

1 **Non-Uniformity of Projection Distributions Attenuates** 2 **Resolution in Cryo-EM**

3
4 Philip R. Baldwin¹ and Dmitry Lyumkis¹

5
6 ¹The Salk Institute for Biological Studies, La Jolla, CA 92037, USA

7
8 Correspondence: pbaldwin@salk.edu, dlyumkis@salk.edu

9
10
11 **Keywords** - Fourier Shell Correlation, single particle analysis, preferred orientation, anisotropy

12

1 **Abstract**

2 Virtually every single-particle cryo-EM experiment currently suffers from specimen adherence to
3 the air-water interface, leading to a non-uniform distribution in the set of projection views.
4 Whereas it is well accepted that uniform projection distributions can lead to high-resolution
5 reconstructions, non-uniform (anisotropic) distributions can negatively affect map quality,
6 elongate structural features, and in some cases, prohibit interpretation altogether. Although some
7 consequences of non-uniform sampling have been described qualitatively, we know little about
8 how sampling quantitatively affects resolution in cryo-EM, especially given the numerous
9 different projection schemes that can arise in experimental situations. Here, we show how
10 inhomogeneity in any projection distribution scheme attenuates the global Fourier Shell
11 Correlation (FSC) in relation to the number of particles and a single geometrical parameter, which
12 we term the sampling compensation factor (SCF). The reciprocal of the SCF is defined as the
13 average over Fourier shells of the reciprocal of the per-particle sampling and normalized to unity
14 for uniform distributions. The SCF therefore ranges from one to zero, with values close to the latter
15 implying large regions of poorly sampled or completely missing data in Fourier space. Using two
16 synthetic test cases, influenza hemagglutinin and human apoferritin, we demonstrate how any
17 amount of sampling inhomogeneity always attenuates the FSC compared to a uniform distribution.
18 We advocate quantitative evaluation of the SCF criterion to approximate the effect of non-uniform
19 sampling on resolution within experimental single-particle cryo-EM reconstructions.

20

1	Introduction
2	Section 1. Summary of the major findings
3	Section 2. Decrement of SSNR due to sampling inhomogeneity
4	2.1. Derivation of the Sampling Compensation Factor (SCF)
5	2.2. An adjusted formula for SSNR for half maps with unmeasured data
6	Section 3. Numerical and analytical forms for the sampling function, and expressions for
7	the SCF geometrical factor
8	3.1 Discrete and continuum approaches to the sampling function
9	3.1.1. Discrete treatment for sampling
10	3.1.2. A continuum treatment for sampling
11	3.1.3. Consistency between numerical and analytical expressions for
12	sampling
13	3.2 Sampling Function for three different distributions in continuum representation
14	3.2.1. Side-like cases
15	3.2.2. Modulated side-like cases
16	3.2.3. Top-like cases
17	3.3. The Sampling Compensation Factors for the three different distributions.
18	Section 4. Relationship between SSNR and the number of particles N in a reconstruction
19	4.1. Linear dependence of SSNR on N
20	4.2. Number of Particles Necessary for Reconstruction
21	4.3. Comparison of graphical methods (Guinier, Reslog, and per-particle SSNR
22	curve)
23	Section 5. Decrement of SSNR through non-uniform sampling

1	5.1. Methods
2	5.1.1 Generation of projection distributions
3	5.1.2 Synthetic data generation with distinct projection distributions
4	5.2. Results comparing decrements in SSNR predicted by sampling non-uniformity
5	5.2.1. Side-like and side-modulated sampling cases
6	5.2.2. Top-like sampling cases for varying cone sizes
7	5.2.3. Top-like sampling cases for fixed cone size and varying fraction of
8	randomly sprinkled projections
9	Discussion
10	

1 Glossary

- 2 FSC(k) is the Fourier shell correlation of the reconstruction at Fourier frequency k .
3 SSNR(k) is the spectral signal to noise ratio of the reconstruction at Fourier frequency k .
4 $ssnr(k)$ is the per particle SSNR, used in the discussion in Section 4.
5 L , the side of the real space box
6 N , the number of particles in the reconstruction
7 k is Fourier magnitude
8 \vec{k} is a 2D or 3D point in Fourier space
9 SCF is the sampling compensation factor, characterizes effect of sampling on SSNR
10 $\mathcal{N}(k)$ is noise-to-signal power (sections 1,2)
11 $Sp(\vec{k}), sp(\vec{k})$ is the sampling function (and per particle sampling function) defined at a 3D
12 lattice site \vec{k}
13 $F_j(\vec{k})$ is the Fourier value of the j^{th} projection at the point \vec{k} .
14 $M_j(k)$ is the effect of the microscope (CTF) on the j^{th} projection
15 $X(\vec{k}), \hat{X}(\vec{k})$ is the target model and a running estimate of the model
16 R_j is a 3D rotation matrix describing the projection, j .
17 $N_j(\vec{k})$ is the noise added to the projection, j .
18 $E(k)$ is the total envelope that attenuate the image due to microscope and misalignment
19 effects.
20 $\hat{N}(\vec{k})$ is the effective noise at 3D lattice sites after regrouping from projections
21 $N_2(k) \equiv \langle |N(\vec{k})|^2 \rangle$ is the power of the noise
22 F, G are half maps used to derive FSC relations in section 2
23 P_k, Q_k the number of measured and unmeasured voxels, on a Fourier shell of radius k , when
24 there is missing data.
25 \hat{n} is used as a unit vector demarcating a projection
26 $\Theta(x)$ indicator function, which is 1 if the condition x is true, 0 otherwise
27 λ is the amplitude of the modulation for the modulation of side views (section 3)
28 α is a cone half angle: for top-like views, projections are inside cone;
29 for side like, projections are outside the cone.
30 ϵ is the fraction of projections that are not restricted to be in the main cone of half angle α
31 (Section 3)
32 Euler Angle is one of the three angles used to describe rotation matrices (θ is rotation around Z-
33 axis, ϕ is rotation around Y-axis, ψ is in-plane rotation).
34

1 **Introduction**

2

3 Single-particle cryo-electron microscopy (cryo-EM) has gained increasing popularity for
4 structural analysis of macromolecules and macromolecular assemblies. Numerous technical
5 advances have contributed to improvements in resolution [1-3], throughput [4], and overall
6 usability of the approaches, leading to a wealth of novel insights pertaining to macromolecular
7 structure and function [5]. Although many steps in the single-particle workflow are becoming more
8 streamlined and automated, a principal remaining challenge pertains to problems resulting from
9 non-uniform projection distributions contributing to reconstructed density maps.

10

11 Non-uniformity in the distribution of projection orientations recorded in a single-particle imaging
12 experiment originates from adherence of the specimen to one of two interfaces (top or bottom) of
13 the grid. The interfaces, which could be air-water or support-water (e.g. thin carbon), cause
14 specimens to stick in one of several “preferential orientations”. It is now clear that virtually every
15 specimen prepared for single-particle imaging using conventional blotting techniques adopts a
16 preferential orientation on cryo-EM grids [6]. The reason for this is that macromolecules, which
17 continuously undergo rapid thermal motion, adhere to interfaces on a time scale that is orders of
18 magnitude shorter than the time to blot off excess sample. Recent inkjet dispensing technologies
19 have ameliorated some of the effects of preferential specimen orientation by attempting to out-run
20 sample adherence to interfaces and by minimizing the amount of time between sample application
21 and plunging into liquid ethane [7]. However, such devices do not yet eliminate preferential
22 orientation in its entirety and depend heavily on high sample concentration. Furthermore, the
23 increase in interest in specimen supports, like graphene [8, 9], which also cause preferential

1 orientation, indicates that the effects of non-uniform sampling on final reconstructions will remain
2 problematic for many single-particle experiments.

3

4 Numerous approaches have been devised to estimate the quality of angular distributions and their
5 effects on a reconstructed density. These ideas are primarily developed in conjunction with some
6 anisotropic measure. One measure derives from the application of a 3D point spread function to
7 estimate the strength of signal above some significance criterion, in all directions of the 3D Fourier
8 transform [10]. In another approach, the 3D spectral signal-to noise ratio (SSNR) is used to define
9 directional resolution differences [11], with the SSNR bearing a direct relationship to the Fourier
10 Shell Correlation (FSC), the conventional means for measuring resolution in single-particle cryo-
11 EM. Multiple groups also described the use of conical FSCs to evaluate anisotropic resolution for
12 tomographic reconstructions [12, 13], as well as our and others' work on evaluating anisotropic
13 resolution in single-particle analysis [14, 15]. More recently, the “efficiency” metric [16] was
14 introduced to characterize an orientation distribution, based on the observed relationship between
15 orientation distribution and experimental resolution. We proposed that an evaluation of anisotropy
16 in cryo-EM experiments should be standard for every cryo-EM reconstruction [17].

17

18 The consequences of sampling non-uniformity on a reconstructed density map can vary and
19 depend on the extent and distribution of projection views. In many experimental cases, one might
20 see a few distinct preferential orientations across the Euler distribution profile, but the resulting
21 map may look reasonable, and is readily interpretable with an atomic model. In the more severe
22 cases, an anisotropic distribution may lead to apparent elongation of structural features within the
23 map. In such cases, the interpretation of the map may be affected, sometimes severely, due to the

1 appearance of artefactual density parallel to the dominant view [15]. In the most severe cases,
2 structure determination may be stifled altogether. Some hallmarks of pathologically anisotropic
3 distributions include inflated Fourier Shell Correlation (FSC) curves, elongated features beyond
4 interpretability, an inability to converge on a final structure, and/or the appearance of false positive
5 orientations in the course of refinement [15]. All these factors can reinforce problems in the
6 density. One interesting observation was that anisotropic orientation distributions lead to an
7 increase in the temperature factor associated with the data, thereby also affecting global resolution
8 [16]. However, a derivation from standard models has not been established.

9
10 While different measures have been introduced to evaluate the effect of anisotropic distributions
11 on directional viewings of the reconstructed density map, the effect of sampling on global
12 resolution has largely been neglected. Furthermore, there remains no systematic, quantitative study
13 of the effects of inhomogeneous projection distributions on cryo-EM reconstructions. Here, we
14 examine the relationship between non-uniform angular sampling and global resolution, as
15 measured using conventional analyses in cryo-EM. A major conclusion from our work is that *any*
16 inhomogeneity, and especially missing information in Fourier space, directly attenuates global
17 resolution in 3D reconstructions, and thus impedes the single-particle experiment.

18

19

20

21

1 Section 1. Summary of the major findings

2
3 Given a set of projection views, we develop an assessment of the quality of the sampling. We
4 chose this assessment based on the expected effect on the spectral signal to noise ratio (SSNR)
5 defined through the FSC. We show that the angular average of the reciprocal of the sampling forms
6 a quantity whose reciprocal attenuates the SSNR, if we consider the other aspects of the problem
7 associated with the overall experimental envelope to be held constant. More specifically, we argue:

$$9 \quad \text{SSNR}(k) = N \frac{\text{SCF}}{2k} \frac{1}{\mathcal{N}(k)}, \quad (1.1)$$

10
11 where N is the number of particles, k is spatial frequency, $\mathcal{N}(k)$ is a noise-to-signal power, and
12 SCF is what we term the “sampling compensation factor” and is defined to be

$$14 \quad \text{SCF} \equiv \frac{1}{\langle 1/(2k \text{ sp}) \rangle}. \quad (1.2)$$

15
16 Here, $\langle \cdot \rangle$ means the average in Fourier space over the nonzero values of shells at (approximately)
17 fixed spatial frequency, k , and sp is the amount of sampling per-particle, determined from the Euler
18 angle assignments. Notably, one must compensate for the geometry of the sampling to correctly
19 estimate the SSNR: hence the name, “sampling compensation factor”. One notes from (1.1) that
20 the number of particles necessary to perform a reconstruction also depends inversely on the SCF,
21 with smaller SCFs requiring larger numbers of particles.

22

1 In section 2, we derive all the formulae relating sampling to SSNR, including the case with missing
2 data, which requires special handling. In section 3, we derive analytical solutions to the sampling
3 and SCF for a variety of different cases. In section 4, we discuss the linear dependence of the SSNR
4 on N , as well as estimating the number of particles to perform a reconstruction. In section 5, we
5 show the correspondence between the proposed decrement of signal based on sampling and the
6 actual decrement in the SSNR when reconstructions are performed for two different proteins.

7

8

9

1 **Section 2. Decrement of SSNR due to sampling inhomogeneity**

2

3 In this section, we derive Eq (1.1), which provides an expression for the SSNR where all the
4 aspects of the sampling have been incorporated specifically into two parameters: the number of
5 particles, and a single geometrical factor. We assume that the effects of the microscope and the
6 effects of the noise can be approximately decoupled, in a manner that has otherwise been typically
7 assumed in the literature [18-20]. In section 2.1, we first consider the cases where the voxels in 3D
8 Fourier space are completely measured and derive the SSNR relationship, Eqs (1.1) and (1.2),
9 which is the main result of this paper. In section 2.2, we extend these derivations to cases when
10 there is missing data, by which we arrive at the adjusted formulae for resolution (2.30). We refer
11 to other sources, as necessary (Sorzano, [20] and Penczek [18]), for more detail on the aspects that
12 are not central to the derivations given here.

13 **2.1 Derivation of the Sampling Compensation Factor (SCF):**

14 The generally accepted understanding of 2D projection data after orientation assignment in cryo-
15 EM single-particle analysis is given by:

16

$$17 \quad F_j(\vec{k}) = M_j(k) X(R_j^T \vec{k}_3) + N_j(\vec{k}), \quad (2.1)$$

18

19 Here \vec{k} is a point in 2D Fourier space as measured on the projection j , where the projection j has
20 data $F_j(\vec{k})$ on the 2D grid point labeled by \vec{k} (see Figure 1). This is the usual Fourier space
21 description of a “single particle”. Eq (2.1) is our statement of the projection slice theorem: the

1 measured data should be a slice out of the true 3D map, X , but that has been modified in the
2 microscope by a transfer function, $M_j(\mathbf{k})$ and corrupted by $N_j(\vec{k})$, which is identically distributed
3 noise with mean zero and a variance that is independent of direction. This is the same set of
4 arguments that appear starting at Eq (7) from [20], as well as other places.

5
6 The 3D rotation, R_j^T , that appears in Eq (2.1) is the mapping from the 3D version, $\vec{k}_3 \equiv (k_x, k_y, 0)$
7 of the 2D point, $\vec{k} \equiv (k_x, k_y)$ to the 3D point, $R_j^T \vec{k}_3$, on the map, X , which is being reconstructed
8 (Figure 1). The “Euler angles” for the projection, j , are the angles that appear in the conventional
9 ZYZ representation of the rotation R_j . The factor $M_j(\mathbf{k})$, has been extensively described (Sorzano
10 [20], Penczek [18]) and should be an oscillating sinusoidal function (CTF) with a frequency-
11 dependent attenuation caused by various envelope effects. Eq. (2.1) is the generally accepted
12 starting point for cryo-EM data.

13
14 We next redefine \vec{k} to represent points on the 3D grid, and we shift our attention to the
15 reconstruction of the map in 3D. In the reasoning of direct Fourier reconstruction, we can form the
16 average over the samples that are used to reconstruct each 3D grid point \vec{k} to arrive at an estimate
17 of the 3D data point after reconstruction within the map $\hat{X}(\vec{k})$:

$$\hat{X}(\vec{k}) \equiv \frac{1}{\text{Sp}(\vec{k})} \sum_{j=1}^{\text{Sp}(\vec{k})} F_j(\vec{k}), \quad (2.2)$$

18
19
20
21 where $\text{Sp}(\vec{k})$ is the number of times that the particular point (in 3D) \vec{k} has been measured (by
22 means of projections as described above). In a conventional direct Fourier reconstruction, both the

1 running estimate of the reconstruction and the total weights that have been used for interpolation
2 (that is, $\text{Sp}(\vec{k})$) are kept as projection data is added.

3

4 One key observation is that, after substituting (2.1) into (2.2), the resulting noise is always down
5 by a factor of one over the square root of the amount of sampling (see [20]):

6

7
$$\hat{X}(\vec{k}) \equiv E(k)X(\vec{k}) + \frac{1}{\sqrt{\text{Sp}(\vec{k})}}\hat{N}(\vec{k}), \quad (2.3)$$

8 where the “renormalized” noise, $\hat{N}(\vec{k})$, has mean zero and the same variance as the average of the
9 variances of the constituent noise variables $N_j(\vec{k})$. Eq. (2.3) has been written so that the variance
10 of $\hat{N}(\vec{k})$ does not depend on the sampling. This is parallel to the argument which appears in [20]
11 at about Eq. (11). We have introduced $E(k)$, which is an effective envelope and the average over
12 the samples of the microscope influences ($M_j(k)$) as well as misalignment effects. Strictly
13 speaking, Eq (2.3) can only be approximate, but it is consistent with other approximate analyses
14 [21].

15

16 In the typical evaluation of cryo-EM resolution, two independent reconstructions are performed to
17 arrive at half maps which we can write in Fourier space as:

18

$$F(\vec{k}), G(\vec{k}) \text{ half maps} \quad (2.4)$$

1 We are interested in the FSC of half maps drawn from the same statistical ensemble as in (2.3).
 2 Therefore, we consider two maps assembled as in (2.3) and then we calculate the FSC. Each half
 3 map, F , G should therefore be of the form given as in Eq. (2.3):

$$4 \quad F(\vec{k}) = E(k)X(\vec{k}) + \frac{1}{\sqrt{\text{Sp}(\vec{k})}} \hat{N}_F(\vec{k}), \quad (2.5)$$

$$6 \quad G(\vec{k}) = E(k)X(\vec{k}) + \frac{1}{\sqrt{\text{Sp}(\vec{k})}} \hat{N}_G(\vec{k}).$$

7 The normal prescription is to introduce the correlation of these half maps at a discrete set of
 8 wavevector magnitudes, and then examine the functional behavior of this scalar correlation as a
 9 function of this wave-vector:

$$11 \quad \text{FSC}(k) \equiv \frac{\sum_{|\vec{k}'| \approx k} F(\vec{k}')G^*(\vec{k}')}{\text{Norm}_F(\vec{k}) \text{Norm}_G(\vec{k})}, \quad (2.6)$$

$$13 \quad \text{Norm}_F(k) \equiv \sqrt{\sum_{|\vec{k}'| \approx k} F(\vec{k}')F^*(\vec{k}')}, \quad (2.7)$$

$$14 \quad \text{Norm}_G(k) \equiv \sqrt{\sum_{|\vec{k}'| \approx k} G(\vec{k}')G^*(\vec{k}')}$$

16 Since F and G are assumed to be statistically similar, we can write (2.6) in short hand as

$$18 \quad \text{FSC}(k) \equiv \frac{\langle FG^* \rangle(k)}{\langle |F|^2 \rangle(k)} \quad (2.8)$$

19

1 Where we have used $\langle \cdot \rangle$ to mean the angular averages, (and equating angular and ensemble
2 averages). Very crudely, it is the cross correlation divided by the self-correlation. For more rigor,
3 see Sorzano et al [20], or Penczek. [18, 22].

