# Random Conical Tilt Reconstruction without Particle Picking in Cryo-electron Microscopy

Ti-Yen Lan,[a]* Nicolas Boumal[b] and Amit Singer[a,c]

[a]*Program in Applied and Computational Mathematics, Princeton University, Princeton, NJ 08544, USA,* [b]*Institute of Mathematics, EPFL, CH-1015 Lausanne, Switzerland,* and [c]*Department of Mathematics, Princeton University, Princeton, NJ 08544, USA. E-mail: tiyenlan@princeton.edu*

**Cryo-EM**; **random conical tilt**; **autocorrelation analysis**; **structure reconstruction**

## Abstract

We propose a method to reconstruct the 3-D molecular structure from micrographs collected at just one sample tilt angle in the random conical tilt scheme in cryo-electron microscopy. Our method uses autocorrelation analysis on the micrographs to estimate features of the molecule which are invariant under certain nuisance parameters such as the positions of molecular projections in the micrographs. This enables us to reconstruct the molecular structure directly from micrographs, completely circumventing the need for particle picking. We demonstrate reconstructions with simulated data and investigate the effect of the missing-cone region. These results show promise to reduce the size limit for single particle reconstruction in cryo-electron microscopy.

## 1. Introduction

Random conical tilt (RCT) (Radermacher *et al.*, 1987; Radermacher, 1988; Sorzano *et al.*, 2015) is an important technique in single-particle cryo-electron microscopy (cryo-EM) to generate a *de novo* 3-D reconstruction, which provides an unbiased initial model for a subsequent iterative refinement process to determine high-resolution structures. The technique applies to molecules that have a preferred orientation to the 2-D substrate they are deposited on and random in-plane rotations. The standard data collection scheme of RCT involves measuring pairs of images, or micrographs, of the same field of view: one with a large sample tilt angle (Figure 1(a)), and one with no tilt (Figure 1(b)). Since the micrograph pairs contain projections of each molecule at two views that are physically related, one can first estimate the in-plane rotation of each molecule by aligning the molecular projections measured in the untilted micrographs and then assemble the corresponding molecular projections recorded in the tilted micrographs to reconstruct the 3-D molecular structure, as shown in Figure 1(c).

However, some limitations exist for the RCT method. The design of the sample holder restricts the maximum tilt angle to about 60°, which makes a considerable fraction of information about the molecular structure inaccessible to the technique: this is the so-called "missing-cone" problem. Another limitation is the need to collect data from the same field of view at two different sample tilt angles. For each of the two tilt angles, the signal-to-noise ratio (SNR) must be high enough so that it is possible to reliably locate the molecular projections (that is, pick particles) in the noisy micrographs. This essentially doubles the required electron dose on the sample. Meanwhile, the molecule must be large enough so that the irreversible structural damage caused by incident electrons is limited enough to allow for particle picking. Indeed, this has led to the common belief that small biological molecules are out of the reach for cryo-EM (Henderson, 1995).

In this study, we develop an approach to reconstruct the 3-D molecular structure from data collected at just one large sample tilt angle, as depicted in Figure 2(a). More importantly, our approach circumvents the need for particle picking to reconstruct the molecular structure directly from the micrographs. The main idea is to first estimate features of the molecule that are invariant to the 2-D positions of molecular projections in the micrographs. The estimation is done through a variant of Kam's autocorrelation analysis (Kam, 1980). We subsequently determine the molecular structure by fitting the estimated invariants through an optimization problem. We address the problem of missing information by adding a regularizer in the optimization. Assuming white noise, this approach can in principle handle cases of arbitrarily low SNR as long as sufficiently many micrographs are used to estimate the invariants. Figure 2(b) shows one such noisy micrograph where particle picking becomes challenging. This observation notably suggests that the feasibility of particle picking does not limit the smallest usable molecule size in single-particle cryo-EM.

Kam's autocorrelation analysis was also applied for analyzing X-ray single particle imaging data (Kam, 1977; Saldin *et al.*, 2010; Donatelli *et al.*, 2015; von Ardenne *et al.*, 2018). In particular, Saldin *et al.* (2010) considered the problem of reconstructing the top-down projection of molecules randomly oriented about a single axis, which is similar to the case of no tilt in RCT. Subsequently, Elser (2011) designed an algorithm to reconstruct the 3-D structure of such partially oriented molecules from a tilt series. Kam's method was recently demonstrated with actual data collected from randomly oriented virus particles (Kurta *et al.*, 2017; Pande *et al.*, 2018).

This work belongs to a methodical program to develop algorithms to reconstruct molecular structures without the need for particle picking, which was first proposed in Bendory *et al.* (2018). The development started with the studies of a simplified 1-D model, where multiple copies of a target signal occur at unknown locations in a noisy

long measurement (Bendory *et al.*, 2018; Bendory *et al.*, 2019; Lan *et al.*, 2020). The extension to the 2-D case, where multiple copies of a target image are randomly rotated and translated in a large noisy measurement image, was later studied in Marshall *et al.* (2020) and Bendory *et al.* (2021). These results can be used to reconstruct the top-down molecular projection from the micrographs collected at no tilt in the RCT scheme.

We organize the rest of the paper as follows. We describe the data simulation procedure in Sections 2.1 to 2.3. The details of our approach are discussed in Sections 2.4 and 2.5. In Section 3, we study the effect of the missing-cone region on the quality of reconstruction and present the reconstructions of two molecular structures from simulated noisy micrographs. The computational details are described in the appendix.

## 2. Methods

### 2.1. Image formation model

In the cryo-EM imaging process, the incident electrons are scattered by the 3-D Coulomb potential of the sample $f_s(x, y, z)$. We define the coordinate system for data collection $S$ by the orthogonal $\mathbf{x}$- and $\mathbf{y}$-axes along the edges of the detector and the normally incident electron beam, as the $\mathbf{z}$-axis. Under the weak-phase object approximation, the micrograph recorded by an $m \times m$ pixelated detector can be modeled as

$$M(x_i, y_i) = (h * \mathcal{P} f_s)(x_i/\xi, y_i/\xi) + \varepsilon(x_i, y_i), \tag{1}$$

where $i \in \{1, \ldots, m^2\}$, $(x_i, y_i) \in \{-\lfloor m/2 \rfloor, \ldots, \lceil m/2 - 1 \rceil\}^2$ is the 2-D coordinate of the $i^{\text{th}}$ pixel, and $\xi$ denotes the pixel sampling rate. The operator $\mathcal{P}$ generates the tomographic projection of $f_s$ along the $\mathbf{z}$-axis by

$$(\mathcal{P} f_s)(x, y) = \int_{-\infty}^{\infty} f_s(x, y, z) \, dz. \tag{2}$$

The 2-D function $h(x, y)$ represents the point spread function of the imaging system, and the operator $*$ denotes the 2-D convolution, where

$$(h * g)(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(u, v) g(x - u, y - v) \, du \, dv \qquad (3)$$

for any 2-D function $g(x, y)$. Finally, the measurement noise is modeled by the additive random variable $\varepsilon(x_i, y_i)$.

