](https://doi.org/10.1016/S0041-5553(89)80020-5).
- [21] A. DUTT AND V. ROKHLIN, *Fast Fourier transforms for nonequispaced data*, SIAM J. Sci. Comput., 14 (1993), pp. 1368–1393, <https://doi.org/10.1137/0914081>.
- [22] H. ERICKSON AND A. KLUG, *Measurement and compensation of defocusing and aberrations by Fourier processing of electron micrographs*, Philos. Trans. Royal Soc. B, 261 (1971), pp. 105–118, <https://doi.org/10.1098/rstb.1971.0040>.
- [23] J. A. FESSLER, S. LEE, V. T. OLAFSSON, H. R. SHI, AND D. C. NOLL, *Toeplitz-based iterative image reconstruction for MRI with correction for magnetic field inhomogeneity*, IEEE Trans. Sig. Process., 53 (2005), pp. 3393–3402, <https://doi.org/10.1109/TSP.2005.853152>.
- [24] J. FRANK, *Three-dimensional electron microscopy of macromolecular assemblies*, Academic Press, 2006.
- [25] M. GAVISH AND D. L. DONOHO, *Optimal shrinkage of singular values*, IEEE Trans. Inf. Theory, 63 (2017), pp. 2137–2152, <https://doi.org/10.1109/TIT.2017.2653801>.



- [26] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, Johns Hopkins University Press, 4th ed., 2013.
- [27] L. GREENGARD AND J.-Y. LEE, *Accelerating the nonuniform fast Fourier transform*, *SIAM Rev.*, 46 (2004), pp. 443–454, <https://doi.org/10.1137/S003614450343200X>.
- [28] N. GRIGORIEFF, *FREALIGN: High-resolution refinement of single particle structures*, *J. Struct. Biol.*, 157 (2007), pp. 117 – 125, <https://doi.org/10.1016/j.jsb.2006.05.004>.
- [29] M. GUERQUIN-KERN, D. V. D. VILLE, C. VONESCH, J.-C. BARITAUX, K. P. PRUESSMANN, AND M. UNSER, *Wavelet-regularized reconstruction for rapid MRI*, in *Proc. ISBI, IEEE*, June 2009, pp. 193–196, <https://doi.org/10.1109/ISBI.2009.5193016>.
- [30] G. HARAUZ AND M. VAN HEEL, *Exact filters for general geometry three dimensional reconstruction*, *Optik*, 73 (1986), pp. 146–156.
- [31] G. T. HERMAN, *Fundamentals of computerized tomography: Image reconstruction from projections*, Springer-Verlag London, 2009, <https://doi.org/10.1007/978-1-84628-723-7>.
- [32] M. R. HESTENES AND E. STIEFEL, *Methods of conjugate gradients for solving linear systems*, *J. Res. Natl. Bur. Stand.*, 49 (1952), <https://doi.org/10.6028/jres.049.044>.
- [33] N. J. HIGHAM, *Functions of matrices: Theory and computation*, Society for Industrial and Applied Mathematics, 2008, <https://doi.org/10.1137/1.9780898717778>.
- [34] P. JAIN, P. NETRAPALLI, AND S. SANGHAVI, *Low-rank matrix completion using alternating minimization*, in *Proc. STOC, ACM*, 2013, pp. 665–674, <https://doi.org/10.1145/2488608.2488693>.
- [35] Q. JIN, C. SORZANO, J. DE LA ROSA-TREVÍN, J. BILBAO-CASTRO, R. NÚÑEZ-RAMÍREZ, O. LLORCA, F. TAMA, AND S. JONIĆ, *Iterative elastic 3D-to-2D alignment method using normal modes for studying structural dynamics of large macromolecular complexes*, *Structure*, 22 (2014), pp. 496–506, <https://doi.org/10.1016/j.str.2014.01.004>.