4

5 Starting from (2.8), we can perform the familiar sort of calculation [18, 20, 23]

6

$$7 \quad \langle FG^* \rangle \approx E^2(k) \langle |X|^2 \rangle, \quad (2.9)$$

8

$$9 \quad \langle |F|^2 \rangle \approx E^2(k) |X|^2(\vec{k}) + \langle \frac{|\hat{N}(\vec{k})|^2}{\text{Sp}(\vec{k})} \rangle, \quad (2.10)$$

10

$$11 \quad \approx E^2(k) |X|^2(\vec{k}) + N_2(k) \langle \frac{1}{\text{Sp}(\vec{k})} \rangle, \quad (2.11)$$

12

13 where:

$$14 \quad N_2(k) \equiv \langle |N(\vec{k})|^2 \rangle, \quad (2.12)$$

15

16 and we have decoupled the noise variance from the sampling in going from Eqs (2.10) to (2.11).

17 There is no a priori reason to anticipate that the noise variances are related to the Euler angle

18 assignments, so the decoupling implicit in going from (2.10) to (2.11) is consistent with standard

19 assumptions. For the half maps, this leads to the following approximate estimate for the FSC using

20 the above (2.8), (2.9), (2.11):

$$21 \quad \text{FSC}(\mathbf{k}) = \frac{E^2(k) \langle |X|^2 \rangle}{E^2(k) \langle |X|^2 \rangle + N_2(k) \langle \frac{1}{\text{Sp}(\vec{k})} \rangle}, \quad (2.13)$$

22

1
$$= \frac{1}{1 + \frac{N_2(k)}{E^2(k) \langle |X|^2 \rangle (k)} \langle \frac{1}{\text{Sp}(\vec{k})} \rangle} , \quad (2.14)$$

2
3
$$= \frac{1}{1 + \left(2k \frac{\mathcal{N}(k)}{N} \right) \langle \frac{1}{2k \text{Sp}(\vec{k})/N} \rangle} , \quad (2.15)$$

4
5
$$= \frac{1}{1 + 2k \frac{\mathcal{N}(k)}{N \text{SCF}}} , \quad (2.16)$$

6
7 Where, in going from (2.14) to (2.15), we have defined

8
9
$$\mathcal{N}(k) \equiv N_2(k) / (E^2(k) \langle X^2(\vec{k}) \rangle) , \quad (2.17)$$

10
11 which is a noise-to-signal power ratio. In going from (2.15) to (2.16), we have defined:

12
13
$$\frac{1}{\text{SCF}} \equiv \langle \frac{1}{2k \text{Sp}(\vec{k})/N} \rangle . \quad (2.18)$$

14
15 The expression (2.18) is the same as (1.2), after identifying $\text{sp}(\vec{k}) \equiv \text{Sp}(\vec{k})/N$, the per-particle
16 version of the sampling. The expression for $\mathcal{N}(k)$ is the effective noise-to-signal ratio. Notably,
17 all the effects of sampling anisotropy are gathered into a single term: the SCF as given by (2.18).

18
19 Following previous formulations, we can define the spectral signal-to-noise ratio (SSNR):

20
21
$$\text{SSNR}(k) \equiv \frac{\text{FSC}(k)}{1 - \text{FSC}(k)} . \quad (2.19)$$

1

2 Substituting (2.16) into (2.19), we arrive at Eq (1.1):

3

4
$$\text{SSNR}(\vec{k}) = N \frac{\text{SCF}}{2k \mathcal{N}(k)} \quad , \quad (2.20)$$

5

6 where N is the number of particles, SCF is the geometrical factor, k is spatial frequency and $\mathcal{N}(k)$
7 is the effective noise to signal power (given by (2.18)), whose inverse would act like the
8 predominant component of the envelope. The reason for the regrouping of the factor $2k$ into the
9 SCF expression is that the SCF is then unity in a continuum calculation for the average density of
10 sampling for distributions that are uniform, as we will show in section 3. Eq (2.20) is essentially
11 the same expression that appears in [23] except for the appearance of the SCF term.

12

13 Under certain circumstances, the reconstructed volume may have regions of Fourier space that
14 have not been sampled. Two typical causes for this are: 1). The set of projection views are not well
15 distributed (such as top views), such that Fourier voxels, even very near the Fourier origin, have
16 not been filled. 2). The set of projection views are reasonably well distributed, but as one moves
17 further from the Fourier origin, there are lattice sites that are not sampled. Because Fourier voxels
18 not receiving information during the reconstruction procedure are traditionally left as zeros, there
19 will be voxels that do not contribute to the angular averages of (2.9) – (2.11). A careful
20 recalculation shows that the amended formula for the SSNR, as defined by (2.19), should still be
21 (2.20), except that the angular average involved with (2.18) for the evaluation of SCF should only
22 take place over non-zero voxels:

23
$$\frac{1}{\text{SCF}} \equiv \left\langle \frac{1}{2k \text{ sp}(\vec{k})} \right\rangle_{\text{non-zero voxels}} \quad . \quad (2.21)$$

1
2 Eq (2.20) along with (2.21) are our Eq. (1.1). As an aside, we show later that $\text{sp}(\vec{k})$ goes like $1/2k$,
3 so that the total sampling follows $N/2k$. Therefore, in the standard cryo-EM experiment, the total
4 sampling typically will not thin to zero, and the only zeros are the result of deficient projection
5 distributions.

6 **2.2 An adjusted formula for SSNR for half maps with unmeasured data**

7 The SSNR, based on half maps, has a drawback when some of the Fourier voxels have been left
8 unmeasured. The voxels in each half map are typically set to zero, which leads to smooth, but
9 artefactual maps, and may yield artificially high resolution measures. We see this in detail late in
10 Section 5, when we look at reconstructions that are performed from projections in a 45° cone, and
11 a percentage of randomly distributed extra projections is decreased in the sequence 10%, 3%, 1%,
12 0%. There is a sudden increase in the improperly defined FSC resolution measure at 0%. We will
13 defer discussion of the reconstructed data to that time. Here, we seek an adjusted SSNR expression,
14 which allows variance to be assigned to regions of Fourier space that have been unmeasured.
15 Consider the simplest situation, as in Figure 2, where we have represented some shell of Fourier
16 space by P measured values having mean T , and variance per voxel, var_N , given by the reciprocal
17 sampling at each measured point. There are also Q unmeasured voxels, that are assigned 0 values
18 in Figure 2A, and contribute neither to the signal nor the variance. Then, the ratio of signal to
19 variance is shown in the figure: $\text{SSNR} = N \frac{T^2}{\langle \frac{1}{\text{sp}} \rangle}$ where the average is taken of the reciprocal per-
20 particle sampling over the measured values. However, if one assigns a variance of 1 to the Q
21 unmeasured voxels, and repeats the same calculation, one arrives at $\text{SSNR} = N \frac{T^2}{\langle \frac{1}{\text{sp}} \rangle + N \frac{Q}{P}}$. In

1 particular, when N is large, the behavior is completely different: in Figure 2A the SSNR increases
 2 without bound, and in Figure 2B the SSNR plateaus to a finite value and is proportional to the area
 3 of measured to unmeasured region.

4

5 Generalizing the scenario in Figure 2, we consider a Fourier shell at Fourier radius k , and let P_k
 6 be the number of voxels that have non-zero sampling and let Q_k be the number of voxels with
 7 missing data. The total number of voxels therefore is then $P_k + Q_k$. We calculate the adjusted
 8 values of the quantities in (2.9) and (2.10), assuming that the data with missing voxels should be
 9 allowed to have variance. Then:

$$10 \quad \langle FG^* \rangle \approx \frac{P_k}{P_k + Q_k} E^2(k) \langle |X|^2 \rangle \quad (2.22)$$

11

$$12 \quad \langle |F|^2 \rangle = \frac{P_k}{P_k + Q_k} (E^2(k) |X|^2(\vec{k}) + N_2(k) \langle \frac{1}{\text{sp}(\vec{k})} \rangle) + \frac{Q_k}{P_k + Q_k} N_2(k) \quad (2.23)$$

13

14 Where E is defined through (2.3), N_2 is the noise variance defined as in (2.12), and $X(k)$ is the
 15 target structure. Our approach for the missing data is now clear: missing voxels take on a single
 16 unit of noise unattenuated by any sampling. The fairest assignment for such voxels is one unit of
 17 variance and zero units of signal. The adjusted formula, FSC^* , for the FSC then becomes:

18

$$19 \quad \text{FSC}^*(k) = \frac{P_k E^2(k) \langle |X|^2 \rangle}{P_k (E^2(k) \langle |X|^2 \rangle + N_2(k) \langle \frac{1}{\text{sp}(\vec{k})} \rangle) + Q_k N_2(k)} \quad (2.24)$$

20

1 This leads to an adjusted SSNR, which we develop by starting with its reciprocal:

2

$$3 \quad \frac{1}{\text{SSNR}^*(k)} \equiv \frac{1 - \text{FSC}^*(k)}{\text{FSC}^*(k)} = \frac{P_k N_2(k) \langle \frac{1}{\text{Sp}(\bar{k})} \rangle + Q_k N_2(k)}{P_k E^2(k) \langle |X|^2 \rangle}, \quad (2.25)$$

4

$$5 \quad = \frac{N_2(k) \langle \frac{1}{\text{Sp}(\bar{k})} \rangle}{E^2(k) \langle |X|^2 \rangle} + \frac{Q_k N_2(k)}{P_k E^2(k) \langle |X|^2 \rangle}, \quad (2.26)$$

6

$$7 \quad = \frac{1}{\text{SSNR}(k)} + \frac{Q_k N(k)}{P_k}. \quad (2.27)$$

8

9 Thus:

10

$$11 \quad \frac{\text{SSNR}(k)}{\text{SSNR}^*(k)} = 1 + \frac{Q_k N_2(k)}{P_k \langle \frac{1}{\text{Sp}(\bar{k})} \rangle}, \quad (2.28)$$

12

13

$$14 \quad = 1 + N \frac{Q_k}{P_k} \frac{\text{SCF}}{2k}, \quad (2.29)$$

15

16 which may be rewritten:

17

$$18 \quad \text{SSNR}^*(k) = \frac{\text{SSNR}(k)}{1 + N \frac{Q_k}{P_k} \frac{\text{SCF}}{2k}}, \quad (2.30)$$

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

$$SCF^*(k) \equiv \frac{SSNR^*(k)}{SSNR(k)} \cdot SCF = \frac{SCF}{1 + N \frac{Q_k SCF}{P_k 2k}},$$

where SCF^* is the expression to use in the adjusted version of Eq. (1.1): $SSNR^*(k) \cong N \frac{SCF^*}{2k} \frac{1}{N(k)}$. Eq. (2.30) gives an expression for reevaluating the SSNR for half maps when the original half maps have missing data. One way to think of Eq (2.30) is it shows how the conventionally constructed SSNR is inflated due to not assigning any variance to missing data. Eq. (2.30) also yields a condition by which a correction is necessitated. The ratio of occupied to unoccupied voxels at some Fourier wavevector is typically only a weak function of Fourier magnitude. This means it is also a geometrical parameter, similar to the SCF. Therefore, when

$$Q_k \gtrsim \frac{P_k}{N}, \tag{2.31}$$

then one should have to correct with the factor in the denominator of Eq. (2.30), to obtain a more realistic value of the SSNR. The condition (2.31) is the condition that the unmeasured variance is similar in magnitude with the measured variance, which is sampled in proportion to the number of particles. Another way to write it is that we must make an adjustment when $Q_k \gtrsim \frac{P_k}{N}$. If there is a sufficiently narrow gap, then we can ignore the adjustment. In practice, if there a sampling geometry that produces true missing gaps, then any number of particles should necessitate the alternate formula.

1 If N is sufficiently large, then what limits the resolution is solely the gap. Adding more particles
2 will not improve the SSNR, because additional particles will not better resolve the missing voxels,
3 and the already measured region is sufficiently well resolved. The expression for the adjusted
4 SSNR is most readily read off from (2.27), when the unadjusted value becomes large. Then the
5 first term on the right-hand side can be neglected, and the reciprocal of the remaining terms taken
6 to find the limit of large particle numbers, but with missing data:

7

$$8 \quad N \rightarrow \infty, \quad \text{SSNR}^*(k) \cong \frac{P(k)}{Q(k)} \frac{1}{N(k)} ; \quad \text{SCF}^*(k) \cong \frac{1}{N} \frac{P(k)}{Q(k)} 2 k . \quad (2.32)$$

9

10 Thus, in the limit of large particle numbers, the adjusted SSNR plateaus to a value, which is the
11 per particle envelope multiplied by the ratio of measured to unmeasured voxels. For positive k , the
12 expression implies that the FSC^* quickly drops from unity for even small $k > 0$: $\text{FSC}^*(k) \approx 1 -$
13 $\frac{Q_k}{P_k} \mathcal{N}(k)$ and is not improved by adding more particles. In this case, the measured voxels are
14 perfectly well sampled, and all the variance is due to the missing values.

15

16 To summarize, we derived the relationship between the SSNR and the type of sampling distribution
17 that is involved in the reconstruction. The latter enters the formula solely as a single geometrical
18 factor, the SCF, given by Eqs (2.20) and (2.21) (which reiterates Eq. (1.1)), the main result of this
19 work. In section 3, we derive analytical expressions for the SCF, and in section 5, we evaluate the
20 efficacy of (1.1) using simulated cryo-EM datasets. In the case of missing data, we suggest an
21 adjusted expression for the SSNR from what is usually used. This is the formula
22 for $\text{SSNR}^*(k)$ given by Eq. (2.30).

1 **Section 3. Numerical and analytical forms for the sampling function, and** 2 **expressions for the SCF geometrical factor**

3
4 In Section 1 and 2, we showed that the entire effect of the sampling inhomogeneity on the SSNR
5 could be incorporated into a single geometrical coefficient, the SCF. In this section, we provide
6 numerical and analytical forms for the sampling function, as well as the geometrical SCF factor
7 that causes decrement to SSNR curves. In section 3.1 we explain our numerical and analytical
8 approaches for evaluating the sampling and show that they evaluate identically for appropriate
9 cases. In section 3.2, we give continuum expressions for the sampling for several families of
10 distributions: 1) a one parameter family of distributions with an axial symmetry, that span the
11 complement to cones, which we term “side-like”; 2) a one parameter family of side-like
12 distributions modulated by fluctuations in the phi angle; 3) A two parameter family of projection
13 views that are constrained to fall within a cone of half-angle α , and that have, in addition, a
14 fraction, ϵ , of views that are randomly scattered through the remainder of Euler space. In section
15 3.3, we calculate analytically the SCF for each of these distributions using the continuum
16 formalism that we developed, which is valid when the sampling is not too small. The range of
17 values of the SCF for “side like views” ranges from 1 (the maximum, corresponding to uniform)
18 to $\frac{8}{\pi^2} = .81$ (side views). For the modulated side view cases, the SCF decreases as $\frac{8}{\pi^2} \sqrt{1 - \lambda^2}$,
19 where λ is the magnitude of the modulation, and we restrict the modulation >1 . This gives us a
20 complete parametrization of reasonable sampling where the SCF decreases from 1 to .81 (side) to
21 0. For the poorly sampled top-like views, we give a closed form integral expression for $SCF(\alpha, \epsilon)$
22 and evaluate the expression graphically. In the case when $\epsilon=0$, we point out that there are typically
23 missing values and the usual expression for the SSNR is not logical, as it neglects the variance that

1 can be estimated for the unmeasured voxels, by using the data already measured on the same shells
2 of Fourier space. Using the expressions that we developed in Section 2, we show theoretically that
3 properly defined SSNR curves should always improve after increasing the sampling (by increasing
4 the percentage of uniformly distributed views that lead to measured data in the unmeasured
5 region). All figures of the SCF curves and dependencies on control parameters are provided
6 accordingly.

7

8 **3.1 Discrete and continuum approaches to the sampling function**

9

10 **3.1.1 Discrete treatment for sampling**

11 The projection-slice theorem [24] states that a 2D projection from a direction \hat{n} of a 3D map, is a
12 slice out of the Fourier volume of the plane perpendicular to \hat{n} and passing through the origin, as
13 shown in Figure 1. If we think of the map as rotated by R before the projection (along \hat{z}), then
14 what we term the projection direction, \hat{n} , is (approximately) perpendicular to the sampled points,
15 and is given by $\hat{n} \equiv R^T \hat{z}$. As suggested in Figure 1, each projection, \hat{n}_j , samples the set of points \vec{k}
16 satisfying

17

$$18 \quad |\hat{n}_j \cdot \vec{k}| \leq \frac{1}{2}. \quad (3.1)$$

19

20 The totality of the discretely sampled points form a lattice as shown in Figure 3. Here, a single
21 projection (in Fourier space) is taken in the \hat{z} direction with Fourier magnitudes less than the real
22 space box, L . Lattice sites, shown as blue dots in the $k_z = 0$ plane are considered to be sampled.
23 Each sampled plane selects a lattice of points in this manner. Our numerical algorithm hinges on

1 finding lattice sites that satisfy (3.1) for each projection. As we sum over projections, we increase
2 the totality of “viewings” of each lattice site. In direct Fourier inversion, this integer number of
3 “viewings” will correspond roughly to the reconstruction weights.

4
5 The number of times a particular 3D point, \vec{k} is sampled, we term $\text{Sp}(\vec{k})$, and is therefore given by
6 the cardinality of the set of the projections, that for a given \vec{k} , satisfy the criterion of Eq. (3.1).