In this work, we consider the simplified scenario where we ignore the effect of the point spread function by making the idealistic assumption that it is a 2-D Dirac delta function, namely, $(h * g)(x, y) = g(x, y)$. Moreover, we assume that the random noise $\varepsilon$ is drawn from an i.i.d. Gaussian distribution with zero mean and variance $\sigma^2$. The arising challenges beyond these assumptions will be discussed in Section 4.

## 2.2. Random conical tilt

The sample used in RCT consists of multiple copies of partially oriented molecules. Specifically, the molecules adsorb to a 2-D substrate such that a particular axis within the molecules aligns with the substrate normal. The molecular orientations are limited to rotations about the particular body axis by angles uniformly drawn from $[0, 2\pi)$. Let $S''$ be the body frame of one particular molecule, where the $\mathbf{z}''$-axis coincides with its body rotation axis. We further define another reference frame $S'$ fixed on the 2-D substrate such that the $\mathbf{x}'$-axis coincides with the tilt axis of the substrate and the $\mathbf{z}'$-axis aligns with the substrate normal. In the following, we also assume that the $\mathbf{x}$-axis of the lab frame is parallel to the $\mathbf{x}'$-axis. After specifying these reference frames, we define the substrate tilt angle $\theta$ as the angle between the $\mathbf{z}$- and $\mathbf{z}'$-axes. The rotation angle $\alpha$ of the particular molecule with respect to its body rotation axis is defined as the angle between the $\mathbf{x}'$- and $\mathbf{x}''$-axes. The relationships between the reference frames are shown in Figure 3.

Let $f(x'', y'', z'')$ be the 3-D Coulomb potential of the particular molecule in its own

body frame $S''$. Hereafter, we refer to $f$ as the structure of the molecule. From the geometries shown in Figure 3, the coordinate transformation between $S$ and $S''$ is given by

$$
\begin{aligned}
\mathbf{r} &= \begin{bmatrix} \cos\alpha & -\sin\alpha & 0 \\ \cos\theta\sin\alpha & \cos\theta\cos\alpha & -\sin\theta \\ \sin\theta\sin\alpha & \sin\theta\cos\alpha & \cos\theta \end{bmatrix} \mathbf{r}'' + \mathbf{t} \\
&= R_\alpha^\theta \mathbf{r}'' + \mathbf{t},
\end{aligned}
\tag{4}
$$

where $\mathbf{r} = [x, y, z]^T$, $\mathbf{r}'' = [x'', y'', z'']^T$, $R_\alpha^\theta$ is the rotation matrix that aligns the axes of $S$ with the axes of $S''$, and $\mathbf{t} = [t_x, t_y, t_z]^T$ is the vector pointing from the origin of $S$ to the origin of $S''$. We can therefore express the molecular structure in the lab frame $S$ by $f((R_\alpha^\theta)^T(\mathbf{r} - \mathbf{t}))$, and its tomographic projection along the $\mathbf{z}$-axis is given by $\mathcal{I}_\alpha^\theta(x - t_x, y - t_y)$, where

$$
\mathcal{I}_\alpha^\theta(x, y) = \int_{-\infty}^{\infty} f((R_\alpha^\theta)^T \mathbf{r})\ dz.
\tag{5}
$$

Taking the 2-D Fourier transform on both sides of (5), with the Fourier slice theorem, we get

$$
\hat{\mathcal{I}}_\alpha^\theta(k_x, k_y) = \hat{f}((R_\alpha^\theta)^T [k_x, k_y, 0]^T),
\tag{6}
$$

where $\hat{f}(k_x, k_y, k_z)$ denotes the 3-D Fourier transform of $f(x, y, z)$. As a result, a projection image contains the same information as the central slice of the 3-D Fourier transform that is perpendicular to the direction of projection. Since the molecular orientations are limited to in-plane rotations on the 2-D substrate, which is itself tilted by an angle $\theta$, the corresponding Fourier slices fill the whole 3-D Fourier space except for the region within a double cone, whose axis coincides with the body rotation axis of the molecules. The double cone has an opening angle $2\theta$ and the region within the missing cone represents the inaccessible information of the molecular structure in the setting of RCT.

## 2.3. Micrograph simulation

Before discussing our model for simulating micrographs, we first consider the computation of the molecular projection images. Let $F$ be the discretization of the molecular structure $f$ that is defined on a cubic grid $(x, y, z) \in \{-2r, \ldots, 2r\}^3$ by

$$F(x, y, z) = f(x/\xi, y/\xi, z/\xi). \tag{7}$$

The integer $r$ represents the radius of a spherical support such that $F(x, y, z)$ is negligible for $(x^2 + y^2 + z^2)^{1/2} \geq r$. In addition, we define the discretization of the molecular projection $\mathcal{I}_\alpha^\theta$ by

$$I_\alpha^\theta(x, y) = \mathcal{I}_\alpha^\theta(x/\xi, y/\xi), \tag{8}$$

where $(x, y) \in \{-2r, \ldots, 2r\}^2$, and it immediately follows that $I_\alpha^\theta$ has a circular support of radius $r$. From the Fourier slice theorem, we can compute the discrete Fourier transform (DFT) of $I_\alpha^\theta$ from the DFT of $F$ by

$$\hat{I}_\alpha^\theta(k_x, k_y) = \hat{F}((R_\alpha^\theta)^T [k_x, k_y, 0]^T), \tag{9}$$

where $(k_x, k_y) \in \{-2r, \ldots, 2r\}^2$. To reduce the interpolation error, we use the FIN-UFFT package (Barnett *et al.*, 2019; Barnett, 2021) to evaluate $\hat{F}$ on the non-Cartesian grid points. Finally, we obtain the molecular projections $I_\alpha^\theta$ by the inverse DFT of $\hat{I}_\alpha^\theta$.

We simulate the micrographs measured in a RCT experiment at the substrate tilt angle $\theta$ by

$$M(x_i, y_i) = \sum_{j=1}^{n_p} I_{\alpha_j}^\theta(x_i - t_{x_j}, y_i - t_{y_j}) + \varepsilon(x_i, y_i), \tag{10}$$

where $(x_i, y_i) \in \{-\lfloor m/2 \rfloor, \ldots, \lceil m/2 - 1 \rceil\}^2$, $n_p$ is the number of molecular projections in the micrograph, $\alpha_j$ is the in-plane rotation of the $j^{\text{th}}$ molecule that is uniformly drawn from $[0, 2\pi)$, $(t_{x_j}, t_{y_j}) \in \{-\lfloor m/2 \rfloor + r, \ldots, \lceil m/2 - 1 \rceil - r\}^2$ is the center of the tomographic projection of the $j^{\text{th}}$ molecule, and $\varepsilon(x_i, y_i)$ is i.i.d. Gaussian noise with zero mean and variance $\sigma^2$. For a reason that will be clear in Section 2.4, we further