- [36] I. M. JOHNSTONE, *On the distribution of the largest eigenvalue in principal components analysis*, *Ann. Statist.*, 29 (2001), pp. 295–327, <http://www.jstor.org/stable/2674106>.
- [37] S. JONIĆ, *Computational methods for analyzing conformational variability of macromolecular complexes from cryo-electron microscopy images*, *Curr. Opin. Struct. Biol.*, 43 (2017), pp. 114–121, <https://doi.org/10.1016/j.sbi.2016.12.011>.
- [38] E. KATSEVICH, A. KATSEVICH, AND A. SINGER, *Covariance matrix estimation for the cryo-EM heterogeneity problem*, *SIAM J. Imaging Sci.*, 8 (2015), pp. 126–185, <https://doi.org/10.1137/130935434>.
- [39] M. KHOSHOUEI, M. RADJAINIA, W. BAUMEISTER, AND R. DANEV, *Cryo-EM structure of haemoglobin at 3.2 Å determined with the volta phase plate*, *Nat. Commun.*, 8 (2017), <https://doi.org/10.1038/ncomms16099>.
- [40] D. KIMANIUS, B. O. FORSBERG, S. H. SCHERES, AND E. LINDAHL, *Accelerated cryo-EM structure determination with parallelisation using GPUs in RELION-2*, *eLife*, 5 (2016), <https://doi.org/10.7554/eLife.18722>.
- [41] A. KLUG AND R. A. CROWTHER, *Three-dimensional image reconstruction from the viewpoint of information theory*, *Nature*, 238 (1972), pp. 435–440, <https://doi.org/10.1038/238435a0>.
- [42] W. KÜHLBRANDT, *The resolution revolution*, *Science*, 343 (2014), pp. 1443–1444, <https://doi.org/10.1126/science.1251652>.
- [43] S. KUNIS AND D. POTTS, *Fast spherical Fourier algorithms*, *J. Comput. Appl. Math.*, 161 (2003), pp. 75–98, [https://doi.org/10.1016/S0377-0427\(03\)00546-6](https://doi.org/10.1016/S0377-0427(03)00546-6).
- [44] R. R. LEDERMAN AND A. SINGER, *A representation theory perspective on simultaneous alignment and classification*. Submitted, 2016, <https://arxiv.org/abs/1607.03464>.
- [45] R. R. LEDERMAN AND A. SINGER, *Continuously heterogeneous hyper-objects in cryo-EM and 3-D movies of many temporal dimensions*. Submitted, 2017, <https://arxiv.org/abs/1704.02899>.
- [46] H. LIAO AND J. FRANK, *Classification by bootstrapping in single particle methods*, in *Proc. ISBI, IEEE*, April 2010, pp. 169–172, <https://doi.org/10.1109/ISBI.2010.5490386>.
- [47] H. Y. LIAO, Y. HASHEM, AND J. FRANK, *Efficient estimation of three-dimensional covariance and its application in the analysis of heterogeneous samples in cryo-electron microscopy*, *Structure*, 23 (2015), pp. 1129–1137, <https://doi.org/10.1016/j.str.2015.04.004>.
- [48] M. LIAO, E. CAO, D. JULIUS, AND Y. CHENG, *Structure of the TRPV1 ion channel determined by electron cryo-microscopy*, *Nature*, 504 (2013), pp. 107–112, <https://doi.org/10.1038/nature12822>.
- [49] W. LIU AND J. FRANK, *Estimation of variance distribution in three-dimensional reconstruction. I. Theory*,

- J. Opt. Soc. Am. A, 12 (1995), pp. 2615–2627, <https://doi.org/10.1364/JOSAA.12.002615>.
- [50] S. LLOYD, *Least squares quantization in PCM*, IEEE Trans. Inf. Theory, 28 (1982), pp. 129–137, <https://doi.org/10.1109/TIT.1982.1056489>.
- [51] S. MALLAT, *A wavelet tour of signal processing*, Academic Press, 3rd ed., 2008.