7 Therefore:

8
9
$$\text{Sp}(\vec{k}) \equiv \sum_{j=1}^N \Theta(|\hat{n}_j \cdot \vec{k}| \leq \frac{1}{2}) , \quad (3.2)$$

10
11
12 where Θ is the indicator function (see Glossary). The per-particle sampling function we define as:

13
14
$$\text{sp}(\vec{k}) \equiv \frac{1}{N} \text{Sp}(\vec{k}) . \quad (3.3)$$

15
16 Eq. (3.1) is what is used numerically to find the sampling at each voxel, wherein the vector to each
17 voxel is checked against every projection to see if the dot product between this vector and the unit
18 direction given by the projection is sufficiently small (less than $\frac{1}{2}$ in magnitude).

19
20 We investigate suitably many approximations that the sampling function emerges as a quantity
21 that independently affects the SSNR (and only coupled to average microscope effects: not
22 individual CTFs per particle, for example). It is our hypothesis that this level of approximation is
23 sufficiently useful to enable understanding the effect of anisotropy on resolution.

1

2 **3.1.2 A continuum treatment for sampling.**

3 We wish to formulate the expressions analytically whenever possible. Toward this end, we recast
4 (3.2) using Dirac delta functions, which will provide continuum calculations that are both useful
5 and accurate. For a single projection in the z-direction, we would like to employ

$$6 \quad \text{sp}(\vec{k}) = \delta(\vec{k} \cdot \hat{z}) \quad . \quad (3.4)$$

7 Generally the Dirac delta function is considered to be $M \Theta\left(|k_z| < \frac{1}{2M}\right)$, in the limit that the
8 parameter M becomes arbitrarily large, whereas we have taken M as simply unity in (3.1). The
9 delta function analytical approximation is crude, but satisfies the proper normalization.

10 The generalization of (3.1) for continuum calculations using the idea in (3.2) yields

$$11 \quad \text{Sp}(\vec{k}) \equiv \sum_{j=1}^N \delta(\hat{n}_j \cdot \vec{k}) \quad , \quad (\text{for analytical evaluations}) \quad (3.5)$$

$$12 \quad \equiv \int \rho(\hat{n}) \delta(\hat{n} \cdot \vec{k}) \quad , \quad (3.6)$$

13 where $\rho(\hat{n})$ is a measure on the distributions of projections parametrized by \hat{n} : (in this case,
14 $\rho(\hat{n}) = \sum_{j=1}^N \delta(\hat{n}_j - \hat{n})$, which is discrete, but generally $\rho(\hat{n})$ may be continuous). Eq. (3.5) is
15 the continuum approximation when the length of the side of the box, which we will use as L, can
16 be considered to be much larger than 1. This is sufficient for many of our analytical treatments and
17 development of formulae, since we are often working far from the Fourier origin. In Eqs. (3.5), we
18 consider Fourier space to be dimensionless (unitless), which is a common practice. To reintroduce
19 units, if one has, in 1D, 200 voxels of voxel size 1 Å per side, then each Fourier space voxel will

1 have width $\frac{1}{200 \text{ \AA}}$ and the largest distance from the Fourier origin will be $\frac{1}{2 \text{ \AA}}$ (this is the Nyquist
 2 frequency). The average sampling across a shell at fixed Fourier magnitude can be derived using
 3 our continuum treatment. Starting from Eq. 3.4:

$$4 \quad \langle \text{sp}(\vec{k}) \rangle = \frac{1}{4\pi} \int_0^{2\pi} d\phi \int_0^\pi d\theta \sin \theta \delta(k \cos \theta) , \quad (3.7)$$

$$5 \quad = \frac{1}{2} \int_0^\pi d\theta \sin \theta \delta(k \cos \theta) = \frac{1}{2k} \quad (3.8)$$

6 This is a natural result: placing planes (Fourier slices) into volumes, the density must fall off as
 7 one over the Fourier radius. A more thorough derivation is given in Appendix A.1, including an
 8 interpretation of the geometrical factor 2, which couples to the k dependence. Eq. (3.8) is also the
 9 sampling per particle for a uniform set of projections, but Eqs (3.7) and (3.8) hold for any
 10 distribution of projections.

11 **3.1.3 Consistency between numerical and analytical expressions for sampling**

12 As a check of both our code and analytical implementation, we tested the total amount of sampling
 13 in our volume by placing 50,000 projections in a box of size $L \times L \times L$ with $L=41$. We calculated
 14 the integer sum, S , over all the sampling at all the points, and evaluated $S/4L^2$ numerically to be
 15 1.19. To develop an analytical expression of this idea, we can write

$$16 \quad S = \int_{-L}^L dk_x \int_{-L}^L dk_y \int_{-L}^L dk_z \left(\frac{1}{2k} \right) , \quad (3.9)$$

17 where k is spatial frequency. This is the average amount of intersection of arbitrarily oriented
 18 planes with a cube of side $2L$. In the appendix A, we show that the integral evaluates to

1
$$\frac{S}{4L^2} = 3 \left(-\frac{\pi}{12} + \ln(1 + \sqrt{3}) - \ln\sqrt{2} \right) = 1.19. \quad (3.10)$$

2 This corroborates our numerical result described above.

3

4 **3.2: Sampling Function for three different distributions in continuum representation**

5

6 We calculate the sampling function for three different distributions. The first is the case of the
7 complement to a cone (which we term side-like). The second is for side views with a modulation
8 in the azimuthal Euler angle. Finally, we also calculate the top-like cases, where a certain fraction
9 of uniform views is also included. The side-like views and side-modulated views are each
10 governed by single parameters: i) the cone half angle, α and ii) the modulation parameter, λ . The
11 top-like family of distributions is governed by two parameters: once again, the cone half-angle, α ,
12 and a parameter, ϵ , to cover uniform projections in the complement to this region.

13

14 Figure 4 shows the projection distributions that will be described in this section, including a
15 schematic representation of the projection distribution (top row), the population of Fourier space
16 through slice insertion (middle row), and the experimental sampling map derived from 10,000
17 insertions (bottom row). These are displayed for different sampling schemes, including the
18 uniform (Figure 4A), side-like or complement to cone (Figure 4B), side (Figure 4C), side-
19 modulated (Figure 4D), and top-like (Figure 4E).

20

21 **3.2.1 side-like cases (α)**

22 For the side-like case, we have

1

2

$$\text{sp}(\vec{k}) = \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} \delta(\hat{n} \cdot \vec{k}) / C_N, \quad (3.11)$$

3

4 where C_N is a constant to ensure the normalization (3.8), leading to (see Appendix B):

5

$$\text{sp}(k, \theta) = \frac{1}{k} \frac{\sin^{-1}\left(\frac{\cos \alpha}{\sin \theta}\right)}{\pi \cos \alpha}, \quad \left| \frac{\pi}{2} - \theta \right| < \alpha; \quad (3.12)$$

8

$$\frac{1}{2k \cos \alpha} \cdot \left| \frac{\pi}{2} - \theta \right| \geq \alpha \quad (\text{"side-like"}) \quad (3.13)$$

9

10 The distribution for side views can be selected by taking the $\alpha \rightarrow \frac{\pi}{2}$ limit to arrive at:

$$\text{sp}(k, \theta) = \frac{1}{\pi k \sin \theta}, \quad \theta > 0, \quad (\text{"side-view"}) \quad (3.14)$$

12

13 Along the z-axis (that is, $k \sin \theta = 0$), sp should have the same value as at the origin, which is 1.

14

15 **3.2.2 modulated side-views (λ):** A set of modulated side view projections can be written as a
 16 density distributions:

17

$$\rho(\phi_n) = (1 + \lambda \cos 2\phi_n) \quad , \quad (3.15)$$

19 where ϕ_n is the azimuthal angle for a projection direction \hat{n} . This gives rise, therefore to a sampling
 20 given by:

21

$$\text{sp}(k, \theta, \phi) = \frac{\int d\hat{n} \delta(\hat{n} \cdot \hat{z}) (1 + \lambda \cos 2\phi_n) \delta(\hat{n} \cdot \vec{k})}{C_N},$$

22

2

$$1 \quad = \frac{(1-\lambda \cos 2\phi)}{\pi k \sin \theta} \quad (\text{"modulated side views"}) \quad (3.16)$$

3

4 We describe in appendix A.3 how to select a set of projections with this form, using the cumulative
5 distribution function.

6

7 **3.2.3 top-like cases (α, ϵ)**

8 Finally, we consider sampling for the top-like cases:

9

$$10 \quad \text{sp}(\vec{k}) = \int^{|\hat{n} \cdot \vec{z}| > \cos \alpha} d\hat{n} \delta(\hat{n} \cdot \vec{k}) / C_N \quad , \quad (3.17)$$

11

12 which leads to (see Appendix B):

$$13 \quad \text{sp}(k, \theta) = \frac{1}{2k} \frac{\cos^{-1}\left(\frac{\cos \alpha}{\sin \theta}\right)}{\pi \sin^2(\alpha/2)} \quad , \quad \left| \frac{\pi}{2} - \theta \right| < \alpha, \quad (3.18)$$

15

$$14 \quad (\text{"top-like"}) \quad 0, \quad \left| \frac{\pi}{2} - \theta \right| \geq \alpha. \quad (3.19)$$

16

17 Taking $\alpha \ll 1$ leads to arbitrarily large values of sp: $2k \text{sp}(k, \pi/2 - \eta) = 4 \sqrt{\alpha^2 - \eta^2} /$
18 $\pi \alpha^2$ (for $|\eta| < \alpha$). Once again, the sampling needs to be truncated to unity, when $\alpha = 0, \theta = \pi/2,$
19 that is in the xy plane, if the top-view is taken along the z-direction.

20

21 These distributions have missing data, and so we can calculate, for each shell, the ratio of filled,
22 P, to unfilled voxels, Q.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

$$\frac{P}{Q} = \frac{\int_{\pi/2-\alpha}^{\pi/2} \sin \theta \, d\theta}{\int_0^{\pi/2-\alpha} \sin \theta \, d\theta} = \frac{\sin \alpha}{1 - \sin \alpha} \quad (3.20)$$

We will need Eq. (3.20), to compute $SSNR^*$, as argued in section 2, because there is missing data and develop an adjusted formula for top-like SCF distributions in the next subsection.

We also can evaluate the top-like cases, when we add a random distribution of projections, so as to fill in the missing data. Then:

$$2k \, sp(k, \theta) = (1 - \epsilon) \frac{\cos^{-1}\left(\frac{\cos \alpha}{\sin \theta}\right)}{\pi \sin^2(\alpha/2)} + \epsilon \quad , \quad \left| \frac{\pi}{2} - \theta \right| < \alpha; \quad (3.21)$$

$$\epsilon \quad , \quad \left| \frac{\pi}{2} - \theta \right| \geq \alpha \quad (\text{"top - like with uniform"})$$

where ϵ is the fraction of projections that are distributed randomly, and the rest fall in the original cone of half-angle α .

Section 3.3 The Sampling Compensation Factors for the three different distributions

The SCF is defined via:

$$1 \quad \frac{1}{\text{SCF}} \equiv \left\langle \frac{1}{2k_{\text{sp}}(\mathbf{k})} \right\rangle \quad (3.22)$$

2
3 where $\langle \cdot \rangle$ is the average over solid angle regions that have non-zero values of the sampling. We
4 can evaluate this numerically for the “top-like” and “top-like with uniform” distributions, and
5 analytically for the “side-like” and “modulated side-view” cases.

6
7 In Appendix A.5, we evaluate Eq. (3.22) using (3.14) to arrive at:

$$9 \quad 1/\text{SCF}_{\text{side modulated}}(\alpha) = \frac{\pi^2}{8} \frac{1}{\sqrt{1-\lambda^2}} \quad , \quad (3.23)$$

10
11 which ranges from arbitrarily large values to $\frac{\pi^2}{8}$ (when $\lambda = 0$; no modulation). In practice, the
12 sampling never achieves a continuum of values, so the expression given by (3.23) cannot be used
13 for λ very close to 1.

14
15 One can also show:

$$17 \quad 1/\text{SCF}_{\text{side-like}}(\alpha) = \cos \alpha (1 - \sin \alpha) + \frac{\pi}{2} (\cos \alpha) \int_0^\alpha \frac{\cos \theta d\theta}{\arcsin\left(\frac{\cos \alpha}{\cos \theta}\right)} \quad , \quad (3.24)$$

18
19 ranges from $\frac{\pi^2}{8}$ ($\alpha = \frac{\pi}{2}$ side views) to a value of 1 ($\alpha = 0$, uniformly distributed views).

20
21 Figure 5 shows schematically the behavior of the SCF for the side-modulated and side-like cases.
22 It is shown there how to continuously vary the SCF from its low values (corresponding to side-

1 modulated cases with sampling suffering from deep pockets) to its highest value unity (for uniform
2 views).

3

4 Finally, we want to calculate the quantity that represents the top-like situation. Since much of the
5 region is zero, the normalization from (3.22) must be carefully calculated and leads to:

6

$$7 \quad 1/\text{SCF}_{\text{top-like}}(\alpha) = \frac{\pi}{2} \tan \alpha/2 \int_0^\alpha \frac{\cos \theta d\theta}{\arccos\left(\frac{\cos \alpha}{\cos \theta}\right)} \quad (3.25)$$

8

9 The right-hand side of (3.25) ranges from 1 ($\alpha = \frac{\pi}{2}$, uniformly distributed views) to 0 (for $\alpha = 0$,
10 purely top views). The asymptotics are $1 - 2\left(\frac{\pi}{2} - \alpha\right)$ for $\frac{\pi}{2} - \alpha \ll 1$, and $\frac{\pi^2}{8}\alpha$ for small α . Thus,
11 we have a set of analytical expressions for SCF that can run from arbitrarily small to unity and
12 from unity to arbitrarily large levels. However, any distribution with $\text{SCF} > 1$, involves
13 distributions with missing data. Ultimately, the more relevant attribute, will be SCF^* which relates
14 how the correctly adjusted SSNR^* is decremented due to the sampling. Thus, the SCF^* is bounded
15 by $0 \leq \text{SCF}^* \leq 1$.

16

17 Repeating with the additional random projections gives a drastically different value for the SCF
18 for the singular change of adding partially uniform perturbations, because now all (or most) of the
19 Fourier points have at least some sampling.

20

$$21 \quad \frac{1}{\text{SCF}_{\text{top-like}, \epsilon}(\alpha)} = \frac{\pi}{2} \sin^2 \frac{\alpha}{2} \int_0^\alpha \frac{\cos \theta d\theta}{\pi \epsilon \sin^2 \frac{\alpha}{2} + (1-\epsilon) \arccos\left(\frac{\cos \alpha}{\cos \theta}\right)} + \frac{1-\sin \alpha}{\epsilon} \quad (3.26)$$

22

1 For such top-like distributions, it is interesting to compare $\text{SCF}_{\text{top-like}}^*(\alpha)$ from (2.31) with
 2 $\text{SCF}_{\text{top-like},\epsilon}(\alpha)$ for small but finite ϵ from Eq. (3.26). From Eq (2.31), (3.20) and (3.25), we can
 3 derive the $\epsilon=0$ quantity:

$$4 \quad \frac{1}{\text{SCF}_{\text{top-like}}^*(\alpha)} = \frac{1}{\text{SCF}_{\text{top-like}}(\alpha)} + N \frac{1-\sin \alpha}{\sin \alpha} \frac{1}{2k}, \quad (3.27)$$

$$5 \quad = \frac{\pi}{2} \tan \alpha/2 \int_0^\alpha \frac{\cos \theta d\theta}{\arccos\left(\frac{\cos \alpha}{\cos \theta}\right)} + N \frac{1-\sin \alpha}{\sin \alpha} \frac{1}{2k}. \quad (3.28)$$

6
 7
 8
 9 Note that $\text{SCF}_{\text{top-like}}^* = \text{SCF}_{\text{top-like},\epsilon}(\alpha)$, when $\alpha = \frac{\pi}{2}$. For small ϵ , but large N , the second
 10 terms of both (3.26) and (3.28) dominate and we get:

$$11 \quad \left(\alpha < \frac{\pi}{2}, N \gg 1, \epsilon \ll 1 \right) \text{SCF}_{\text{top-like}}^*(\alpha) \approx \frac{2k}{N} \frac{\sin \alpha}{1-\sin \alpha}, \quad (3.29)$$

$$12 \quad \text{SCF}_{\text{top-like},\epsilon}(\alpha) \approx \frac{\epsilon}{1-\sin \alpha} \quad (3.30)$$

13
 14
 15
 16 The crossover between these expressions occurs approximately when

$$17 \quad \epsilon(\alpha) \cong \sin(\alpha) \frac{2k}{N}. \quad (3.31)$$

18
 19
 20 The situation for the decrement in the correctly adjusted SSNR is depicted In Figure 6, for the
 21 poorly sampled cases. The Eq. (3.28) is the lower bounding curve in gray ($\epsilon = 0$). Otherwise, the

1 curves represent Eq. (3.26). There is no crossover, unless epsilon is sufficiently small: in the figure
2 there are only three crossings of the curves. For $k = 15$, and $N = 10^4$, Eq. (3.31) implies $\epsilon =$
3 $0.03 * \sin \alpha$, $SCF^* = 0.03 \frac{\sin \alpha}{1 - \sin \alpha} =$ is the crossover between curves.

4
5 The last expressions tell the entire story of missing data. If data is missing in some sizeable region,
6 the adjusted SSNR is drastically reduced. However, even a small fraction of random perturbations
7 starts to quickly reintroduce signal. If there is a gap, the SCF is increased by a factor

8
9
$$\frac{SCF^*(\epsilon)}{SCF^*(\epsilon=0)} = \frac{N\epsilon}{2k \sin \alpha}, \quad (\text{“ratio of top-like SCF at finite } \epsilon, \text{ to } \epsilon = 0, \text{ top-like”}) \quad (3.32)$$

10
11 by adding back a fraction ϵ worth of random perturbations. For $N = 10^4, k = 15, \alpha = 45^\circ, \epsilon =$
12 0.01 , the RHS becomes 4.7 , which is a huge jump over such a small change in ϵ . Conversely,
13 having an empty region of Fourier space gives much lower SCF^* than a lightly sampled Fourier
14 space.