IUCr macros version 2.1.11: 2020/04/29

assume that

$$((t_{x_j} - t_{x_k})^2 + (t_{y_j} - t_{y_k})^2)^{1/2} > 4r \quad \text{for } j \neq k \tag{11}$$

such that the molecular projections are well separated in the micrographs. Figure 4 shows a sample micrograph with SNR = 1. We define SNR as the ratio of the mean squared pixel values of molecular projections to the noise variance. Specifically,

$$\text{SNR} = \frac{1}{2\pi} \int_0^{2\pi} d\alpha \; \frac{1}{\pi r^2} \sum_{x_i^2 + y_i^2 < r^2} |I_\alpha^\theta(x_i, y_i)|^2 \bigg/ \sigma^2. \tag{12}$$

*2.4. Autocorrelation analysis*

The standard data processing pipelines in single-particle cryo-EM start with the step of particle picking to locate the molecular projections in the noisy micrographs, which is equivalent to determining the 2-D vector $[t_{x_j}, t_{y_j}]^T$ for each molecular projection. This task, however, becomes challenging when the noise level is high. An alternative is to extract from the data quantities that are invariant to the 2-D translations of molecular projections in the micrographs. We achieve this through the approach of autocorrelation analysis.

Consider an $n \times n$ image $g(\mathbf{x})$. We define its autocorrelation function of order $q = 1, 2, \ldots$ for any 2-D translations $\mathbf{x_1}, \ldots, \mathbf{x_{q-1}} \in \mathbb{Z}^2$ by

$$a_g^q(\mathbf{x_1}, \ldots, \mathbf{x_{q-1}}) = \frac{1}{n^2} \sum_{\mathbf{x}} g(\mathbf{x})g(\mathbf{x} + \mathbf{x_1}) \cdots g(\mathbf{x} + \mathbf{x_{q-1}}), \tag{13}$$

where $\mathbf{x} \in \{-\lfloor n/2 \rfloor, \ldots, \lceil n/2 - 1 \rceil\}^2$ and $g(\mathbf{x})$ is zero-padded for arguments out of the range. In the context of this study, we set $n = m$ when $g$ represents a micrograph $M$ and $n = 4r + 1$ when $g$ represents a molecular projection $I_\alpha^\theta$.

Under the assumption that the molecular projections are well separated, as in (11), the autocorrelations of a micrograph with 2-D translations $\mathbf{x_1}, \ldots, \mathbf{x_{q-1}}$, where $|\mathbf{x_1}|, \ldots, |\mathbf{x_{q-1}}| \leq 2r$, are insensitive to the locations of molecular projections in the micrograph. As a result, the micrograph autocorrelations can be directly related to

the autocorrelations of molecular projections, which provide information about the molecular structure.

In this work, we consider the micrograph autocorrelations up to the third order. This choice is based on the number of equations provided by the autocorrelations. The first, second and third order autocorrelations provide $1, O(r^2)$ and $O(r^4)$ equations respectively. Our goal is to estimate the $O(r^3)$ voxel values of the molecular structure $F$. Hence, we need to go to at least the third order. Using autocorrelations of even higher orders may provide additional information about $F$, but it also requires more data and computational resources to accurately estimate their values.

Under the additional assumption that the density of molecular projections $\gamma = n_p(4r+1)^2/m^2$ is fixed, it is straightforward to show that (see for example in Bendory $et\ al.$ (2018))

$$\mathbb{E}\{a_M^1\} = \gamma \ \langle a_{I_\alpha^\theta}^1 \rangle_\alpha \tag{14}$$

$$\mathbb{E}\{a_M^2(\mathbf{x_1})\} = \gamma \ \langle a_{I_\alpha^\theta}^2(\mathbf{x_1}) \rangle_\alpha + \sigma^2 \delta(\mathbf{x_1}) \tag{15}$$

$$\mathbb{E}\{a_M^3(\mathbf{x_1}, \mathbf{x_2})\} = \gamma \ \langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2}) \rangle_\alpha$$
$$+ \ \gamma \ \langle a_{I_\alpha^\theta}^1 \rangle_\alpha \ \sigma^2 \big( \delta(\mathbf{x_1}) + \delta(\mathbf{x_2}) + \delta(\mathbf{x_1} - \mathbf{x_2}) \big) \tag{16}$$

for any fixed level of noise and $|\mathbf{x_1}|, |\mathbf{x_2}| \leq 2r$. Here $\mathbb{E}\{\cdot\}$ represents the expectation over the distributions of the random Gaussian noise and the in-plane rotations of molecules, and $\langle \cdot \rangle_\alpha$ denotes the angular average over $\alpha \in [0, 2\pi)$. The delta functions, defined by $\delta(0) = 1$ and $\delta(\mathbf{x} \neq 0) = 0$, are due to the autocorrelations of the random Gaussian noise.

We estimate the expectations in (14)-(16) by averaging autocorrelations computed from many micrographs. In practice, $\sigma^2$ and $\gamma \ \langle a_{I_\alpha^\theta}^1 \rangle_\alpha$ can be estimated from the micrographs: $\sigma^2$ can be estimated by the variance of micrograph pixel values in the low SNR regime; $\gamma \ \langle a_{I_\alpha^\theta}^1 \rangle_\alpha$ can be estimated by the empirical mean of micrographs. As a result, we can estimate the autocorrelations $\langle a_{I_\alpha^\theta}^1 \rangle_\alpha$, $\langle a_{I_\alpha^\theta}^2(\mathbf{x_1}) \rangle_\alpha$ and $\langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2}) \rangle_\alpha$

up to the constant factor $\gamma$. For simplicity, we assume that $\sigma^2$ and $\gamma$ are known to us.

*2.5. Regularized optimization*

In this section, we design an optimization problem to reconstruct the molecular structure $F$ from autocorrelations. We start by expressing $F$ in a non-redundant representation. Recall that $F$ is defined on a cubic grid of size $4r+1$ and has a spherical support of radius $r$. We represent $F$ by a vector $\mathbf{u}$ of length $n_r$, where $n_r$ denotes the number of voxels within the support. Furthermore, we define the linear operator $\mathcal{A}$ that maps $\mathbf{u}$ to $F$ by $F = \mathcal{A}\mathbf{u}$.

In our optimization problem, we estimate $\mathbf{u}$ by fitting the rotationally averaged 3rd order autocorrelation $\langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2})\rangle_\alpha$. As will be seen later, $\langle a_{I_\alpha^\theta}^1\rangle_\alpha$ is used to generate the initial guess for $\mathbf{u}$, and $\langle a_{I_\alpha^\theta}^2(\mathbf{x_1})\rangle_\alpha$ is used to build the regularizer in the optimization. For computational efficiency, we construct the cost function with the DFT of $\langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2})\rangle_\alpha$, where

$$
\begin{aligned}
s_F^3(\mathbf{k_1}, \mathbf{k_2}) &= \mathcal{F}\{\langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2})\rangle_\alpha\}(\mathbf{k_1}, \mathbf{k_2}) \\
&= \frac{1}{2\pi} \int_0^{2\pi} d\alpha \; \hat{I}_\alpha^\theta(\mathbf{k_1})\hat{I}_\alpha^\theta(\mathbf{k_2})\hat{I}_\alpha^{\theta^*}(\mathbf{k_1} + \mathbf{k_2}) \\
&\approx \frac{1}{n_\alpha} \sum_{i=0}^{n_\alpha - 1} \hat{I}_{\alpha_i}^\theta(\mathbf{k_1})\hat{I}_{\alpha_i}^\theta(\mathbf{k_2})\hat{I}_{\alpha_i}^{\theta^*}(\mathbf{k_1} + \mathbf{k_2}),
\end{aligned}
\tag{17}
$$

where $*$ denotes the complex conjugate. In the last step, we replace the integration with a discrete sum over $n_\alpha$ samples, where $\alpha_i = 2\pi i/n_\alpha$.