- [52] V. A. MARČENKO AND L. A. PASTUR, *Distribution of eigenvalues for some sets of random matrices*, Math. USSR Sb., 1 (1967), p. 457, <https://doi.org/10.1070/SM1967v001n04ABEH001994>.
- [53] E. MICHIELSEN AND A. BOAG, *A multilevel matrix decomposition algorithm for analyzing scattering from large structures*, IEEE Trans. Antennas Propag., 44 (1996), pp. 1086–1093, <https://doi.org/10.1109/8.511816>.
- [54] J. L. MILNE, M. J. BORGNIA, A. BARTESAGHI, E. E. TRAN, L. A. EARL, D. M. SCHAUDER, J. LENGYEL, J. PIERSON, A. PATWARDHAN, AND S. SUBRAMANIAM, *Cryo-electron microscopy—A primer for the non-microscopist*, FEBS Journal, 280 (2013), pp. 28–45, <https://doi.org/10.1111/febs.12078>.
- [55] J. A. MINDELL AND N. GRIGORIEFF, *Accurate determination of local defocus and specimen tilt in electron microscopy*, J. Struct. Biol., 142 (2003), pp. 334–347, [https://doi.org/10.1016/S1047-8477\(03\)00069-8](https://doi.org/10.1016/S1047-8477(03)00069-8).
- [56] B. MUSICUS, *Levinson and fast Cholesky algorithms for Toeplitz and almost Toeplitz matrices*, Tech. Report 538, MIT Research Laboratory of Electronics, 1988.
- [57] F. NATTERER, *The mathematics of computerized tomography*, Society for Industrial and Applied Mathematics, 2001, <https://doi.org/10.1137/1.9780898719284>.
- [58] D. PAUL, *Asymptotics of sample eigenstructure for a large dimensional spiked covariance model*, Statist. Sinica, 17 (2007), pp. 1617–1642, <http://www.jstor.org/stable/24307692>.
- [59] P. PENCZEK, M. KIMMEL, AND C. SPAHN, *Identifying conformational states of macromolecules by eigenanalysis of resampled cryo-EM images*, Structure, 19 (2011), pp. 1582–1590, <https://doi.org/10.1016/j.str.2011.10.003>.
- [60] P. A. PENCZEK, *Variance in three-dimensional reconstructions from projections*, in Proc. ISBI, 2002, pp. 749–752, <https://doi.org/10.1109/ISBI.2002.1029366>.
- [61] P. A. PENCZEK, *Image restoration in cryo-electron microscopy*, in Cryo-EM, Part B: 3-D Reconstruction, G. J. Jensen, ed., vol. 482 of Methods Enzymol., Academic Press, 2010, pp. 35–72, [https://doi.org/10.1016/S0076-6879\(10\)82002-6](https://doi.org/10.1016/S0076-6879(10)82002-6).
- [62] P. A. PENCZEK, J. FRANK, AND C. M. SPAHN, *A method of focused classification, based on the bootstrap 3D variance analysis, and its application to EF-G-dependent translocation*, J. Struct. Biol., 154 (2006), pp. 184–194, <https://doi.org/10.1016/j.jsb.2005.12.013>.
- [63] P. A. PENCZEK, C. YANG, J. FRANK, AND C. M. SPAHN, *Estimation of variance in single-particle reconstruction using the bootstrap technique*, J. Struct. Biol., 154 (2006), pp. 168–183, <https://doi.org/10.1016/j.jsb.2006.01.003>.
- [64] A. PUNJANI, J. L. RUBINSTEIN, D. J. FLEET, AND M. A. BRUBAKER, *cryoSPARC: algorithms for rapid unsupervised cryo-EM structure determination*, Nat. Methods, 14 (2017), pp. 290–296, <https://doi.org/10.1038/nmeth.4169>.