1 **Section 4. Relationship between SSNR and the number of particles N in a** 2 **reconstruction**

3 In Section 2, we derived the relationship (2.20) (or equivalently (1.1)), which is the estimate of the
4 SSNR in terms of the sampling. There are two aspects of the latter: the cumulative extrinsic effect
5 due to the number of particles in the data, and the shape of the distribution of the sampling (or
6 projection directions), an intrinsic quality. When Fourier space is reasonably sampled everywhere,
7 we can assign a single parameter to each of the extrinsic and intrinsic qualities of the sampling: N ,
8 the number of particles, and SCF, the sampling compensation factor, defined as in Eq (1.2). The
9 SSNR is seen to be proportional to each quantity, with the SCF attaining its maximum value of
10 unity when the distributions of projections are uniform.

11

12 In this section, we revisit the dependence of the SSNR on N , the number of particles, when every
13 other aspect of the problem is held constant:

14

$$15 \quad \text{SSNR}(N, k) = N \text{ssnr}(k) \quad , \quad (4.1)$$

16

17 for some function, ssnr , which is the form of Eq. (1.1) with $\text{ssnr}(k) \equiv \frac{\text{SCF}}{2k} \frac{1}{\mathcal{N}(k)}$. Eq. (4.1) is the
18 familiar way that the signal in a noisy system should accrue, if N represents the total number of
19 measurements. The per-particle SSNR, which depends on many factors that corrupt the final
20 reconstruction, is observed to be quite rigid and independent of N , as previously noted [19]. The
21 resulting universal curve, $\text{ssnr}(k)$, includes multiple components inherent within the cryo-EM
22 pipeline that attenuate resolution: attenuation due to the microscope transfer function, detector

1 noise, incorrect image orientation assignment, structural heterogeneity, among others. The
2 consequence is that the number of particles needed to obtain a higher resolution using the same
3 collection scheme can be determined from a single SSNR curve, provided that the curve is
4 sufficiently smooth at the desired resolution: indeed, smoothness of the SSNR curve might be
5 another possible criterion for resolution. The universality of the SSNR/ N curves is akin to the
6 familiar Reslog [25] or Guinier [19] analyses.

7

8 **4.1 Linear dependence of SSNR on N**

9

10 According to Eq. (4.1), dividing the SSNR by the number of particles results in a universal per-
11 particle curve. To test this idea, we looked at sequences of FSC, equivalently SSNR, curves for
12 reconstructions using successively larger number of particles, N , for data from an experimental
13 dataset contributing to a 2.9 Å reconstruction of the eukaryotic large ribosomal subunit [26]. Figure
14 7A shows a total of seventeen experimental FSC curves, from $N=7000$ to 70000 particles. The
15 series of FSC curves collapse to a universal curve via SSNR/ N , as predicted by 4.1, where
16 $SSNR(N, k) \equiv \frac{FSC(N, k)}{(1-FSC(N, k))^2}$, as shown in Figure 7B. Although this idea has appeared formally in
17 many places [19, 25, 27, 28], we have not noted the explicit construction of such universal curves,
18 as highlighted here. For smaller values of particle number, N , the $ssnr(k)$ curve loses continuity
19 at smaller values of resolution and limits our ability to calculate the necessary number of particles
20 to achieve higher resolutions, as described below.

21

22 **4.2 Number of Particles Necessary for Reconstruction**

23

1 Eq. 4.1 can be used to predict the number of particles necessary to attain a given resolution for a
2 general envelope function, $ssnr(k)$, derived from a single SSNR curve. A common scenario that
3 is encountered during cryo-EM data collection is one in which the experimentalist asks whether
4 the current approach is conducive toward achieving a target resolution, given a fixed amount of
5 collection time. Our claim is that, there is some N_0 so that for $N = N_0$, we can construct the curve
6 $ssnr(k) \equiv SSNR(N_0, k)/N_0$ and arrive at a reasonable estimate predicting the necessary number
7 of particles (it is conceivable to make a lower estimate for the necessary N_0 , but this is beyond the
8 scope of the current discussion). Thus, for a resolution criterion, $FSC = FSC^*$, one arrives at an
9 implied criterion, $SSNR = SSNR^* = FSC^*/(1 - FSC^*)$ (If $FSC^* = 0.143$, or 0.5 then $SSNR^* =$
10 0.167 or 1.0 respectively). Next, one defines k_T to be the target resolution. Then the necessary
11 number of particles, N_T , to achieve the target resolution is given by:

$$N_T = SSNR^*/ssnr(k_T) \quad . \quad (4.2)$$

12
13
14
15 Graphically, we can make a construction on a semilog plot of the original SSNR curve, and realize
16 that

$$\log N_T/N_0 = \log(SSNR^*/SSNR(k_T, N_0)) \quad , \quad (4.3)$$

17
18
19
20 which follows from Eq. (4.1), which implies $N_T/N_0 = SSNR(N_T, k_T)/SSNR(N_0, k_T)$, and
21 $SSNR(k_T, N_T) = SSNR^*$. Now, Eq. (4.3) can be used to graphically find the number of particles
22 needed to achieve a target resolution, since the shift from the current resolution to a target
23 resolution gives the ratio in the number of particles to increase to. Unlike other methods discussed

1 in this section, this makes no assumptions about the functional form of the per-particle SSNR,
2 $\text{ssnr}(k)$: it can be exponential (Reslog) or Gaussian (Guinier) or indeed hold to any shape.

3
4 The idea is demonstrated for the ribosome sequence of reconstructions in Figure 7C. For
5 convenience and in line with standard assumptions in the cryo-EM literature, we used the same
6 two FSC criteria described above of 0.5 and 0.143, which is equivalent to an SSNR condition of
7 $\text{SSNR}^* = 1$ and 0.167, respectively, and analyzed the SSNR curve corresponding to 7000 particles.
8 Using $\text{SSNR}^* = 1$ (or equivalently $\text{FSC}=0.5$), the resolution is measured to be 7.9 Å. To obtain the
9 necessary ratio of number of particles required for reaching the target resolution of 4.2 Å, and
10 using this same criterion, we can measure the difference on the log plot, which is 2.3 or $\log(10)$,
11 that is one decade. Therefore, the prediction is that 10 times the original number of particles are
12 necessary to obtain a reconstruction at 4.2 Å. When the orange dotted curve that corresponds to
13 10x particles is then inspected on the plot, the prediction is corroborated, since the resolution of
14 the 70K particle FSC curve, where the orange dotted curve intercepts the SSNR^* condition,
15 matches to the predicted 4.2 Å. The identical analysis is repeated using the $\text{FSC}=0.143$ criterion
16 in Figure 7D.

17
18 Finally, we note that the SSNR is inversely proportional to the geometrical SCF factor, so that
19 distributions with lower SCF (more fluctuations in the sampling) require larger numbers of
20 particles. Under typical data collection procedures, the SCF is fixed by the sample preparation and
21 microscope conditions, and one cannot easily consider the use of the SCF as an independent control
22 parameter that can be conveniently varied. The exception would be to tilt the specimen, which
23 would alter the orientation distribution, and thus the SCF [15].

1

2 **4.3 Comparison of graphical methods (Guinier, Reslog, and per-particle SSNR curve)**

3

4 The Guinier [19] and Reslog [25] formulations are popular for extrapolating the number of
5 particles necessary for reconstruction. We would like to understand the relationship between these
6 graphical constructions and the per-particle SSNR curves. We give a thorough analysis of the
7 Guinier analysis, and see that the Guinier assumptions essentially also imply (4.1), but restrict
8 $\text{ssnr}(k)$ to a Gaussian form. Our method is seen to be slightly more general, but in typical usage,
9 identical to these, based on the argument below.

10

11 The prescription in Guinier analysis [19] is to estimate the number of particles needed to achieve
12 a given resolution, and mark this on a semilog plot of N as a function of the square of the spatial
13 frequency, and repeat. This procedure is presumed to form a line, which can be extrapolated to
14 find the number of particles to achieve a desired resolution. That is, knowing N_1 , define k_1
15 implicitly by $\text{SSNR}(N_1, k_1) = \text{SSNR}^*$, where SSNR^* is the fixed value of SSNR that demarcates
16 resolution as described above, and define k_2 similarly. The Guinier assumption is that:

17

$$18 \quad (k_2^2 - k_1^2)\lambda^2/2 = \log N_2/N_1 \quad , \quad (4.4)$$

19

20 for some constant λ , for resolutions of interest corresponding to k_1, k_2 . That is, along the fixed
21 contours of SSNR, the change in the square of the resolution is proportional to the logarithm of
22 the ratio of the number of particles used to achieve the SSNR criterion. By means of such a

1 construction, one can estimate the number of particles needed to achieve a higher resolution. Eq
2 (4.4) is easily solved formally as

3

$$4 \quad C_1 = \log N(k) - k^2 \lambda^2 / 2 \quad , \quad (4.5)$$

5

6 where C_1 acts like a constant of integration, which depends on the SSNR, which is held fixed in
7 the construction. This implies (exponentiate 4.5):

8

$$9 \quad \text{SSNR}(N, k) = H(N \exp(-\frac{1}{2} \lambda^2 k^2)) \quad (\text{Guinier}), \quad (4.6)$$

10

11 for some function H . Every set of SSNR curves of the form (4.6) will yield (4.4). The only
12 reasonable choice for H is linear, which matches our result (4.1), when $\text{ssnr}(k)$ from (4.1) is a
13 Gaussian. We should point out that in light scattering, Guinier plots are used as a low frequency
14 approximation where the various physical parameters can rigorously be argued to hold to the
15 damped Gaussian format indicated by (4.6) [29].

16

17

18 The Reslog analysis [25] is very similar and leads to:

19

$$20 \quad \text{SSNR}(N, k) = H(N \exp(-k/c)) \quad (\text{Reslog}), \quad (4.7)$$

21

22 for some constant wavevector c . Once again, the only reasonable choice is a linear function, H ,
23 leading to Eq. (4.1) with an exponential form for $\text{ssnr}(k)$ in Eq. 4.1.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

Heymann [23] made an identical argument to arrive at our Eq. (4.1) and used the Guinier analysis, based partly on the formal results on blurring [30] and other envelopes [21]. As an aside, much like multiple time scales [31] can create an effective 1/frequency noise in physical systems rich with multiple time scales, with so many differing sources of noise in cryo-EM, it may be that, depending on the experimental circumstances, the linear behavior is equally valid to the quadratic behavior for governing the log of the envelope. In any case, Heymann suggests the Gaussian form for $ssnr(k)$ above: $ssnr(k) = \exp(-\frac{1}{2}\lambda^2 k^2)$, consistent with the Guinier analysis. Although Heymann arrives at Eq. (4.1), he does not arrive at our Eq. (1.1), because the possible anisotropy is not discussed, and therefore he uses the expression for the uniform distribution of the per-particle sampling ($1/2k$) which is our Eq. (3.8) and Eq. (9) of Heymann [23] .

1 **Section 5. Decrement of SSNR through non-uniform sampling**

2

3 In Section 4, we observed that the SSNR depends in two ways with the sampling; the extrinsic part
4 governed by the number of particles N (as already has been discussed in the literature) as well as
5 the type of the sampling governed by the geometrical factor of the sampling map, which we have
6 termed the SCF. In this section, we test whether the SCF (or SCF^{*}) has the predicted effect on the
7 SSNR as described by Eq. (1.1) and explained in section 3.3. We look at sequences of
8 reconstructions of two proteins that vary in their size and shape: the influenza hemagglutinin trimer
9 and human apoferritin, for all the situations for which we calculated the SCF (or SCF^{*}) values in
10 Section 3.3. In each case, we compare the SSNR curves of reconstructions versus the baseline
11 case, which is a set of uniformly distributed views.

12

13 **5.1 Methods**

14

15 **5.1.1 Generation of projection distributions**

16

17 We generated a set of 10,000 projection Euler angles for sequences of different sampling
18 distributions, each of which is described in section 3.2. We evaluated three different schemes for
19 modulating the projection distribution and comparing to the uniform distribution, as depicted in
20 Figure 8. For the well-sampled side-like sequence, we used pure side views and modulated side
21 views with a set of modulation parameters given by $\lambda = 0.4, 0.6, 0.8$, and 1.0 , (Figure 8A). For
22 the first of the more poorly-sampled cases, we selected top-like projections, distributed within
23 varying cone sizes of half angular width (5° , 30° and 45°), and fixed a small amount of random

1 projections (3%) distributed evenly across the rest of Euler space (Figure 8B). This scenario
2 evaluates the effect of increasing cone size. For the second of the more poorly-sampled cases, we
3 fixed the cone size to be 45° and added random assignments of 0%, 1%, 3% and 10% evenly
4 distributed projections across the rest of Euler space (Figure 8C). This scenario evaluates the effect
5 of increasing the amount of random projection “sprinkling” in the presence of an otherwise fixed
6 distribution.

7

8 **5.1.2 Synthetic data generation with distinct projection distributions**

9

10 To test our idea relating the effect of a single geometrical parameter and the SSNR, we generated
11 synthetic datasets corresponding to two proteins of varying size and shape, namely the
12 hemagglutinin (HA) trimer and apoferritin. The synthetic data generation followed previously
13 described protocols [15, 32-34]. Briefly, 10K projections were generated from cryo-EM maps of
14 either HA or apoferritin, according to the viewing directions that were described in Section 5.1
15 above. These projections were shifted and rotated at random and noise was added. Next a
16 distribution of CTFs were applied to the 2D projections, followed by an additional layer of noise
17 to arrive at an SNR approximately equal to 0.05. This SNR is consistent with experimental cryo-
18 EM data [35]. A reconstruction was performed with the known orientations using the Frealign
19 software, and the usual FSC was calculated between half maps. In parallel, the angle assignments
20 were used to calculate sampling maps, as described in Section 3 (and shown graphically in Figure
21 4). From the sampling maps, the SCFs were calculated numerically by implementing (3.22).

22

23 **5.2 Results comparing decrements predicted by sampling and reconstructions.**

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

We tested how well the SCF geometrical parameter, based solely on the projection directions, could predict the decrement of the SSNR, with all other aspects of the problem held constant.

5.2.1 Side-like and side-modulated sampling cases

We first proceeded to test the predictive ability of the SCF on well-sampled cases, where most of the values of the sampling remain reasonably high and each index point is sufficiently sampled, e.g. above 20. From a theoretical perspective, we should expect that the ideas set forth are most accurate in the this scenario. This is a typical case in cryo-EM reconstructions, even if some views are dominant. In the well sampled cases, all the structure factors remain at play, so we expect that the formulae relating SSNR to SCF are reasonably accurate. The situation is presented in Figure 9, where we describe the effect on reconstructions for a uniform case, for side views, and for modulated side-views. For both reconstructions of HA and apoferritin, in comparison to uniform, the SSNR curves are attenuated for side sampling in accordance with the amount of sampling inhomogeneity (Figure 9A-C). Side views have a range of sampling values over the surface of a Fourier sphere of radius, k , from the on-axis values to those on the orthogonal plane with a max-min ratio of πk . For the modulated side-view case, the ratio is even larger: $\pi k / (1 - \lambda)$, where λ is the strength of the modulation. Nevertheless, the agreement between the decrement in SSNR and the SCF, as shown in the table in Figure 9D, is acceptable for both HA and Apoferritin.

5.2.2 Top-like sampling cases for varying cone sizes

1 We then proceeded to describe cases that would be reflective of a predominant top view, and for this
2 reason constructed the two-parameter family of distributions described by α and ϵ , where the
3 former represents the half-angle of the cone from which the projections are drawn, and the latter
4 represents the percentage of uniform projections besides those drawn from the cone. First, we vary
5 the size of the cone, while fixing 3% uniform sampling across the rest of Euler space. Figure 10A-
6 C shows how the SSNR is attenuated for reconstructions generated from such top-like distributions
7 containing a fixed amount of sprinkled projections. In these cases, the maximum to minimum
8 sampling can be so large as $1 + 4/(\pi\epsilon\alpha)$, for small ϵ, α according to the analytical formula.
9 Nevertheless, in Figure 10D, the multiplicative shift determined from SCF (both numerically and
10 from formulae) approximately matches the decrement in SSNR.

11

12 **5.2.3 Top-like sampling cases for fixed cone size and varying fraction of randomly sprinkled** 13 **projections**

14

15 Finally, we took the same two parameter family as in Section 5.2.2, but examined a fixed cone
16 size, and varied the fraction of random projections. Figure 11A-C shows how the SSNR is
17 attenuated for reconstructions generated from such top-like distributions containing a fixed cone
18 and varied number of random projections. The first observation from this data, as we explained in
19 Section 3, is the artefactual increase in the SSNR for cases with completely missing data (black
20 dotted curve in Figure 11A-B). This stems from the singularity in the theory for how the SSNR is
21 typically defined, and a separate formula is needed to properly account for the variation that is
22 implicitly missing, in half maps created from sets of projections with missing data. The adjusted
23 formula from Eq. (3.27) pushes the SSNR curve to the appropriate ordering of the curves, where

1 increasing the sampling always increases the SSNR (black solid curve in Figure 11A-B).
2 Theoretically, the other curves (1%, 3%, 10%) should not be adjusted, because there is sufficient
3 sampling to add information to the missing regions. In practice, there is also a small shift in those
4 curves, which is not shown for the sake of clarity. The second observation from this data is that,
5 for cases with large gaps in Fourier space, a small amount of additional projections goes a long
6 way in increasing the SSNR. This is not surprising. Even in the early days of reconstructions, it
7 was realized that, for under-sampled cases, adding small amounts of information to deficient parts
8 of Fourier space greatly improves the ability to solve the reconstruction problem [36]. The
9 experimental attenuations of the SSNR are also in line with the geometrical decrement of the SCF
10 in continuum calculations (compare Figures 11 and 6). As in the previous cases described above,
11 Figure 11D shows that the multiplicative shift determined from SCF (both numerically and from
12 formulae) approximately matches the decrement in SSNR.

13

14

1 **Discussion**

2 In this work, we show that non-uniformity of the set of projection views drives down properly
3 defined global resolution measures. Our calculations are based on standard assumptions, that there
4 is some envelope that seems to stabilize for values less than 10 Angstroms [19]. The SSNR
5 resolution measure estimates the ratio of the signal power to the signal variance. Using ordinary
6 statistics, we expect that the variance per voxel will be decremented by the sampling. Therefore,
7 if we assume that the noise variance approximately decouples from the sampling, then the average
8 over Fourier shells of the reciprocal sampling arises naturally in the expression for the SSNR,
9 leading to Eqs (1.1) to (1.2). Thus, the measure for the efficacy of sampling that we advocate, the
10 SCF, emerges naturally, if we wish to isolate the effects of the geometry of the sampling on the
11 resolution. The incorporation of the SCF is the step that distinguishes our calculations from similar
12 calculations, such as [23].