The triple product in (17) is the Fourier transform of the 3rd order autocorrelation $a_{I_{\alpha_i}^\theta}^3(\mathbf{x_1}, \mathbf{x_2})$, also known as the bispectrum (Tukey, 1953). Its applications in signal processing can be seen, for instance, in Sadler & Giannakis (1992) and Bendory *et al.* (2017). Since we assume that the information of the molecular projections is preserved only up to the Nyquist frequency due to noise, we only consider spatial frequencies $(\mathbf{k_1}, \mathbf{k_2}) \in \mathcal{V}$, where $\mathcal{V} = \{(\mathbf{k_1}, \mathbf{k_2}) : |\mathbf{k_1}|, |\mathbf{k_2}|, |\mathbf{k_1} + \mathbf{k_2}| < 2r\}$. Let $\tilde{s}_F^3(\mathbf{k_1}, \mathbf{k_2})$ be the

DFT of the estimation of $\langle a_{I_\alpha^\theta}^3(\mathbf{x_1}, \mathbf{x_2})\rangle_\alpha$ from data. We can hence express the sum of least-square errors by $\sum_{(\mathbf{k_1}, \mathbf{k_2}) \in \mathcal{V}} |s_F^3(\mathbf{k_1}, \mathbf{k_2}) - \tilde{s}_F^3(\mathbf{k_1}, \mathbf{k_2})|^2$.

As discussed in Section 2.2, there exists a double-cone region in the Fourier space that cannot be probed in RCT. Therefore, our reconstruction problem is ill-posed in nature, and we must include a regularization term in the cost function to incorporate some prior knowledge of the true solution. Our regularization enforces the smoothness assumption on $F$ and has the form of the weighted sum of squares: $\sum_{\mathbf{q}} |\hat{F}(\mathbf{q})|^2/\tau(\mathbf{q})^2$, where $\mathbf{q} \in \{-2r, \ldots, 2r\}^3$. This regularization is related to the Gaussian prior described in Scheres (2012) in that we expect the scale parameters $\tau(\mathbf{q})^2$ to act as a low-pass filter to reduce high-frequency noise while still preserve some high-resolution features of the molecule.

We estimate the values of $\tau(\mathbf{q})$ based on the observation that the structure factors of proteins obey Wilson statistics (Wilson, 1949). To be more precise, the structure factors within each resolution shell follow the complex normal distribution with mean zero and variance estimated from the mean intensity in the resolution shell (French & Wilson, 1978). Taking the DFT of $\langle a_{I_\alpha^\theta}^2(\mathbf{x_1})\rangle_\alpha$, we obtain

$$
\begin{aligned}
s_F^2(\mathbf{k_1}) = \mathcal{F}\{\langle a_{I_\alpha^\theta}^2(\mathbf{x_1})\rangle_\alpha\}(\mathbf{k_1}) &= \frac{1}{2\pi} \int_0^{2\pi} d\alpha \ \hat{I}_\alpha^\theta(\mathbf{k_1}) \hat{I}_\alpha^{\theta^*}(\mathbf{k_1}) \\
&= \frac{1}{2\pi} \int_0^{2\pi} d\alpha \ |\hat{F}((R_\alpha^\theta)^T[k_{1x}, k_{1y}, 0]^T)|^2,
\end{aligned}
\tag{18}
$$

where $\mathbf{k_1} \in \{-2r, \ldots, 2r\}^2$ and we only consider spatial frequencies within the Nyquist frequency, that is, $|\mathbf{k_1}| < 2r$. Since

$$
(R_\alpha^\theta)^T \begin{bmatrix} k_{1x} \\ k_{1y} \\ 0 \end{bmatrix} = \begin{bmatrix} k_{1x} \cos\alpha + k_{1y} \cos\theta \sin\alpha \\ -k_{1x} \sin\alpha + k_{1y} \cos\theta \cos\alpha \\ -k_{1y} \sin\theta \end{bmatrix},
\tag{19}
$$

we can see that $s_F^2(\mathbf{k_1})$ is the mean intensity over a circle that is perpendicular to the body rotation axis of the molecule and has radius $(k_{1x}^2 + k_{1y}^2 \cos^2\theta)^{1/2}$ and height $-k_{1y} \sin\theta$. Therefore, with appropriate weights, the average of $s_F^2(\mathbf{k_1})$ for all $\mathbf{k_1}$ that

fall into the same annulus $q_{\min} < |\mathbf{k_1}| < q_{\max}$ gives the mean intensity within the resolution shell $q_{\min} < |\mathbf{q}| < q_{\max}$, excluding the spherical caps that lie in the missing-cone region. We represent this weighted average by $\tau_{|\mathbf{q}|}^2$, whose values are in practice computed from $\tilde{s}_F^2(\mathbf{k_1})$, the DFT of the estimation of $\langle a_{I_\alpha^\theta}^2(\mathbf{x_1}) \rangle_\alpha$ from data.

In addition to the scale parameters $\tau_{|\mathbf{q}|}^2$ for $|\mathbf{q}| < 2r$, it is helpful to have regularization outside the Nyquist frequency to limit high-frequency noise. We choose $\tau(\mathbf{q}) = |\mathbf{q}|^{-1}$ for $|\mathbf{q}| \geq 2r$. This choice is based on the identity

$$\int_{\mathbb{R}^3} |\nabla f|^2 \, d^3\mathbf{x} = \int_{\mathbb{R}^3} |\mathbf{q}|^2 |\hat{f}(\mathbf{q})|^2 \, d^3\mathbf{q} \tag{20}$$

such that one can minimize the sum of gradient squares by minimizing $\int_{\mathbb{R}^3} |\mathbf{q}|^2 |\hat{f}(\mathbf{q})|^2 \, d^3\mathbf{q}$.