- [65] J. RADON, *Über die Bestimmung von Funktionen durch ihre Integralwerte längs gewisser Mannigfaltigkeiten*, Berichte Sächsischen Akad. Wissenschaft., Math. Phys. Klass, 69 (1917), pp. 262–277.
- [66] Y. SAAD, *Analysis of some Krylov subspace approximations to the matrix exponential operator*, SIAM J. Numer. Anal., 29 (1992), pp. 209–228, <https://doi.org/10.1137/0729014>.
- [67] Y. SAAD, *Iterative methods for sparse linear systems*, SIAM, 2nd ed., 2003, <https://doi.org/10.1137/1.9780898718003>.
- [68] S. SCHERES, *RELION: Implementation of a Bayesian approach to cryo-EM structure determination*, J. Struct. Biol., 180 (2012), pp. 519–530, <https://doi.org/10.1016/j.jsb.2012.09.006>.
- [69] S. H. SCHERES, *A Bayesian view on cryo-EM structure determination*, J. Mol. Biol., 415 (2012), pp. 406–418, <https://doi.org/10.1016/j.jmb.2011.11.010>.
- [70] J. SCHUR, *Über potenzreihen, die im innern des einheitskreises beschränkt sind.*, J. Reine Angew. Math., 147 (1917), pp. 205–232.
- [71] M. SHATSKY, R. HALL, E. NOGALES, J. MALIK, AND S. BRENNER, *Automated multi-model reconstruction from single-particle electron microscopy data*, J. Struct. Biol., 170 (2010), pp. 98–108, <https://doi.org/10.1016/j.jsb.2010.01.007>.
- [72] J. SHI AND J. MALIK, *Normalized cuts and image segmentation*, IEEE Trans. Pattern Anal. Mach. Intell.,



- 22 (2000), pp. 888–905, <https://doi.org/10.1109/34.868688>.
- [73] Y. SHKOLNISKY AND A. SINGER, *Viewing direction estimation in cryo-EM using synchronization*, SIAM J. Imaging Sci., 5 (2012), pp. 1088–1110, <https://doi.org/10.1137/120863642>.
- [74] F. J. SIGWORTH, *A maximum-likelihood approach to single-particle image refinement*, J. Struct. Biol., 122 (1998), pp. 328–339, <https://doi.org/10.1006/jsbi.1998.4014>.
- [75] A. SINGER, R. R. COIFMAN, F. J. SIGWORTH, D. W. CHESTER, AND Y. SHKOLNISKY, *Detecting consistent common lines in cryo-EM by voting*, J. Struct. Biol., 169 (2010), pp. 312–322, <https://doi.org/10.1016/j.jsb.2009.11.003>.
- [76] D. SLEPIAN, *Prolate spheroidal wave functions, Fourier analysis and uncertainty—IV: Extensions to many dimensions; generalized prolate spheroidal functions*, Bell Syst. Tech. J., 43 (1964), pp. 3009–3057, <https://doi.org/10.1002/j.1538-7305.1964.tb01037.x>.
- [77] G. STRANG, *A proposal for Toeplitz matrix calculations*, Stud. Appl. Math., 74 (1986), pp. 171–176, <https://doi.org/10.1002/sapm1986742171>.
- [78] H. D. TAGARE, A. KUCUKELBIR, F. J. SIGWORTH, H. WANG, AND M. RAO, *Directly reconstructing principal components of heterogeneous particles from cryo-EM images*, J. Struct. Biol., 191 (2015), pp. 245–262, <https://doi.org/10.1016/j.jsb.2015.05.007>.
- [79] L. N. TREFETHEN AND D. BAU, *Numerical linear algebra*, SIAM, 1997.
- [80] M. TYGERT, *Fast algorithms for spherical harmonic expansions, III*, J. Comput. Phys., 229 (2010), pp. 6181–6192, <https://doi.org/10.1016/j.jcp.2010.05.004>.
- [81] E. E. TYRTYSHNIKOV, *Optimal and superoptimal circulant preconditioners*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 459–473, <https://doi.org/10.1137/0613030>.