13

14 A typical cryo-EM reconstruction procedure carries along information that can be represented by
15 three maps: two half-maps and a sampling map that can be created from knowledge of the angle-
16 assignments or that can be taken to be the map of reconstruction weights in a direct Fourier
17 reconstruction. From these maps, one can estimate up to second moments and continue to combine
18 information to arrive at more refined reconstructions. Ultimately, one arrives at a mean map,
19 variance map, and sampling map, or three pieces of information per voxel. If there is missing data,
20 then there is a pathology in the way that SSNR is typically defined. Although defining the mean
21 of the missing values to zero is acceptable (and forms the best estimate of the original structure),
22 setting the variance to zero is illogical, since there is enough information to give a better estimate.

1 We find a self-consistent correction to the ordinary SSNR and showed in section 5 that the
2 redefined SSNR always increases with more uniform sampling, as should be expected.

3

4 We also demonstrated the linear dependence of the SSNR on the total sampling, which is governed
5 by the number of particles. This was implicit in earlier analyses of Guinier or Reslog, as shown
6 in the mathematical description of section 4, but takes on a simpler form here. We show that these
7 latter constructions imply a definite functional form for the SSNR, which is more restrictive than
8 necessary. Indeed, we provide the mathematical argument, that one can estimate the number of
9 particles necessary to achieve a higher resolution, using the same collection strategy, but with a
10 single SSNR curve, provided that the curve is sufficiently continuous over the resolution ranges in
11 question. This has value during data acquisition, since it can inform the experimentalist how a
12 given collection might be altered or abandoned based on the goals of the experiment, and the
13 prediction is achieved without the need to recalculate reconstructions using particle subsets.

14

15 There are several major implications from the current work. Most importantly, the direct
16 relationship between sampling and global resolution in single-particle cryo-EM implies that any
17 deviation from uniformity *always* drives down the SSNR, and thereby leads to an increase in the
18 number of particles that are required for attaining a specified resolution. There is a persistent
19 problem of preferred specimen orientation (and consequently non-uniform projection
20 distributions) that appears to affect the vast majority of single-particle reconstructions [6]. This
21 means that virtually all data sets are characterized by non-ideal imaging and image processing

1 conditions. As dictated by Eq. 1.1 (also Eq. 2.20) the experimental situation therefore requires
2 optimizing two parameters – the experimental “envelope” *as well as* the sampling distribution.
3 Here, we use the term “envelope” in a broad sense to encompass all of the factors that attenuate
4 experimental resolution. These include, but are not limited to, beam coherence, ice thickness (and
5 its effect on the background signal-to-noise ratio), quantum efficiency of the detector, residual
6 specimen movement that is not corrected by motion correction, errors in computational orientation
7 assignment, structural heterogeneity, and any other factors that generally attenuate experimental
8 resolution, as measured by the FSC. In addition to the envelope, the sampling distribution
9 matters. To reach the hypothetical resolution limit for small particles [37], it is therefore essential
10 to not only improve hardware and software, but also techniques for specimen preparation, in order
11 to maximize sampling uniformity on cryo-EM grids. Some effort toward this goal is ongoing [7],
12 but more needs to be done. Along these lines, the more symmetric the particle, the more
13 orientations are sampled during the reconstruction process. Therefore, symmetry does not merely
14 multiply the number of particles in the data in accordance with the symmetry group; the
15 improvement in sampling for symmetric particles also contributes to gains in SSNR by virtue of
16 improvements to the SCF. Thus, symmetry has a dual effect in improving both data quantity and
17 quality. In part for this reason, cases like AAV [2] and Apoferritin [3] have pushed the resolution
18 limits and are associated with very low temperatures factors (or slowly decreasing envelopes) in
19 the data.

20

21 Beyond attenuation of global resolution, the extent to which the map suffers as a consequence of
22 incomplete sampling is currently unclear. Specimens with high C- or D-fold symmetry that are

1 characterized by pure side views are, strictly speaking, anisotropic. However, the effect at the level
2 of the reconstructed map is negligible, and the experimentalist should not notice differences in
3 structural details if one were to directly compare to a map reconstructed from a uniform sampling
4 distribution. Nonetheless, as we show in figure 9 and emphasize throughout this work, the SSNR
5 for pure side views is still attenuated in comparison to uniform by ~20%, and thus the amount of
6 data required for reaching certain resolutions is increased by approximately the same
7 percentage. Beyond the simple cases, there are multiple factors that currently complicate an
8 exhaustive analysis of experimental maps characterized by different symmetries and sampling
9 geometries. First, it will be necessary to decouple the effect of anisotropy (in its strict definition,
10 impacting directional resolution) from the attenuating effect on *global* resolution. More worryingly
11 however, we believe that there may, in certain extreme cases of missing data, be systematic bias
12 in the reported resolution in the field, caused by artefactual inflation in the FSC (for example, as
13 observed in figure 11). In part for this reason, we introduced the FSC* and SSNR* criteria, which
14 compensate for missing views in Euler space and report a more realistic value of resolution for the
15 pathological cases. FSC* and SSNR* can, in principle, be extended to highly under-sampled
16 orientations that may be prevalent in experimental situations. Implementation of these criteria to
17 experimental data, and a careful analysis of the underlying sources and resulting statistics, will be
18 the subject of future work.

19

20 Experimental improvements to the sampling distribution can be achieved by tilting the specimen
21 inside of the electron microscope. However, this comes at a cost of degradation in image quality
22 [15]. The direct relationship between sampling and resolution indicates that any attenuation due to

1 sampling can now be compared with other types of experimental attenuations, for example due to
2 beam-induced movement, ice thickness, errors in the image processing pipeline, etc. Thus, a
3 natural direction will be to quantify the resolution gains caused by improvements in orientation
4 sampling, as compared to resolution losses caused by degradation of image quality during tilted
5 data acquisition. Such studies will help to quantitatively establish an optimal tilt angle for any
6 dataset containing a given sampling distribution.

7

8 Finally, The SCF provides a direct means by which to evaluate a sampling distributions, with an
9 intuitive scale ranging from 0 to 1. We propose the use of the SCF for evaluating Euler angle
10 assignments for sets of particles that produce 3D reconstructions in cryo-EM.

11

12 **Acknowledgements**

13

14 PRB would like to thank Pawel Penczek for conceptual explanations pertaining to the current
15 work. PRB and DL would like to thank David De Rosier for discussions. PRB and DL are
16 supported by grants from the NIH: DP5-OD021396, R01 AI136680, and U54 GM103368.

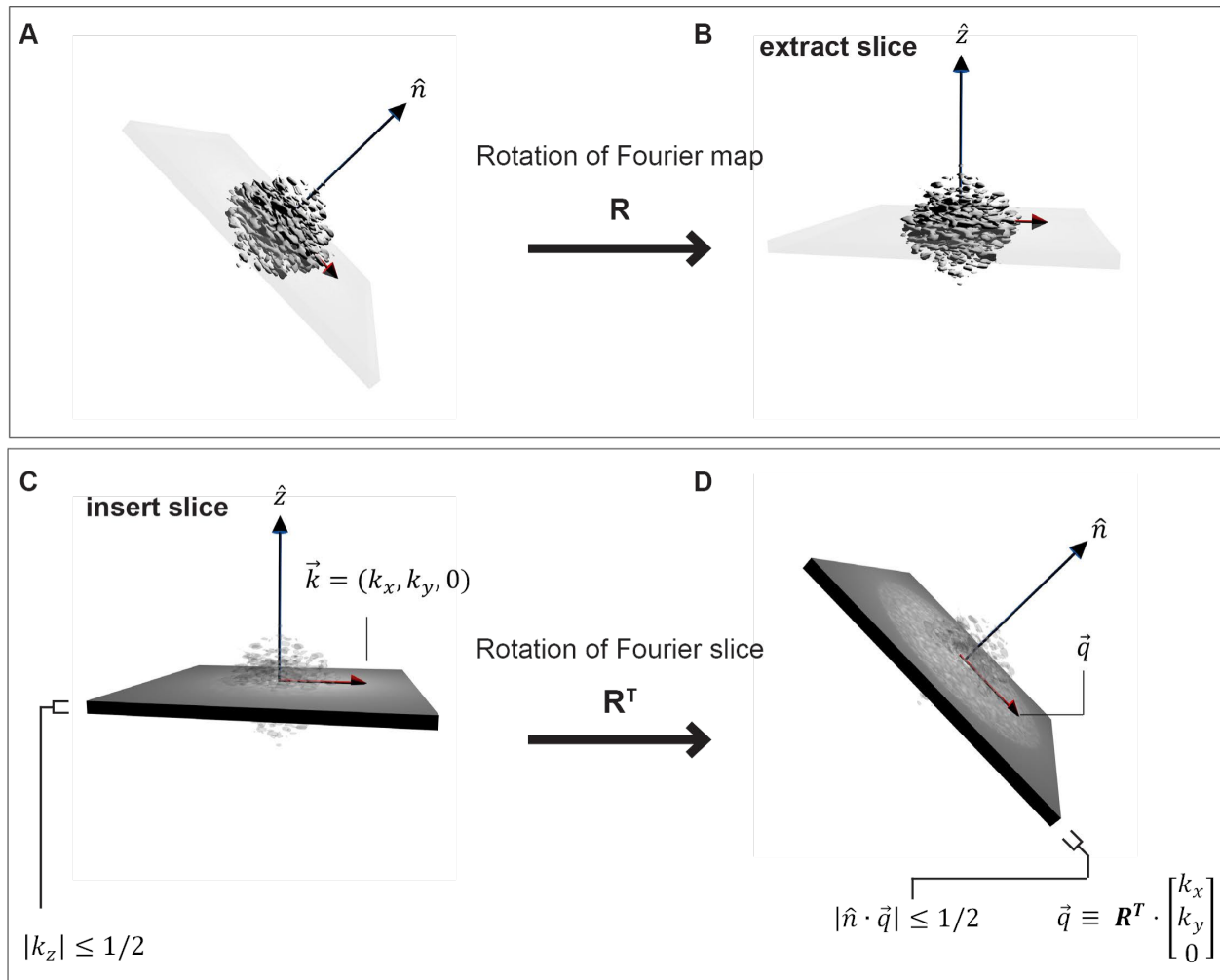
17

18

19

20

21

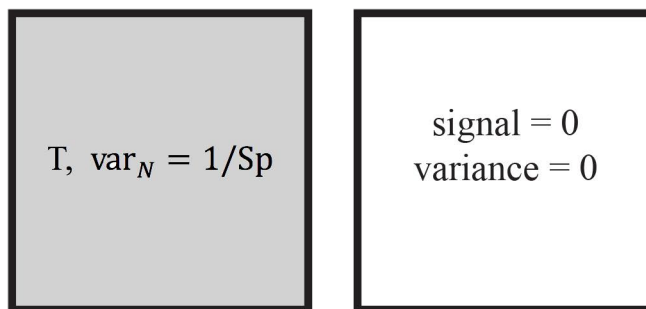


1
2

3 **Figure 1. Geometry of projections in Fourier space.** (A) A 3D object in its Fourier space representation
 4 is rotated by \mathbf{R} , and (B) a slice is extracted from the 3D Fourier transform (FT). Based on the Fourier slice
 5 theorem, selecting a 2D slice out of a 3D FT is equivalent to orthogonally projecting the original real-space
 6 map along the new \hat{z} axis. (C) The data in a projection is contained in a slab of Fourier space of unit height.
 7 When considering what the data in a 2D projection, $F(\vec{k})$, corresponds to in 3D, it is easiest to consider the
 8 mapping \mathbf{R}^T , as shown from (C) to (D). Now, the coordinates on the 3D FT as shown in D are clear: $\vec{q} =$
 9 $(k_x, k_y, 0) \cdot \mathbf{R}$. The slab condition on the projection, $|k_z| < 1/2$ readily translates into the
 10 condition $|\hat{n} \cdot \vec{q}| < \frac{1}{2}$, where $\hat{n} = \mathbf{R}^T \hat{z}$.

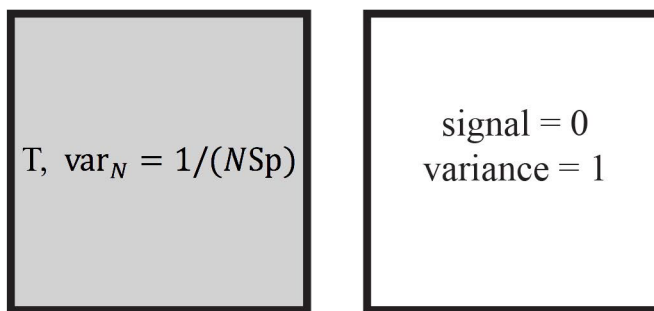
11

A P measured values Q unmeasured values



$$\text{SSNR} = \frac{P T^2}{P \text{var}_N} = \frac{T^2}{\langle \frac{1}{Sp} \rangle_P} = N \frac{T^2}{\langle \frac{1}{sp} \rangle_P}$$

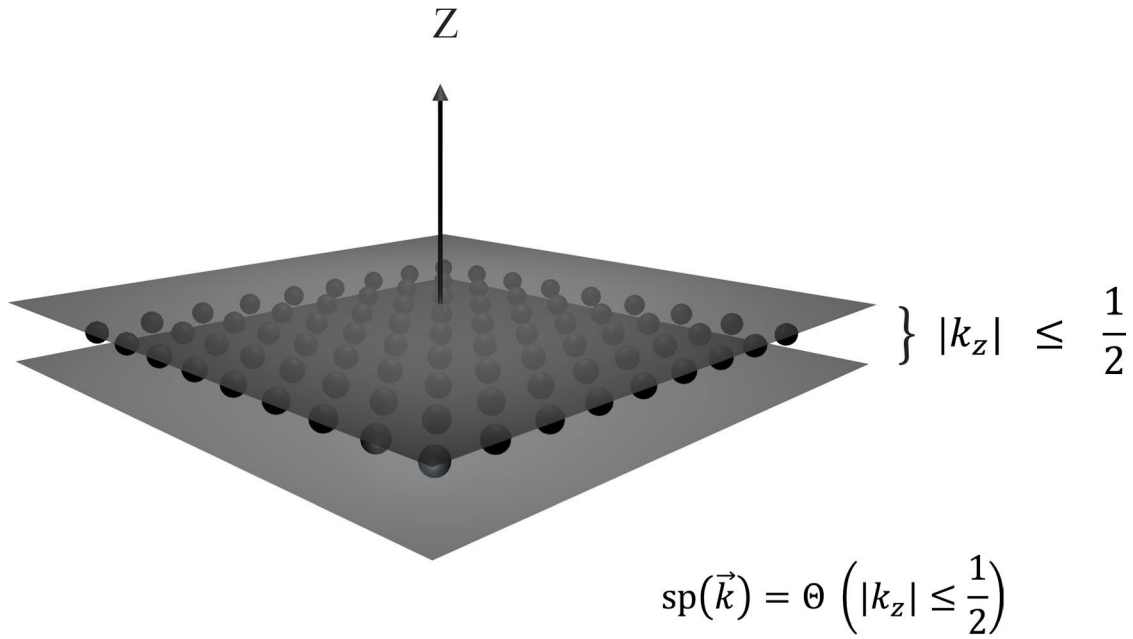
B P measured values Q unmeasured values



$$\text{SSNR}^* = \frac{P T^2}{P \text{var}_N + Q \text{var}_Q} = \frac{P T^2}{P \langle \frac{1}{Sp} \rangle + Q} = \frac{NP T^2}{P \langle \frac{1}{sp} \rangle + NQ}$$

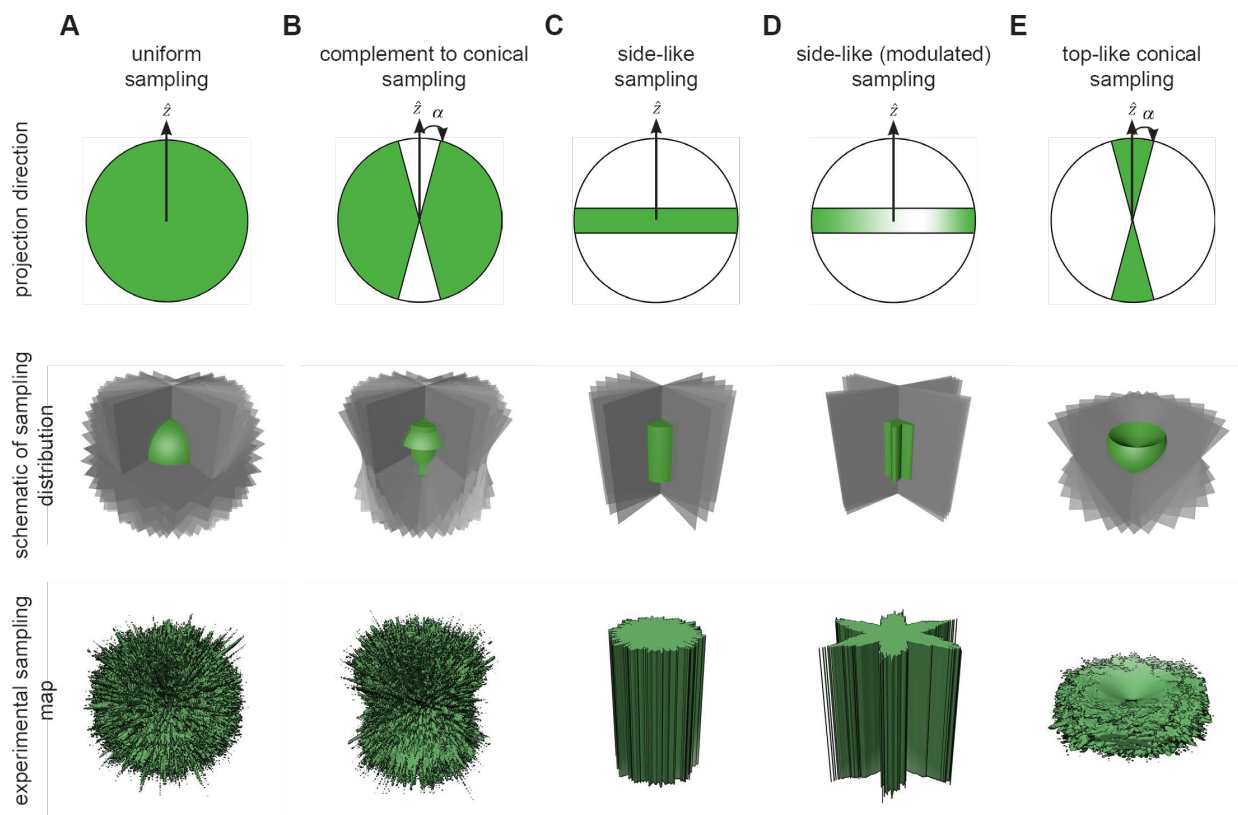
1
2
3
4
5
6
7
8
9
10
11

Figure 2. Schematic of the proposed adjustment to the SSNR formula for cases with unmeasured regions of Fourier space. The shaded areas represent the P values on a Fourier sphere that have been measured for a target value, T , in the measured region. The variance in the measured region is down-weighted by the total sampling, Sp , which is N times the per-particle sampling, sp . Meanwhile, the unshaded region represents Q unmeasured values. In (A) the unmeasured voxels are assigned zero variance, whereas in (B) the voxels are assigned variance consistent with the already measured voxels, resulting in two different expressions for SSNR. The expression in (A) limits to arbitrarily large values as the number of particles increases, whereas in (B) the expression saturates.