Finally, we define the cost function of our optimization problem by

$$\begin{aligned}
C(\mathbf{u}) = &\sum_{(\mathbf{k_1},\mathbf{k_2}) \in \mathcal{V}} \left| s_F^3(\mathbf{k_1},\mathbf{k_2}) - \tilde{s}_F^3(\mathbf{k_1},\mathbf{k_2}) \right|^2 \\
&+ \lambda \left( \sum_{|\mathbf{q}| < 2r} \frac{|\hat{F}(\mathbf{q})|^2}{\tau_{|\mathbf{q}|}^2} + \beta \sum_{|\mathbf{q}| \geq 2r} |\mathbf{q}|^2 |\hat{F}(\mathbf{q})|^2 \right),
\end{aligned} \tag{21}$$

where $\mathbf{u}$ is the non-redundant representation of $F$, $\lambda$ denotes the regularization parameter, and we compute the scale factor $\beta$ such that the two curves $\tau_{|\mathbf{q}|}^2$ and $|\mathbf{q}|^{-2}$ attain the same value at $|\mathbf{q}| = 2r$. In a separate attempt, we have used $\sum_{\mathbf{q}} |\mathbf{q}|^2 |\hat{F}(\mathbf{q})|^2$ as the only regularizer, but the quality of the reconstruction appears to be inferior with significant high-frequency noise (not shown in this study). This result suggests that Wilson statistics is a reasonably good prior for the Fourier components.

The optimization problem shown in (21) is inherently nonconvex due to the term of non-linear least-square errors. We find that the BFGS algorithm, despite being a local search algorithm, works well on the problem when initialized at a reasonable guess. We initialize $\mathbf{u}$ from a (deterministic) 3-D Gaussian profile with variance $r^2$, which is rescaled such that the sum of its 3-D discretization is consistent with $\langle a_{I_\alpha^\theta}^1 \rangle_\alpha$ estimated from data. We run the BFGS algorithm in the tensorflow software library (Abadi *et al.*,

2015) to minimize (21) over a set of regularization parameters $\lambda = 10^{-2}, 10^{-1}, \ldots, 10^7$. From the converged solutions, we choose the optimal value of $\lambda$ using the L-curve method (Hansen, 1992). Our reconstructed structures are the estimates for $F$ with these optimal values of $\lambda$.

## 3. Results

### 3.1. Reconstruction at different substrate tilts

In this section, we explore the effect of the missing-cone region on the quality of reconstruction by considering micrographs measured at different substrate tilt angles $\theta = 60°, 35°$ and $10°$. The molecule used in our simulation is Bovine Pancreatic Trypsin Inhibitor (BPTI), which has size of 35 Å and weight of 6.5 kDa. This molecular size is substantially below the limit (40 kDa) believed to be attainable by single-particle cryo-EM (Henderson, 1995), and our model structure was determined using X-ray crystallography.

We generate the discrete molecular structure $F$ from the PDB entry 1QLQ (Czapinska *et al.*, 2000) using the UCSF Chimera software (Pettersen *et al.*, 2004) at a resolution of 5 Å. The resulting contrast has a spherical support of radius $r = 15$ voxels, and is further zero-padded to be a cubic grid of size 61. From the discrete contrast $F$, we simulate the micrographs as described in Section 2.3. To obtain the baseline results on the effect of the missing cone region, we consider the idealistic scenario that the in-plane rotation of the $j^{\text{th}}$ molecule is given by $\alpha_j = 2\pi j/n_p$, $j \in \{1, \ldots, n_p\}$, and the noise variance $\sigma^2 = 0$. By setting the micrograph length $m = 4096$ pixels and the number of molecules $n_p = 400$, we only simulate one micrograph at each given value of the substrate tilt angle.

From the simulated micrographs, we compute the rotationally averaged autocorrelations of molecular projections and the values of $\tilde{s}_F^3(\mathbf{k_1}, \mathbf{k_2})$ and $\tau_{|\mathbf{q}|}^2$. Figure 5 shows

the comparison of the mean intensities $\langle |\hat{F}(\mathbf{q})|^2 \rangle$ and the scale parameters $\tau_{|\mathbf{q}|}^2$ and $|\mathbf{q}|^{-2}/\beta$ for $\theta = 60°$. We first see that $\tau_{|\mathbf{q}|}^2$ provides a good estimate for $\langle |\hat{F}(\mathbf{q})|^2 \rangle$ up to the Nyquist frequency. On the other hand, the scale parameter $|\mathbf{q}|^{-2}/\beta$ is substantially greater than $\langle |\hat{F}(\mathbf{q})|^2 \rangle$ outside the Nyquist frequency, which may inevitably preserve some high-resolution noise in the reconstruction.

Figure 6(a) shows the comparison of our reconstructed BPTI structures with the ground truth used to simulate the micrographs. As expected, the visual quality of the reconstructions degrades when the sample tilt angle $\theta$ decreases, which results in a larger missing-data region. To assess the reconstructions in more detail, we plot the Fourier shell correlation (FSC) (Harauz & Van Heel, 1986) of the reconstructed structures with the ground truth in Figure 6(b). Although the reconstruction at $\theta = 35°$ correlates to the ground truth worse than the one at $\theta = 60°$, both of them have the same resolution as the ground truth (5 Å) according to the FSC $= 0.5$ criterion. Using the same criterion, the resolution of the reconstruction at $\theta = 10°$ is 8.3 Å.

### 3.2. Reconstruction from noisy micrographs

After having the baseline results for reconstructions from noiseless micrographs, we turn to test our approach on noisy micrographs. At the sample tilt angle $\theta = 60°$, we simulate 500 micrographs of size $m = 4096$ using the same discrete contrast $F$ for BPTI. We adjust the noise level such that the micrographs have SNR $= 1$. By maximizing the density of molecular projections while still preserving the requirement of well separation (11), the resulting micrographs contain $1.4 \times 10^6$ molecular projections in total.

From the noisy micrographs, we compute the estimates for the rotationally averaged autocorrelations of molecular projections. Figure 7(a) shows the reconstruction from these estimates along with the ground truth. The negative effect of noise on the

quality of the reconstruction can best be seen by comparing this reconstruction with its counterpart in Figure 6(a). As plotted in Figure 7(b), we determine the resolution of this reconstructed structure to be 6.5 Å using the FSC = 0.5 criterion.

To demonstrate that our approach applies to other biological molecules, we test our approach on another dataset simulated from the myoglobin molecule, which has size of 40 Å and weight of 17.8 kDa. We generate the discrete molecular structure $F$ for myoglobin from the PDB entry 1MBN (Watson, 1969) using the UCSF Chimera software at a resolution of 5 Å. The resulting contrast has a spherical support of radius $r = 16$ voxels, and is further zero-padded to be a cubic grid of size 65. At the sample tilt angle $\theta = 60°$, we generate 500 micrographs of size $m = 4096$ from $F$. The number of molecular projections in these micrographs totals $1.2 \times 10^6$, and we also set SNR = 1 for the micrographs. The reconstructed myoglobin structure from the noisy micrographs is shown in Figure 8(a) along with the ground truth. We can see that our reconstruction recovers most of the main features of the ground truth. We plot the FSC of our reconstruction with the ground truth in Figure 8(b), and we determine the resolution of the reconstruction to be 7.0 Å according to the FSC = 0.5 criterion.