- [82] M. VAN HEEL, B. GOWEN, R. MATADEEN, E. V. ORLOVA, R. FINN, T. PAPE, D. COHEN, H. STARK, R. SCHMIDT, M. SCHATZ, AND A. PATWARDHAN, *Single-particle electron cryo-microscopy: Towards atomic resolution*, Q. Rev. Biophys., 33 (2000), pp. 307–369, <https://doi.org/10.1017/S0033583500003644>.
- [83] K. R. VINOTHKUMAR AND R. HENDERSON, *Single particle electron cryomicroscopy: Trends, issues and future perspective*, Q. Rev. Biophys., 49 (2016), <https://doi.org/10.1017/S0033583516000068>.
- [84] C. VONESCH, L. WANG, Y. SHKOLNISKY, AND A. SINGER, *Fast wavelet-based single-particle reconstruction in cryo-EM*, in Proc. ISBI, IEEE, March 2011, pp. 1950–1953, <https://doi.org/10.1109/ISBI.2011.5872791>.
- [85] M. VULOVIĆ, R. B. RAVELLI, L. J. VAN VLIET, A. J. KOSTER, I. LAZIĆ, U. LÜCKEN, H. RULLGÅRD, O. ÖKTEM, AND B. RIEGER, *Image formation modeling in cryo-electron microscopy*, J. Struct. Biol., 183 (2013), pp. 19–32, <https://doi.org/10.1016/j.jsb.2013.05.008>.
- [86] R. WADE, *A brief look at imaging and contrast transfer*, Ultramicroscopy, 46 (1992), pp. 145–156, [https://doi.org/10.1016/0304-3991\(92\)90011-8](https://doi.org/10.1016/0304-3991(92)90011-8).
- [87] F. T. A. W. WAJER AND K. P. PRUESSMANN, *Major speedup of reconstruction for sensitivity encoding with arbitrary trajectories*, in Proc. ISMRM, 2001, p. 767.
- [88] L. WANG, Y. SHKOLNISKY, AND A. SINGER, *A Fourier-based approach for iterative 3D reconstruction from cryo-EM images*. Unpublished, 2013, <https://arxiv.org/abs/1307.5824>.
- [89] L. WANG, A. SINGER, AND Z. WEN, *Orientation determination of cryo-EM images using least unsquared deviations*, SIAM J. Imaging Sci., 6 (2013), pp. 2450–2483, <https://doi.org/10.1137/130916436>.
- [90] N. XU, D. VEESLER, P. C. DOERSCHUK, AND J. E. JOHNSON, *Allosteric effects in bacteriophage HK97 procapsids revealed directly from covariance analysis of cryo EM data*, J. Struct. Biol., (2018), <https://doi.org/10.1016/j.jsb.2017.12.013>. In press.
- [91] K. ZHANG, *Gctf: Real-time CTF determination and correction*, J. Struct. Biol., 193 (2016), pp. 1–12, <https://doi.org/10.1016/j.jsb.2015.11.003>.
- [92] Z. ZHAO AND A. SINGER, *Fourier–Bessel rotational invariant eigenimages*, J. Opt. Soc. Am. A, 30 (2013), pp. 871–877, <https://doi.org/10.1364/JOSAA.30.000871>.
- [93] Z. ZHAO AND A. SINGER, *Rotationally invariant image representation for viewing direction classification in cryo-EM*, J. Struct. Biol., 186 (2014), pp. 153–166, <https://doi.org/10.1016/j.jsb.2014.03.003>.
- [94] Y. ZHENG, Q. WANG, AND P. C. DOERSCHUK, *Three-dimensional reconstruction of the statistics of heterogeneous objects from a collection of one projection image of each object*, J. Opt. Soc. Am. A, 29 (2012), pp. 959–970, <https://doi.org/10.1364/JOSAA.29.000959>.