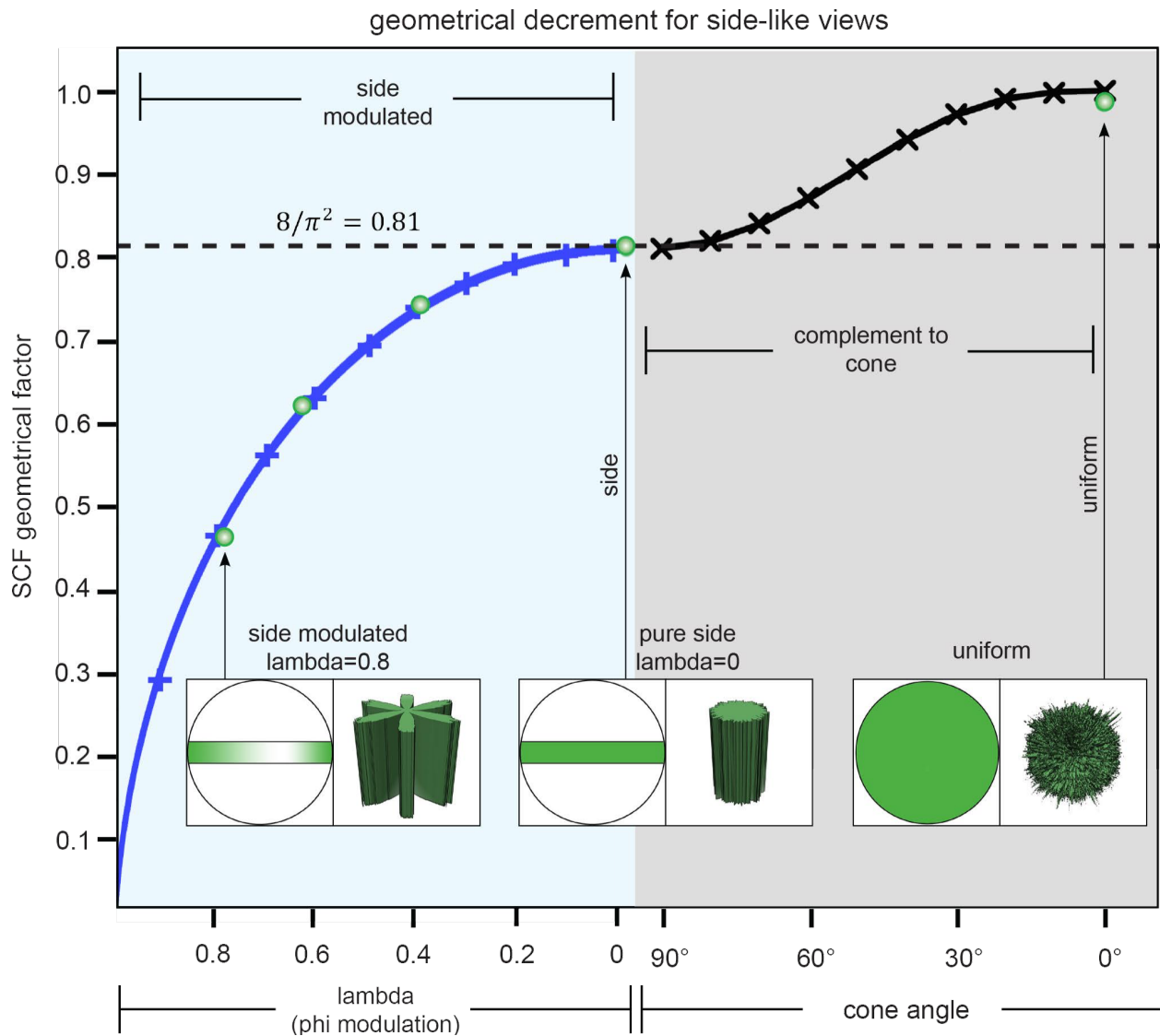
1
2
3
4
5
6
7
8
9
10

Figure 3. Numerical and analytical representations for sampling. A representation in Fourier space of a single projection with Fourier magnitudes less than L . Lattice sites, shown as blue dots, in the $k_z = 0$ plane are considered to be sampled. Our coarse-grained approximation is that the entire slab of unit width containing those points are sampled: $sp(\vec{k}) = \Theta \left(|k_z| \leq \frac{1}{2} \right)$, where Θ is the indicator function. Our numerical algorithm hinges on finding lattice sites that satisfy the generalization of this condition: $sp(\vec{k}) = \Theta \left(|\hat{n} \cdot \vec{k}| \leq \frac{1}{2} \right)$ for any arbitrary projection direction given by \hat{n} . For analytical calculations, we further approximate this as a delta function $sp(\vec{k}) = \delta(\hat{n} \cdot \vec{k})$, which satisfies the proper normalization.



1
2
3 **Figure 4. Projection distributions used to evaluate sampling.** The first row denotes the directions of the
4 projections (green); the middle row provides a schematic of the sampling (Fourier slices are in gray and the
5 sampling map schematic is in green); the last row shows the experimental sampling map from 10,000 slices
6 inserted into the 3D FT. (A) The uniform sampling distribution evenly covers the entirety of reciprocal
7 space. (B) “side-like” projections are uniformly drawn from the complement to a cone of half angle, α . The
8 uncovered region is orthogonal to the cone and lies along the X/Y plane. (C) Projections are drawn from
9 the side, which corresponds to the $\alpha \rightarrow \pi/2$ limit of case B (Euler angle $\theta=90^\circ$). (D) In addition, we
10 incorporated azimuthal oscillations (Euler angle ϕ is modulated), which increases the fluctuation of side
11 view sampling. The depth of the oscillations is governed by a parameter, λ , and we term this scenario
12 modulated side-views. (E) Projections are drawn from within a cone of half angle α . Unlike the other cases,
13 the top-like distributions always have missing conical regions of Fourier space related to the size of the
14 half-angle α . For this reason, we ultimately include an additional parameter, ϵ , which represents the fraction
15 of projections scattered randomly over the projection sphere.
16

1

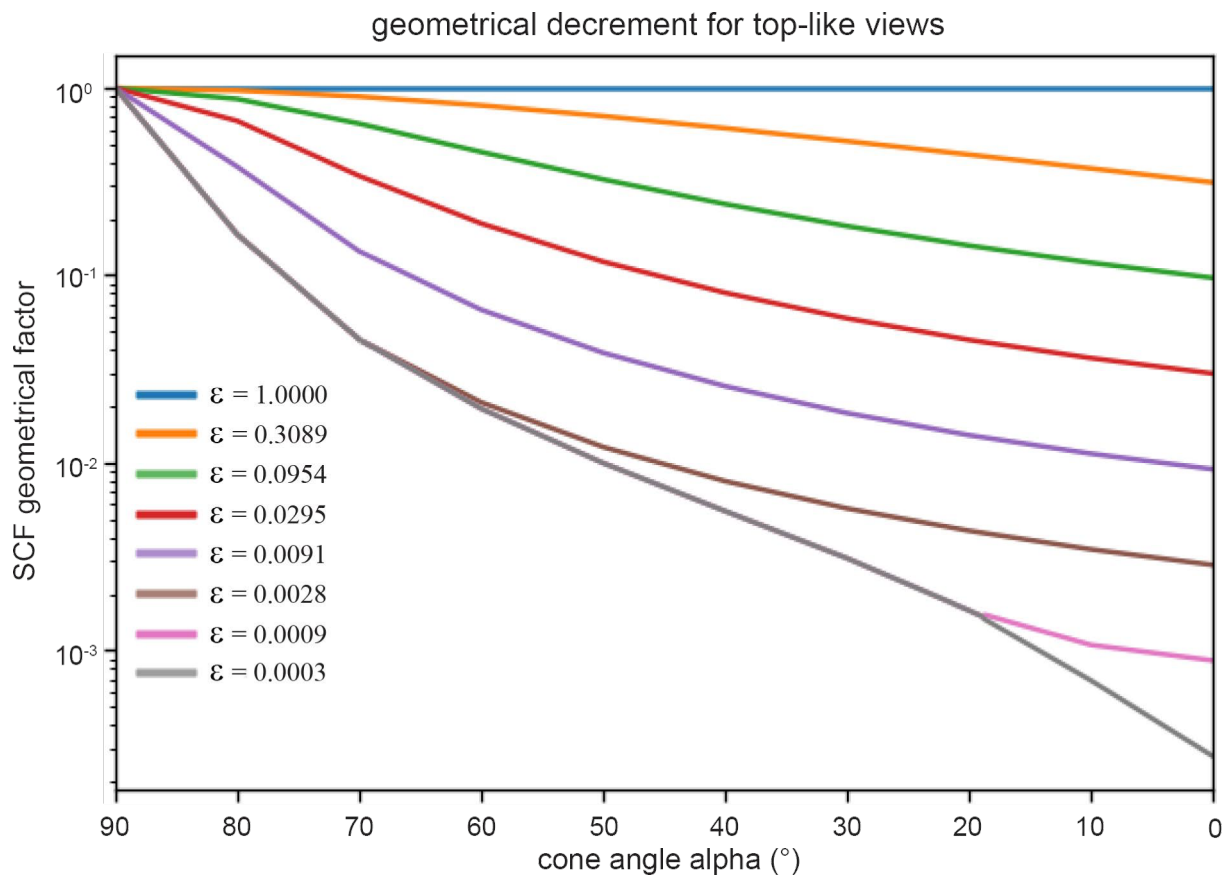


2

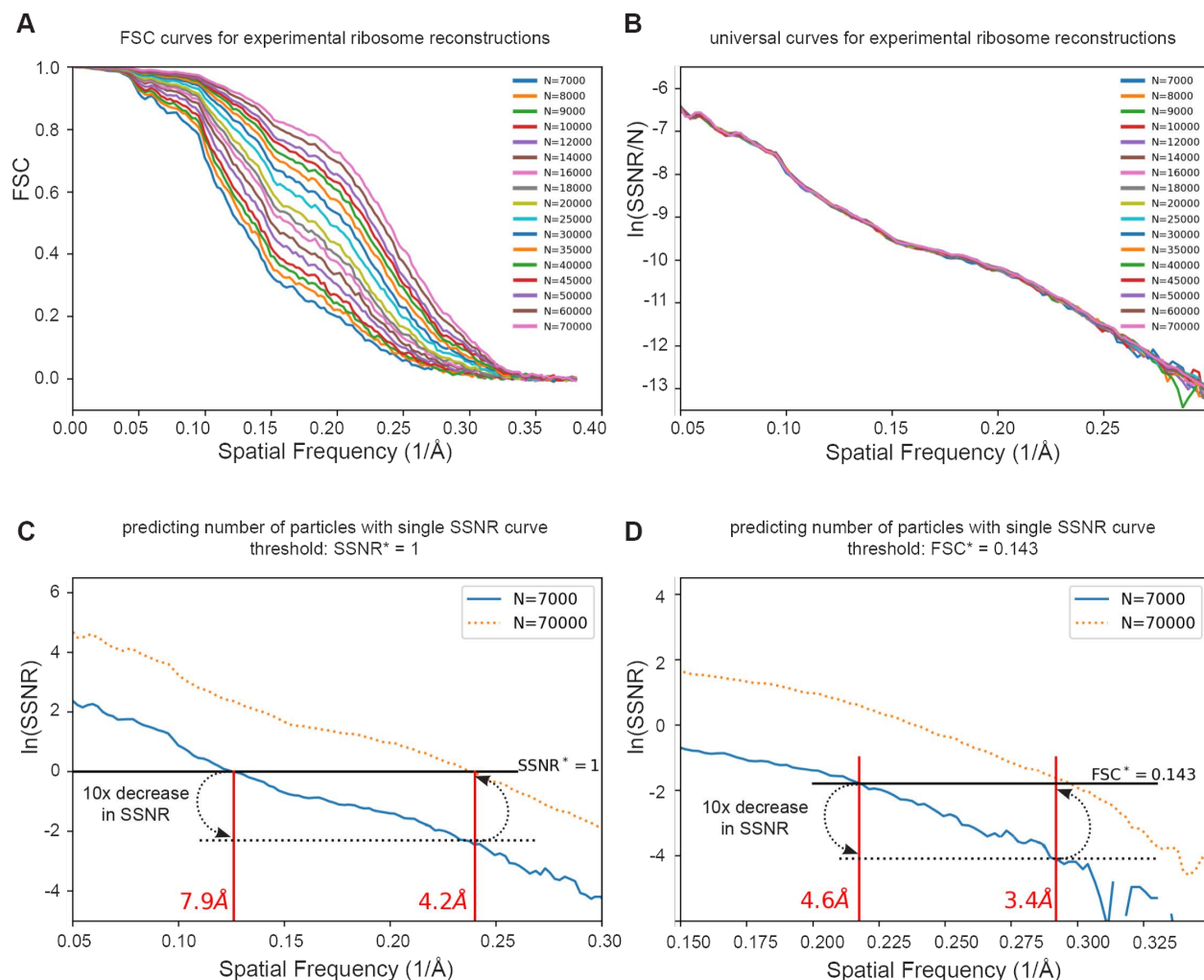
3

4 **Figure 5. Geometrical decrement of the SCF for side-like views characterized by thorough sampling.**

5 The SCF is a single geometrical factor, which forms an approximate estimate for how the SSNR is
6 decremented due to deficits in sampling. We plotted $SCF_{\text{side modulated}}(\lambda)$ for modulated sets of side views
7 (see Eq (3.23)), as well as $SCF_{\text{side-like}}(\alpha)$ for side-like views (see Eq (3.24)). The approximations inherent
8 in Eq. 3.23 are no longer valid for λ very close to 1, and therefore the plot is not shown for $\lambda < 1$, $SCF >$
9 0. The projection views and sampling maps are shown in three typical cases: (i) $\lambda = 0.8$ side-modulated,
10 where there is a deep pocket in the sampling, (ii) pure side views, where the contours of equal sampling are
11 cylinders, and (iii) uniform views, where the SCF attains its maximum value of 1. Green open circles
12 represent the numerical evaluations of the SCF, and their correlation with our continuum calculations
13 (represented by the curves) reinforces the efficacy of both approaches.