## 4. Discussion

In this paper, we present a method to reconstruct the 3-D molecular structure from data collected at just one sample tilt angle in RCT. Our method reduces data to quantities that are invariant to the 2-D positions of molecular projections in the micrographs, which removes the need for particle picking when analyzing data. In order to address the missing data in the double-cone region of the molecule's Fourier transform, we design a regularized optimization problem to reconstruct the molecular structure by fitting the autocorrelations estimated from micrographs. Our numerical studies illustrate the effect of the missing-cone region on the quality of reconstruction.

In addition, we demonstrate structure reconstruction from the autocorrelations computed from noisy micrographs. Since the accuracy of the autocorrelation estimates can be improved by averaging many more micrographs, our results show promise of applying autocorrelation analysis to reconstruct the structures of small biological molecules in the setting of RCT.

A few issues still stand in the way of applying our approach to real RCT data. In Section 2.1, we make the assumption that the point spread function is a 2-D Dirac delta function to ignore its effect. In reality, however, we may have to consider a varying point spread function with respect to the locations on the detector because different regions of the tilted specimen are exposed to the electron beam with different defocus values. Another challenge arises when the noise is colored. In that case, the expectations of products of noise at different pixels are not zero. It will require a more sophisticated model for the noise power spectrum instead of a single parameter $\sigma^2$. Furthermore, structure heterogeneity of the target molecule will be another test for our approach.

Additionally, we assume that the molecular projections are well separated in the micrographs. This assumption enables us to directly relate the micrograph autocorrelations to the autocorrelations of molecular projections. However, it is preferable in practice to have the molecular projections densely packed in micrographs to maximize the available structural information within limited data collection time. We expect to remove this assumption by considering the cross correlations between neighboring molecular projections. A similar idea was recently demonstrated in Lan *et al.* (2020) for the simplified 1-D model.

Another practical concern is the amount of required data. As a proof of concept, we reconstruct the molecular structures from simulated micrographs with SNR = 1. For small biological molecules that challenges particle picking, we expect the SNR of

the micrographs to be much lower. Since our approach uses autocorrelations up to the $3^{\text{rd}}$ order, the sample complexity would scale as $\text{SNR}^{-3}$. This means that we will need $10^3$ times more molecules to estimate the autocorrelations with similar accuracy when the SNR drops from 1 to 0.1. We plan to address this concern in the following ways: First, the third order correlations contain a large degree of redundancy, as they are 4-D functions containing information from a 3-D structure. Ideally, with proper denoising, for example, the Noise2Noise scheme (Lehtinen *et al.*, 2018), the SNR of the correlation function could likely be enhanced by this redundancy factor, ranging from 10 to 100 depending on the resolution. Second, the SNR of the correlation function is proportional to the density of molecular projections present in the micrographs. By enabling the reconstruction from micrographs of densely packed projections, we can boost the SNR of the correlation function by another factor of 10.

In the long run, we would like to extend the approach described here to real cryo-EM data to reconstruct high-resolution structures directly from micrographs, without being restricted to molecules which have a preferred orientation on their substrate.

## 5. Acknowledgements

## Appendix A
## Computational Details

The data simulation and structure reconstruction were performed on an Nvidia Tesla P100 GPU, which has 16 GB RAM. The computation of the micrograph auto-correlations for relevant step sizes took $1.5 \times 10^2$ seconds on average for a $4096 \times 4096$ micrograph. As for the structure reconstruction, it took a few hours for an instance with a given value of the regularization parameter $\lambda$ to converge. Therefore, if one knows the correct $\lambda$ for some setting, it may be advantageous to use the same $\lambda$ in a similar case. The code is publicly available at https://github.com/tl578/RCT-without-detection.

## References

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y. & Zheng, X., (2015). Software available from tensorflow.org.

von Ardenne, B., Mechelke, M. & Grubmüller, H. (2018). *Nature communications*, **9**, 2375.

Barnett, A. H. (2021). *Applied and Computational Harmonic Analysis*, **51**, 1–16.

Barnett, A. H., Magland, J. & af Klinteberg, L. (2019). *SIAM Journal on Scientific Computing*, **41**(5), C479–C504.

Bendory, T., Boumal, N., Leeb, W., Levin, E. & Singer, A. (2018). *arXiv preprint arXiv:1810.00226*.

Bendory, T., Boumal, N., Leeb, W., Levin, E. & Singer, A. (2019). *Inverse Problems*, **35**(10), 104003.

Bendory, T., Boumal, N., Ma, C., Zhao, Z. & Singer, A. (2017). *IEEE Transactions on Signal Processing*, **66**(4), 1037–1050.

Bendory, T., Lan, T.-Y., Marshall, N. F., Rukshin, I. & Singer, A. (2021). *arXiv preprint arXiv:2101.07709*.

Czapinska, H., Otlewski, J., Krzywda, S., Sheldrick, G. & Jaskólski, J. (2000). *Journal of Molecular Biology*, **295**(5), 1237–1249.

Donatelli, J. J., Zwart, P. H. & Sethian, J. A. (2015). *Proceedings of the National Academy of Sciences*, **112**(33), 10286–10291.

Elser, V. (2011). *New Journal of Physics*, **13**, 123014.

French, S. & Wilson, K. (1978). *Acta Crystallographica Section A*, **34**(4), 517–525.

Hansen, P. C. (1992). *SIAM Review*, **34**(4), 561–580.

Harauz, G. & Van Heel, M. (1986). *Optik (Stuttgart)*, **73**(4), 146–156.

Henderson, R. (1995). *Quarterly Reviews of Biophysics*, **28**(2), 171–193.

Kam, Z. (1977). *Macromolecules*, **10**(5), 927–934.

Kam, Z. (1980). *Journal of Theoretical Biology*, **82**(1), 15–39.

Kurta, R. P., Donatelli, J. J., Yoon, C. H., Berntsen, P., Bielecki, J., Daurer, B. J., DeMirci, H., Fromme, P., Hantke, M. F., Maia, F. R. N. C., Munke, A., Nettelblad, C., Pande, K., Reddy, H. K. N., Sellberg, J. A., Sierra, R. G., Svenda, M., van der Schot, G., Vartanyants, I. A., Williams, G. J., Xavier, P. L., Aquila, A., Zwart, P. H. & Mancuso, A. P. (2017). *Physical Review Letter*, **119**, 158102.

Lan, T.-Y., Bendory, T., Boumal, N. & Singer, A. (2020). *IEEE Transactions on Signal Processing*, **68**, 1589–1601.

Lehtinen, J., Munkberg, J., Hasselgren, J., Laine, S., Karras, T., Aittala, M. & Aila, T. (2018). *Proceedings of the 35th international conference on machine learning (ICML)*, p. 2965–2974.

Marshall, N., Lan, T.-Y., Bendory, T. & Singer, A. (2020). In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5780–5784.

Pande, K., Donatelli, J. J., Malmerberg, E., Foucar, L., Bostedt, C., Schlichting, I. & Zwart, P. H. (2018). *Proceedings of the National Academy of Sciences*, **115**(46), 11772–11777.