1
2
3 **Figure 6. Geometrical decrement of the SCF for top-like views characterized by poor sampling.** Plot
4 shows the predicted SCF as a function of cone size and fraction of randomly sampled projections. The
5 attenuation of SCF due to conical sampling is typically more severe than for well-sampled cases. The angle,
6 α , represents the cone (in degrees) that predominantly contains the projections, except for a fraction, ϵ ,
7 which represent the percentage of projections that are distributed uniformly over the projection sphere.
8 Typical distributions are shown later in Figure 8. The curve that bounds the whole set from below, is given
9 by Eq. (3.28), but is otherwise given by Eq. (3.26): the plotted SCF* is the maximum of the two different
10 expressions. This crossover between expressions is given at the bounding curve (gray) when $\epsilon(\alpha) \cong$
11 $\sin(\alpha) \frac{2k}{N}$.
12

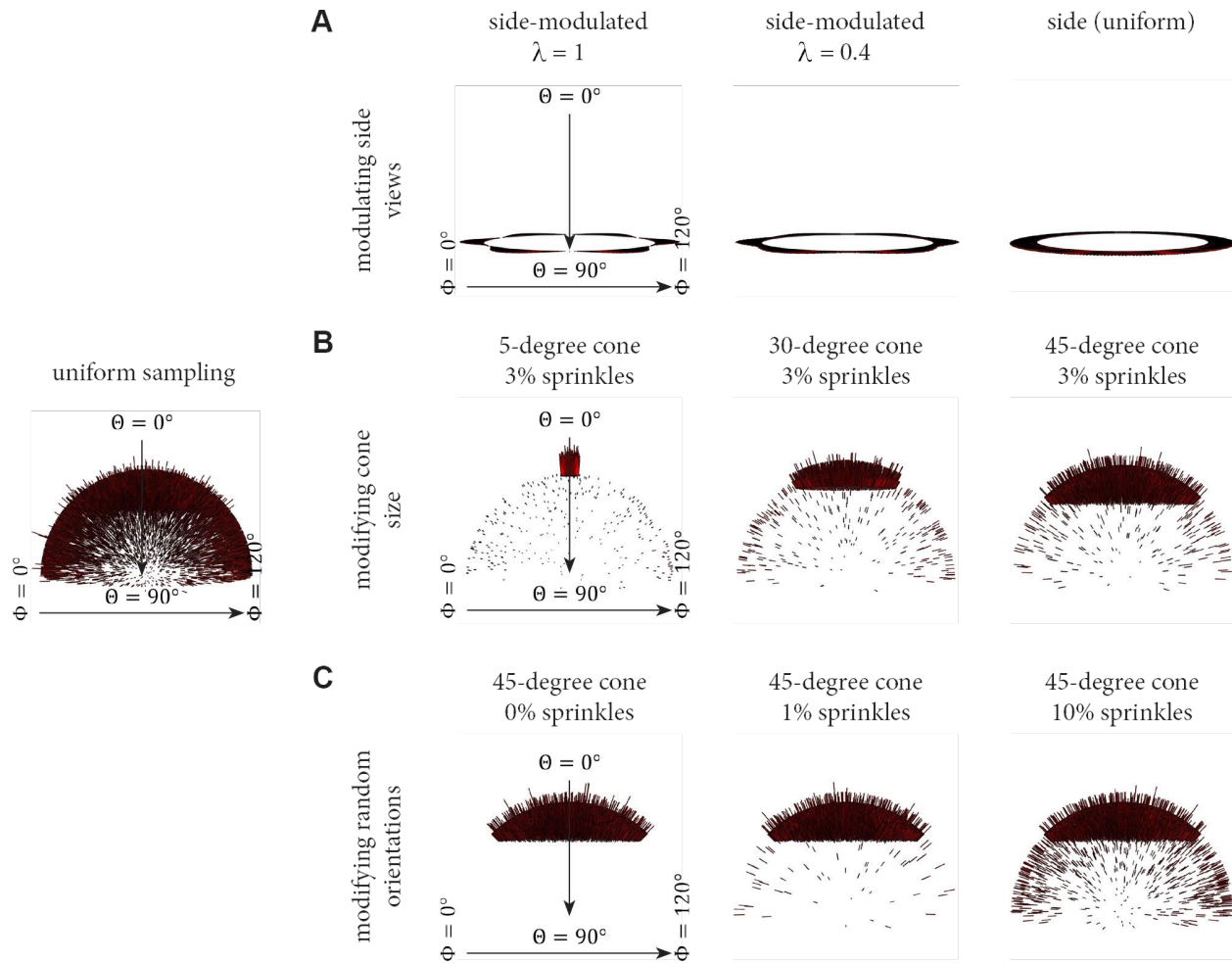


1

2

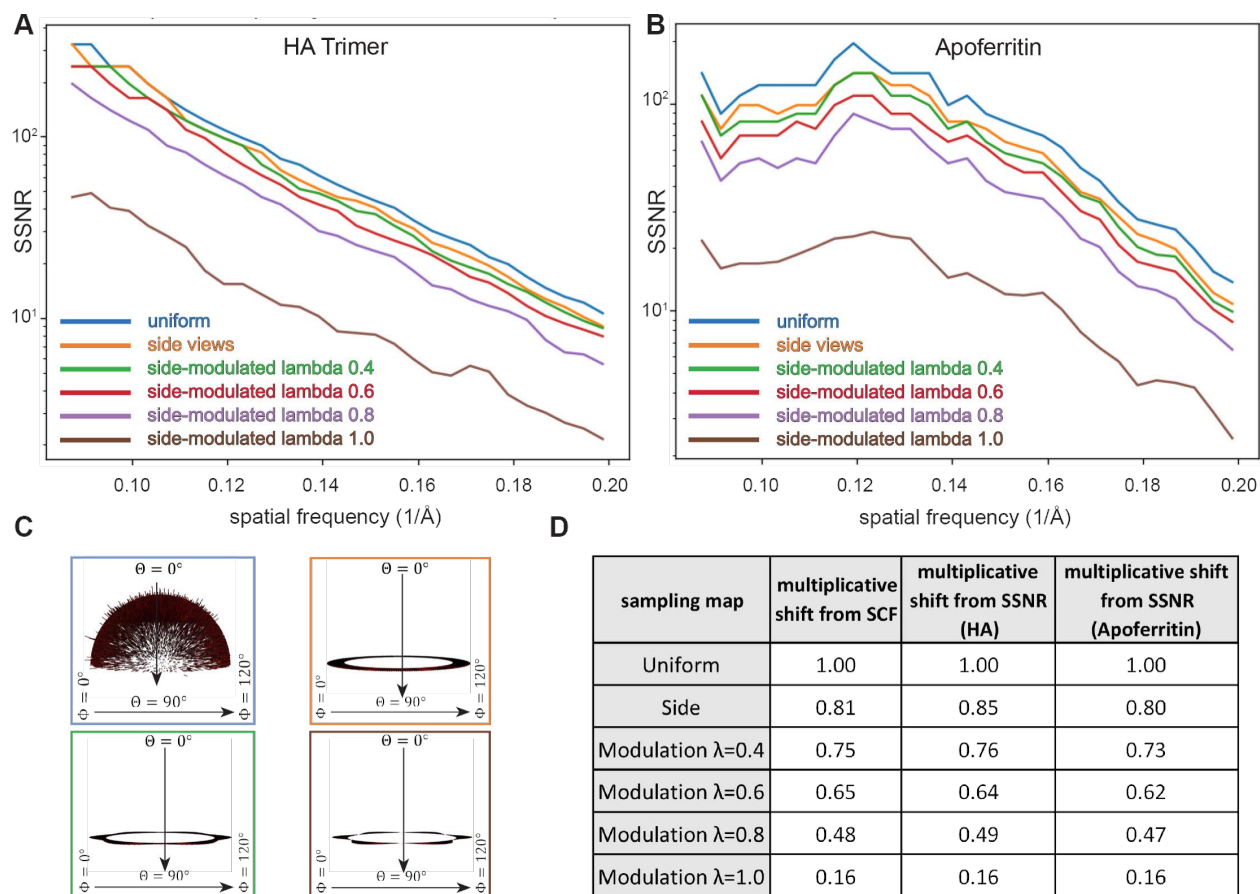
3 **Figure 7. Predicting the number of particles based on per-particle SSNR curves** (A) We reexamined
 4 17 of the FSC curves with the largest numbers of particles from Passos et al [26]. (B) A plot of the per-
 5 particle $\ln(SSNR/N)$ shows that these 17 distinct curves collapsed approximately to a single curve. (C-D)
 6 Since each curve contains essentially the same information, we can estimate the number of particles needed
 7 to achieve a target resolution. For example, one may wish to know the approximate experimental resolution
 8 by increasing the number of particles by tenfold from 7,000 (solid blue line) to 70,000. (C) using the
 9 $SSNR^*=1$ threshold (solid black line), one would find the intercept corresponding to a 10x decrease in
 10 $\ln(SSNR)$ (dotted black line), and plot that back onto the solid $SSNR^* = 1$ line. The experimental $\ln(SSNR)$
 11 curve for 70,000 particles (dotted orange line) shows a correspondence between the prediction and the
 12 experimental intercept. (D) The same argument approximately holds for the $FSC^* = 0.143$ criterion, or for
 13 other thresholds.

14



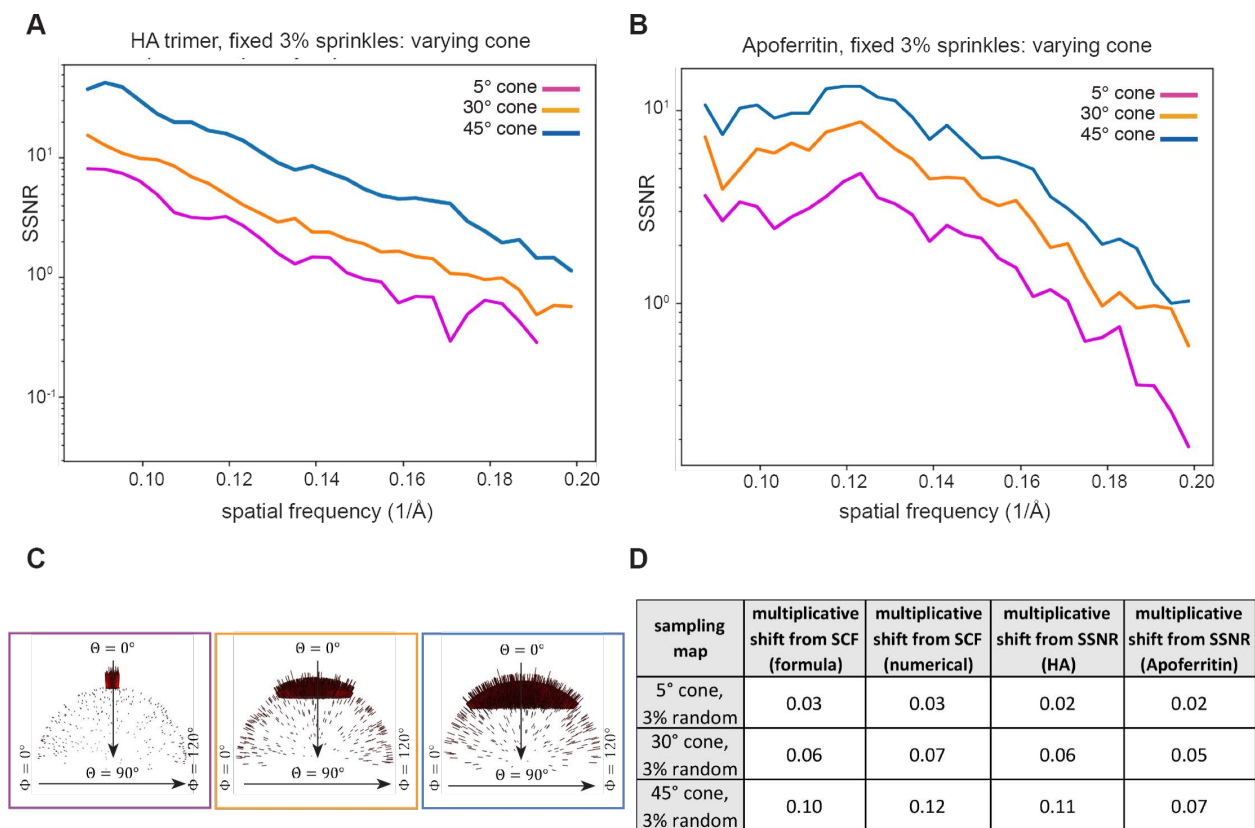
1
2
3
4
5
6
7
8
9
10
11
12
13
14

Figure 8. Euler angle sampling schemes used in this study. The projection distributions described in Figure 4 can be modulated to vary specific parameter and evaluate distinct conditions. For each scenario, Euler angle distributions for 10,000 projections are displayed. A uniform distribution of views across Euler space is shown at left for comparison. (A) The side-modulated case, whereby the Euler angle $\theta=90^\circ$, but ϕ is modulated in accordance with the modulation parameter, λ . This scenario corresponds to the transition between Figure 4C-D. (B) The top-like case, whereby the size of the cone is varied, and there is a fixed amount of randomly distributed views across the rest of Euler space. (C) The top-like scenario, whereby the size of the cone remains constant at 45° , but the amount of randomly distributed views is varied across Euler space. The experimental results corresponding to these three cases will be described in Figures 9-11.



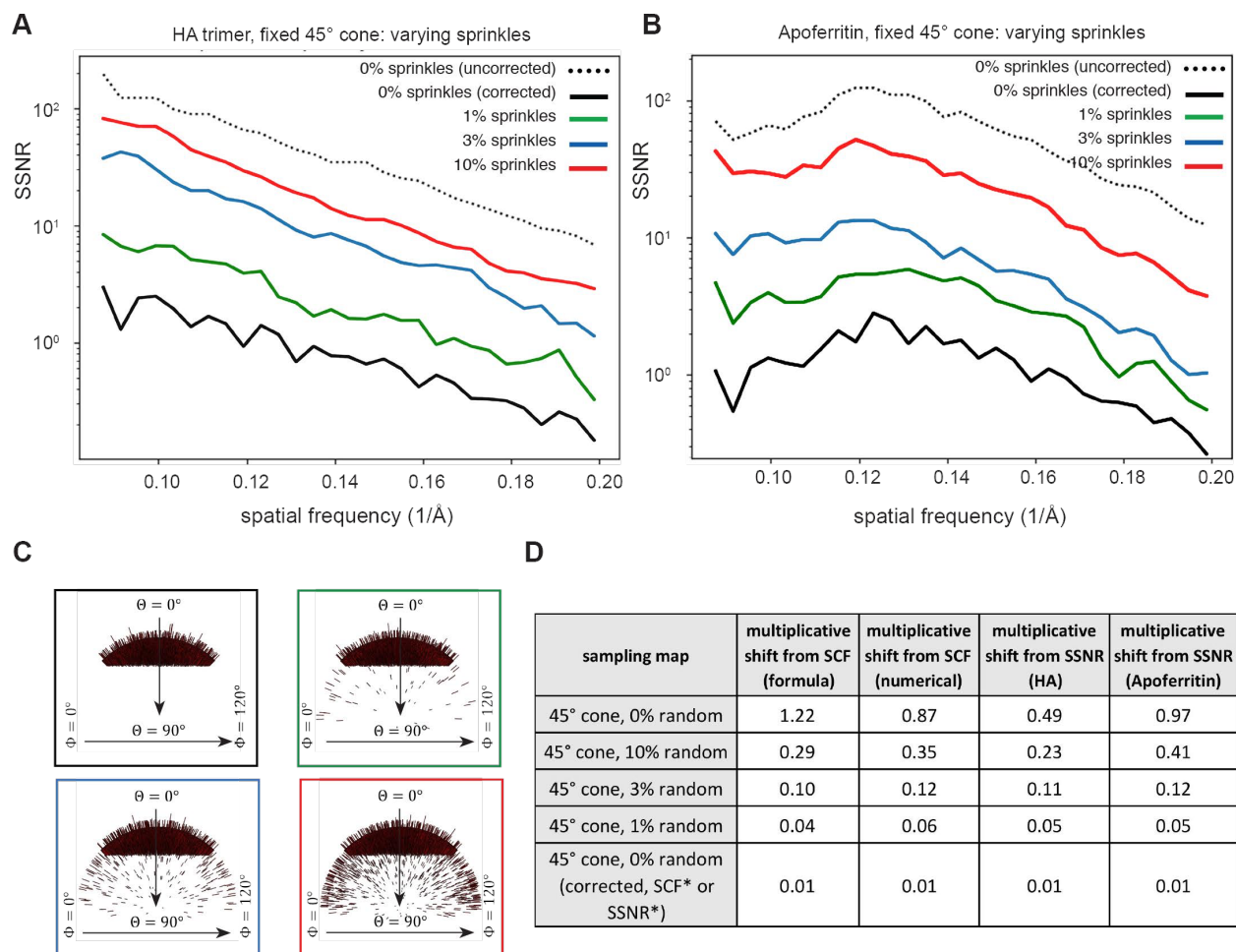
1
2
3 **Figure 9. Attenuation of the SSNR due to the projection distribution for side-like cases characterized**
4 **by thorough sampling.** (A-B) Two synthetic datasets corresponding to (A) the HA trimer and (B)
5 Apoferritin were reconstructed, as described in section 5. The decrement, which is reflected by the shift in
6 the SSNR due to the different type of sampling distributions is displayed. (C) Euler angle distribution
7 profiles corresponding to selected cases from A-B. (D) Table showing the decrement in SCF and the
8 multiplicative shift from SSNR. For each of the 6 different sampling distributions, the numerical and
9 analytical forms for the SCF agree (except for $\lambda = 1$, where we only have a numerical formula), and thus
10 only a single multiplicative shift from SCF is indicated.

11
12
13
14



1
2
3 **Figure 10. Attenuation of the SSNR due to the projection distribution for top-like cases characterized**
4 **by varying cone sizes.** As in Figure 9, for both HA and Apoferritin, we examined the effect of the SCF on
5 the SSNR for poorly sampled cases, where the projection distributions were constrained to varying cone
6 sizes, but the number of random projections was fixed at 3%. This leads to Fourier space being sparsely
7 sampled around the z-axis in each case. **(A-B)** Decrement of the SSNR for **(A)** HA and **(B)** with **(C)** the
8 corresponding Euler angle distribution profiles. Even though there is a low percentage of views in certain
9 regions, the total sampling is not small, because the total number of particles is 10^4 . **(D)** Table showing the
10 decrement in SCF. In each case, there is a crude agreement between the analytical expectation for the SCF,
11 the numerically calculated SCF, and the shifts of the SSNR.

12
13



1
2
3 **Figure 11. Attenuation of the SSNR due to the projection distribution for top-like cases characterized**
4 **by varying random views.** As in Figures 9, for both HA and Apoferritin, we examined the effect of the
5 SCF on the SSNR for poorly sampled cases, where the projection distributions were constrained to the fixed
6 cone size of 45°, but the percent of random projections varies from 0%, 1%, 3% to 10%. **(A-B)** Decrement
7 of the SSNR for **(A)** HA and **(B)** Apoferritin with **(C)** the corresponding Euler angle distribution profiles.
8 **(D)** Table showing the decrement in SCF. For the first four rows in the table (and the corresponding SSNR
9 curves in A-B), the SCF, as defined theoretically by Eq (3.25) for 0%, and Eq (3.26) for the 1%, 3% and
10 10% cases, and numerically by Eq (3.22), approximately describes the observed change in the SSNR, as
11 given by the last two columns. There is the one serious issue, as discussed in the text, that the SSNR with
12 completely empty regions of Fourier space is significantly higher for 0% uniform rather than 1%. The text
13 explains a logical correction, given by the SCF theoretically by (3.27) and numerically by the same
14 algorithm. The correction to the SSNR is shown in the last row and appropriately predicts a large
15 multiplicative shift from uniform, as would be expected for such a poorly sampled case.