Pettersen, E. F., Goddard, T. D., Huang, C. C., Couch, G. S., Greenblatt, D. M., Meng, E. C. & Ferrin, T. E. (2004). *Journal of Computational Chemistry*, **25**(13), 1605–1612.

Radermacher, M. (1988). *Journal of Electron Microscopy Technique*, **9**(4), 359–394.

Radermacher, M., Wagenknecht, T., Verschoor, A. & Frank, J. (1987). *Journal of Microscopy*, **146**(2), 113–136.

Sadler, B. M. & Giannakis, G. B. (1992). *Journal of the Optical Society of America A*, **9**(1), 57–69.

Saldin, D. K., Poon, H. C., Shneerson, V. L., Howells, M., Chapman, H. N., Kirian, R. A., Schmidt, K. E. & Spence, J. C. H. (2010). *Physical Review B*, **81**, 174105.

Scheres, S. (2012). *Journal of Molecular Biology*, **415**(2), 406–418.

Sorzano, C. O. S., Alcorlo, M., de la Rosa-Trevín, J. M., Melero, R., Foche, I., Zaldívar-Peraza, A., del Cano, L., Vargas, J., Abrishami, V., Otón, J., Marabini, R. & Carazo, J. M. (2015). *Scientific Reports*, **5**(14290).

Tukey, J. (1953). *Reprinted in The Collected Works of John W. Tukey*, **1**, 165–184.

Watson, H. (1969). *Progress in Stereochemistry*, **4**, 299–333.

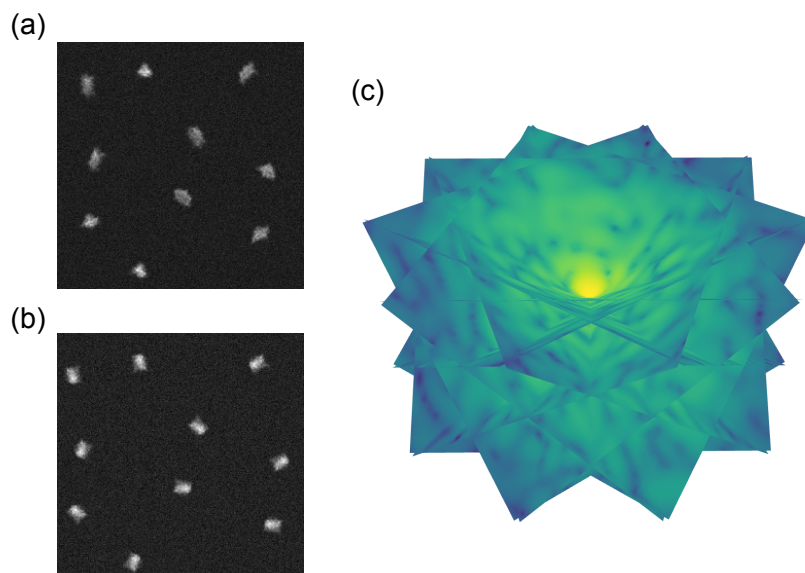Wilson, A. (1949). *Acta Crystallographica*, **2**(5), 318–321.

Fig. 1. The micrographs of the same field of view collected at (a) one large sample tilt angle and (b) no tilt. (c) The Fourier transforms of the molecular projections recorded in (a), which are assembled in Fourier space with respect to their corresponding orientations according to the Fourier slice theorem discussed in Section 2.2.
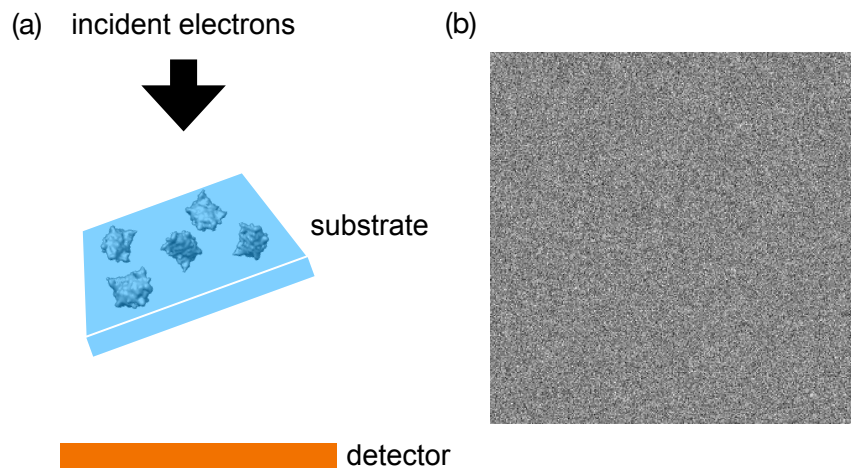
Fig. 2. (a) The data collection scheme of RCT with just one sample tilt angle. (b) A micrograph that is so noisy that picking particles is challenging.
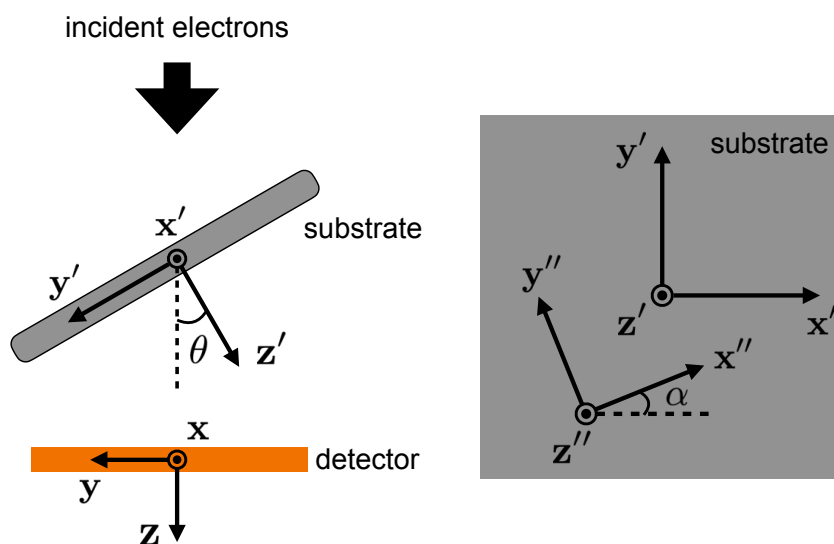


Fig. 3. The relationships between the lab frame $S$, the frame fixed on the 2-D substrate $S'$ and the body frame of one particular molecule $S''$.
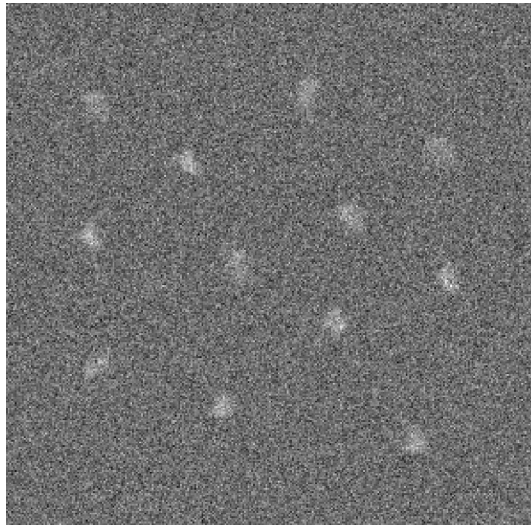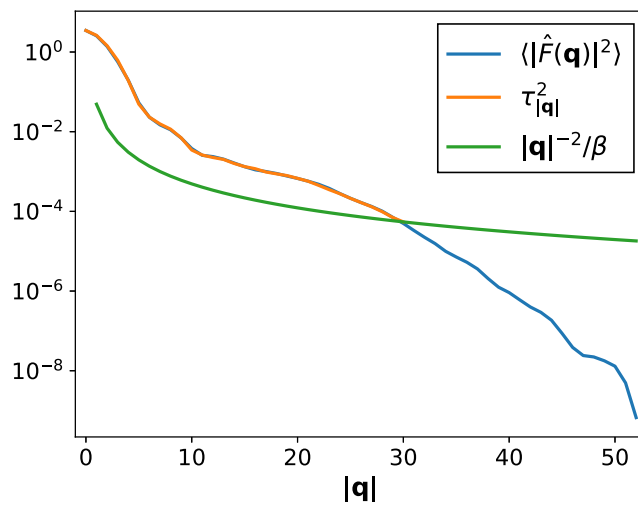
Fig. 4. A sample micrograph with SNR = 1.



Fig. 5. The comparison of the mean intensities $\langle |\hat{F}(\mathbf{q})|^2 \rangle$ and the scale parameters $\tau^2_{|\mathbf{q}|}$ and $|\mathbf{q}|^{-2}/\beta$ for the BPTI molecule at the substate tilt angle $\theta = 60°$.
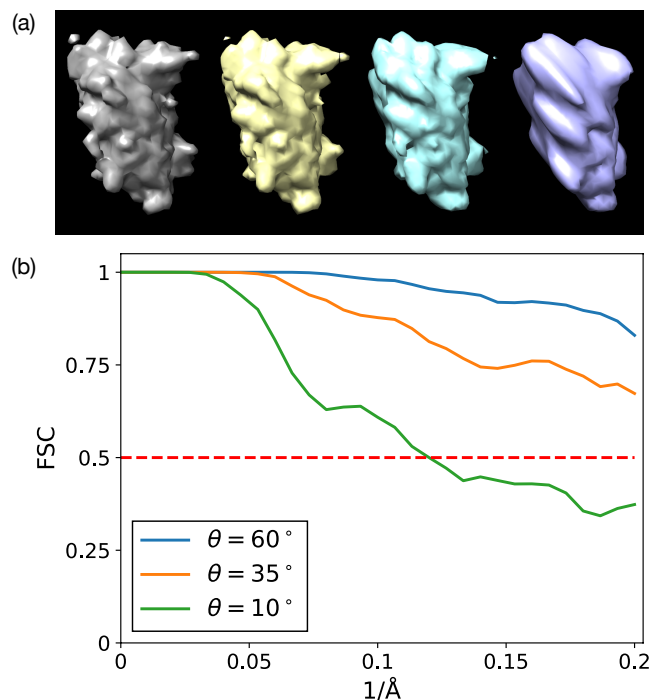
Fig. 6. (a) The reconstructed BPTI structures from noiseless micrographs at different sample tilt angles: $\theta = 60°$ (yellow), $\theta = 35°$ (cyan) and $\theta = 10°$ (purple). The grey one is the ground truth used to simulate the micrographs. (b) The FSC of the reconstructed structures with the ground truth.



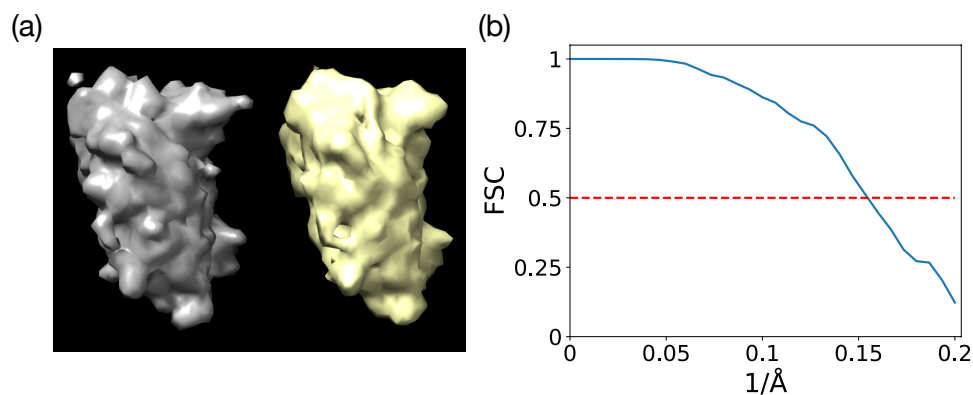Fig. 7. (a) The reconstructed BPTI structure (yellow) from noisy micrographs with SNR = 1 at the sample tilt angle $\theta = 60°$. The ground truth is rendered in grey. (b) The FSC of the reconstructed structure with the ground truth.

(a)
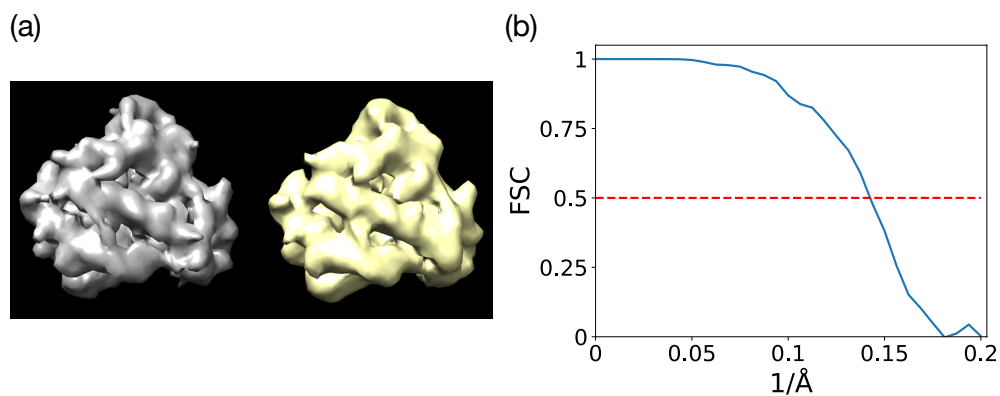
(b)



Fig. 8. (a) The reconstructed myoglobin structure (yellow) from noisy micrographs with SNR = 1 at the sample tilt angle $\theta = 60°$. The ground truth is rendered in grey. (b) The FSC of the reconstructed structure with the ground truth.

**Synopsis**

We describe a method to reconstruct the 3-D molecular structure without the need for particle picking in the random conical tilt scheme in cryo-electron microscopy. Our results show promise to reduce the size limit for single particle reconstruction in cryo-electron microscopy.