1 **Appendix A. Some details for calculations in Section 3: geometrical factor for**
 2 **decay of density Eq (3.5), checking numerical sampling code Eq (3.9), creating**
 3 **distributions according to some prescribed function Eq (3.15), proof that**
 4 **uniform distribution maximizes the SCF Eq (3.22), derivation of Eq. (3.23) for the**
 5 **SCF for modulated side-views**

6
 7 **A.1 A general formula for the projection geometrical factor: Eq. (3.5)**

8 Our claim in Eq (3.5) is that

9
 10
$$\langle \delta(\hat{n} \cdot \hat{k}) \rangle_{\hat{n}} = \frac{1}{c_p k} , \quad (\text{A.1})$$

11 where c_p is a geometrical factor that we wish to calculate in general dimensions, especially for
 12 $D = 2, 3$. By $\langle \cdot \rangle_{\hat{n}}$, we mean the average over the surface of the unit ball in D dimensions. One
 13 easy way is to integrate the above equation over all \vec{k} with $k < L$ in D dimensions. Then on the
 14 left-hand side we get:

15
 16
$$\int^{\vec{k}, k \leq L} \langle \delta(\hat{n} \cdot \hat{k}) \rangle_{\hat{n}} = \langle \left(\int^{\vec{k}, k \leq L} \delta(\hat{n} \cdot \vec{k}) \right) \rangle_{\hat{n}} , \quad (\text{A.2})$$

17
 18
 19
$$= \int^{\vec{k}, k \leq L} \delta(k_z) , \quad (\text{A.3})$$

20
 21
 22
$$= \int^{\vec{k}, k \leq L, D-1 \text{ dim}} 1 , \quad (\text{A.4})$$

23
 24
$$= L^{D-1} V_{D-1} , \quad (\text{A.5})$$

25 where V_{D-1} is the volume of the unit ball in $D - 1$ dimensions. Eq (A.3) holds because the
 26 integrand in (A.2) is no longer dependent on the direction, \hat{n} , so the average over \hat{n} seen in (A.2)
 27 integrates to 1. Moreover \hat{n} where it appears in the integral may be set to \hat{z} for convenience. On
 28 the RHS of (A.1) we also perform the integration over the ball of radius L and get

1
$$\int_{\vec{k}, k \leq L} \frac{1}{c_p k} = \frac{1}{c_p} A_D \cdot \int_0^L dk k^{D-1} \frac{1}{k} , \quad (\text{A.6})$$

2
3
$$= \frac{1}{c_p} A_D \cdot \int_0^L dk k^{D-2} , \quad (\text{A.7})$$

4
5
$$= \frac{1}{c_p} A_D \cdot \frac{L^{D-1}}{D-1} , \quad (\text{A.8})$$

6 where A_D is the surface area of the unit ball in D dimensions .

7

8 Equating the last two expressions shows that

9
$$c_p = \frac{A_D}{(D-1) \cdot V_{D-1}} = \frac{A_D}{A_{D-1}} . \quad (\text{A.9})$$

10 This gives

11
$$c_p(D=2) = \frac{2\pi}{2} = \pi , \quad (\text{A.10})$$

12 and

$$c_p(D=3) = \frac{4\pi}{2\pi} = 2. \quad (\text{A.11})$$

13

14 Therefore, the geometrical factor 2 that appears in Eq 3.5 is simply the ratio of the surface area
15 of a unit ball to the circumference of a great circle of the same ball.

16

17 **A.2 Checking the Sampling Code Eq. (3.9)**

18 We want to evaluate

19
$$S = \int_{-L}^L dk_x \int_{-L}^L dk_y \int_{-L}^L dk_z \left(\frac{1}{2k} \right) \quad (\text{A.12})$$

1
2

$$S = \left(\frac{1}{2}\right) \int k \, dk \int_0^\pi \sin \theta \, d\theta \int_0^{2\pi} d\phi \quad \Theta(|k_x|, |k_y|, |k_z| \leq L) \quad (\text{A.13})$$

3 It is enough to consider the upper quadrant, where all the components are positive.

4 This is where both the azimuthal angle, ϕ , and the spherical angle, θ , are in the range $\left[0, \frac{\pi}{2}\right]$. This

5 gives us a symmetrization factor of 8. However, we may also consider a definite ordering for the

6 k_x, k_y, k_z giving a symmetrization factor of 6. Putting this together we have

8
7

$$S = \left(\frac{8 \cdot 6}{2}\right) \int k \, dk \int_0^{\pi/2} \sin \theta \, d\theta \int_0^{\pi/2} d\phi \quad \Theta(0 \leq k_y \leq k_x \leq k_z \leq L) \quad . \quad (\text{A.14})$$

9 We wish to reorder the integrations: first k , then θ then ϕ . The spherical representations for the

10 components may be written as: $k_x, k_y, k_z \equiv (\sin \theta \cos \phi, \sin \theta \sin \phi, \cos \theta)$

11 Now $0 \leq k_y \leq k_x$ is easily represented by $0 \leq \phi \leq \frac{\pi}{4}$. Let's write down what we have so far:

12

$$S = \left(\frac{8 \cdot 6}{2}\right) \int_0^{\pi/4} d\phi \int_0^{\pi/2} \sin \theta \, d\theta \int k \, dk \quad \Theta(k_x \leq k_z \leq L) \quad . \quad (\text{A.15})$$

13 To ensure the last two inequalities we need $k \leq L/\cos \theta$ and $\tan \theta \cos \phi \leq 1$. This last

14 inequality can be used to govern the upper limit of the θ integration, in place of the $\pi/2$,

15 because $\tan \theta$ can always attain the value $1/\cos \phi$ on the interval $[0, \pi/2]$.

16 Putting this all together and developing we get

$$S = \frac{8 \cdot 6}{2} \int_0^{\pi/4} d\phi \int_0^{\text{atan}\left(\frac{1}{\cos \phi}\right)} \sin \theta \, d\theta \int_0^{\frac{L}{\cos \theta}} k \, dk \quad , \quad (\text{A.16})$$

$$= \frac{8 \cdot 6}{2 \cdot 2} L^2 \int_0^{\pi/4} d\phi \int_0^{\text{atan}\left(\frac{1}{\cos \phi}\right)} \frac{\sin \theta}{\cos^2 \theta} d\theta \quad , \quad (\text{A.17})$$

$$= 4L^2 \cdot 3 \int_0^{\pi/4} d\phi \left(\frac{1}{\cos \theta} \right) \Big|_0^{\text{atan}(\frac{1}{\cos \phi})} , \quad (\text{A.18})$$

$$= 4L^2 \cdot 3 \int_0^{\pi/4} d\phi \left(\frac{\sqrt{1 + \cos^2 \phi}}{\cos \phi} - 1 \right) , \quad (\text{A.19})$$

$$= 4L^2 \cdot 3 I \quad . \quad (\text{A.20})$$

1 To evaluate the last integral, I we use the substitution $\cos \phi = \sqrt{\cos \gamma}$. The limits for γ now

2 become $\phi = 0 \leftrightarrow \gamma = 0, \phi = \pi/4 \leftrightarrow \gamma = \pi/3$. We also have $\sin \phi =$

3 $\sqrt{1 - \cos \gamma}, \sqrt{1 + \cos^2 \phi} = \sqrt{1 + \cos \gamma}$. This leads to $d\phi = \frac{\sin \gamma d\gamma}{2\sqrt{\cos \gamma}\sqrt{1 - \cos \gamma}}$, which can be

4 simplified to $d\phi = \frac{\sqrt{1 + \cos \gamma} d\gamma}{2\sqrt{\cos \gamma}}$. So

$$I = -\frac{\pi}{4} + \int_0^{\pi/3} \frac{\sqrt{1 + \cos \gamma} d\gamma}{2\sqrt{\cos \gamma}} \left(\frac{\sqrt{1 + \cos \gamma}}{\sqrt{\cos \gamma}} \right) , \quad (\text{A.21})$$

$$= -\frac{\pi}{4} + \frac{1}{2} \int_0^{\pi/3} \frac{(1 + \cos \gamma) d\gamma}{\cos \gamma} , \quad (\text{A.22})$$

$$= -\frac{\pi}{4} + \frac{\pi}{6} + \frac{1}{2} \int_0^{\pi/3} \frac{d\gamma}{\cos \gamma} , \quad (\text{A.23})$$

$$= -\frac{\pi}{12} + \frac{1}{2} \int_{\pi/6}^{\pi/2} \frac{d\sigma}{\sin \sigma} , \quad (\text{A.24})$$

$$= -\frac{\pi}{12} + \frac{1}{2} (\ln \tan(\sigma/2)) \Big|_{\pi/6}^{\pi/2} , \quad (\text{A.25})$$

$$= -\frac{\pi}{12} - \frac{1}{2} \ln \tan(\pi/12) , \quad (\text{A.26})$$

$$= -\frac{\pi}{12} + \frac{1}{2} \ln(2 + \sqrt{3}) , \quad (\text{A.27})$$

$$= -\frac{\pi}{12} + \ln(1 + \sqrt{3}) - \ln \sqrt{2} , \quad (\text{A.28})$$

$$= 0.39667956 \quad . \quad (\text{A.29})$$

5 Finally, now

$$\frac{S}{4L^2} = 3I = 1.19 , \quad (\text{A.30})$$

- 1 This factor is the empirically observed excess area of an average plane embedded into a cube.
 2 This is the approximately 20 per cent increase in actively sampled points. It is a larger factor than
 3 the comparable 1.12 that would appear in a similar 2D problem.

4
 5
$$\frac{S}{2L} = \frac{4}{\pi} \log \cot \frac{\pi}{8} = \frac{4}{\pi} \ln(1 + \sqrt{2}) = 1.122 \quad . \quad (\text{A.31})$$

- 6 Another approach to evaluating (A.9) is to introduce an auxiliary variable via $\frac{1}{k} =$
 7 $\frac{2}{\sqrt{\pi}} \int_0^\infty d\alpha \exp -(\alpha^2(k_x^2 + k_y^2 + k_z^2))$. Then there are just a few steps to a single integral and a
 8 numerical evaluation: $\frac{S}{4L^2} = \frac{\pi}{4} \int_0^\infty d\alpha \left(\frac{\text{erf}(\alpha)}{\alpha}\right)^3 = 1.190038$, which is (A.30).

9 **A.3 Creating distributions according to some prescribed function Eq (3.15)**

- 10 In order to create a numerical sampling map for modulated side views, we would like to assign
 11 azimuthal angles to projections such that the oscillatory azimuthal distribution density indicated
 12 by (3.11) is achieved. This is well known how to do: for completeness, we include the argument
 13 here. From the density function (3.11), the cumulative distribution function can be found which
 14 is

15
$$\text{cdf}(\phi) = \int_0^\phi (1 + \lambda \cos 2 \phi) d\phi = \phi + \frac{\lambda}{2} \sin 2 \phi \quad . \quad (\text{A.32})$$

- 16 Now the azimuthal angle should be given by

17
$$\text{cdf}^{-1}(\phi \in [0, 2\pi]) \quad . \quad (\text{A.33})$$

- 18 That is, numbers should be drawn evenly between 0 and 2π , resulting in an array given by (A.33).
 19 These are the angle labels to be given to achieve the desired distribution (3.11). So long as $\lambda < 1$,
 20 this is easy to do, because the distribution is positive and the cdf is monotonically increasing

1 (graphically, the inverse corresponds to flipping across the diagonal, which maps a function into
2 another function due to monotonicity). The python pseudo code would read: phi0= cdf0=
3 np.linspace(0,2*np.pi,NumPoints); cdf = phi0 + $\frac{\lambda}{2} \sin(2 \text{ phi0})$, cdfInv = np.interp(phi0,cdf,
4 cdf0). That is, map the array phi0 to the desired phi (which is the desired cdfInv), in the same
5 manner that cdf was mapped to cdf0, where phi0, cdf0 are both regularly spaced.

6 **A.4 Proof that uniform distribution maximizes the SCF Eq (3.22).**

7
8 Consider a set of positive numbers $\{a_i\}$ that satisfy a constraint $C: \sum_i a_i = M$. The set are to
9 represent the sampling on the unit sphere. We wish to maximize $\sum_i \frac{1}{a_i}$ subject to C . We begin by
10 writing the usual variational:

$$11 \quad \mathcal{L} = \sum_i \frac{1}{a_i} + \mu (\sum_i a_i - M) \quad , \quad (\text{A.34})$$

12 where μ is a Lagrange parameter. Extremizing \mathcal{L} wrt the a_j yields

$$13 \quad \frac{\partial \mathcal{L}}{\partial a_j} = \mu - \frac{1}{a_j^2} = 0 \quad \rightarrow \quad a_j = \sqrt{\mu} \quad . \quad (\text{A.35})$$

14 The second variation is:

$$15 \quad \frac{\partial^2 \mathcal{L}}{\partial a_j^2} = 2 \frac{1}{a_j^3} = 2 \mu^{\frac{3}{2}} > 0 \quad . \quad (\text{A.36})$$

16 Since the second variation is positive, the uniform solution $a_j = \text{constant}$, corresponds to a
17 minimum.

18 The argument supplied here implies why the SCF attains its maximum (1/SCF attains its minimum
19 as in the above calculation), when the sampling (which is a conserved quantity on every shell of

1 Fourier space) is distributed as uniformly as possible, or equivalently the projections are
 2 distributed uniformly.

3

4 **A.5 Derivation of Eq. (3.23); SCF for modulated side-views**

5

6 From Eq. 3.16, we have

7

$$8 \quad 2k \text{ sp}(k, \theta, \phi) = \frac{2(1-\lambda \cos 2\phi)}{\pi \sin \theta} \quad (\text{A.37})$$

9

10

11 Using the definition of SCF, $1/\text{SCF} = \langle (1/(2k \text{ sp})) \rangle$, then (A.37) becomes

12

$$13 \quad \text{“side-modulated”} \quad \frac{1}{\text{SCF}} = \frac{1}{\pi} \int_0^{\pi/2} d\phi \int_0^{\pi} \sin \theta d\theta \frac{\pi \sin \theta}{2(1-\lambda \cos 2\phi)}, \quad (\text{A.38})$$

14

15

16 where the last term in the integrand of (A.38) is the reciprocal of (A.37). The integration over

17 θ , can be easily performed ($\int_0^{\pi} \sin^2 \theta d\theta = \pi/2$) leaving:

18

$$19 \quad \text{“side-modulated”} \quad \frac{1}{\text{SCF}} = \frac{\pi}{4} \int_0^{\pi/2} \frac{d\phi}{(1-\lambda \cos 2\phi)}, \quad (\text{A.39})$$

20

$$21 \quad = \frac{\pi^2}{8} \frac{1}{\pi} \int_0^{\pi} \frac{dv}{(1-\lambda \cos v)}, \quad (\text{A.40})$$

22

23 Integrals of the sort that appear in (A.40) are easily reduced by means of the so-called

24 Weierstrass half angle formula: $t = \tan v/2$; $\cos v = \frac{1-t^2}{1+t^2}$. The integral in (A.40) becomes

$$25 \quad \frac{2}{1-\lambda} \int_0^{\infty} \frac{dt}{1+\frac{1+\lambda}{1-\lambda}t^2} = \frac{2}{\sqrt{1-\lambda^2}} \int_0^{\infty} \frac{dw}{1+w^2} = \frac{2}{\sqrt{1-\lambda^2}} \frac{\pi}{2}. \quad \text{So the expression in (A.40) becomes:}$$

$$26 \quad \text{“side-modulated”} \quad \frac{1}{\text{SCF}} = \frac{\pi^2}{8} \frac{1}{\sqrt{1-\lambda^2}}. \quad (\text{A.41})$$

27

1 which is (3.23).

2

1 **Appendix B Derivation of Eq. (3.12) and (3.18): sampling distributions from**
 2 **projection distributions**

3
 4
 5 Eq (3.11) reads

$$6 \quad \text{sp}(\vec{k}) = \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} \delta(\hat{n} \cdot \vec{k}) / C_{N,\text{side}} \quad , \quad (\text{B.1})$$

7
 8
 9
 10
 11 where integrations are taken over all unit vectors, \hat{n} , in 3D. Also $C_{N,\text{side}}$ is a normalization
 12 constant ensuring Eq (3.8): $\langle 2k \text{sp}(\vec{k}) \rangle_{\hat{k}} = 1$, where $\langle \cdot \rangle_{\hat{k}}$ denotes angular average over the
 13 angles in k with the uniform measure on the sphere. The integration in B.1 is over the set of normal
 14 vectors to the sphere, with the given constraint. Putting this together with B.1 yields:

$$15 \quad C_{N,\text{side}} = C_{N,\text{side}} \langle 2k \text{sp}(\vec{k}) \rangle_{\hat{k}} \quad , \quad (\text{B.2})$$

$$16 \quad = 2k \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} \langle \delta(\hat{n} \cdot \vec{k}) \rangle_{\hat{k}} \quad , \quad (\text{B.3})$$

$$17 \quad = 2 \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} \langle \delta(\hat{n} \cdot \hat{k}) \rangle_{\hat{k}} \quad . \quad (\text{B.4})$$

18
 19 Now $\langle \delta(\hat{n} \cdot \hat{k}) \rangle_{\hat{k}}$ cannot be a function of the direction of \hat{n} . So it can be conveniently
 20 calculated when $\hat{n} = \hat{z}$, which does not depend on an azimuthal angle in the integration over \hat{n} ,
 21 and therefore leads only to the average over the altitude. This leads to:

$$22 \quad \langle \delta(\hat{n} \cdot \hat{k}) \rangle_{\hat{k}} = \frac{1}{2} \int_0^\pi \sin \theta \, d\theta \delta(\hat{k}_z) \quad , \quad (\text{B.5})$$

$$23 \quad = \frac{1}{2} \int_0^\pi \sin \theta \, d\theta \delta(\cos \theta) \quad , \quad (\text{B.6})$$

$$24 \quad = \frac{1}{2} \quad . \quad (\text{B.7})$$

25

1

2 Eq. B.7 is a natural result. It is the ratio of the circumference to the surface area of the unit circle:

3 $2\pi/4\pi = 1/2$. Returning to (B.4) we get:

$$4 \quad C_{N,\text{side}} = \int^{|\hat{n}\cdot\hat{z}| < \cos \alpha} d\hat{n} \quad , \quad (\text{B.8})$$

$$5 \quad = 2\pi \int_0^\pi \sin \theta_n \, d\theta_n \, \Theta(|\cos \theta_n| < \cos \alpha) \quad , \quad (\text{B.9})$$

$$6 \quad = 2\pi \int_\alpha^{\pi-\alpha} \sin \theta_n \, d\theta_n \quad , \quad (\text{B.10})$$

$$7 \quad = 4\pi \cos \alpha \quad . \quad (\text{B.11})$$

8

9 So, substituting (B.11) into (B.1) yields

10

$$11 \quad \text{sp}(\vec{k}) = \frac{1}{4\pi \cos \alpha} \int^{|\hat{n}\cdot\hat{z}| < \cos \alpha} d\hat{n} \, \delta(\hat{n} \cdot \vec{k}) \quad . \quad (\text{B.12})$$

12

13 It is easy to argue that $\text{sp}(\vec{k})$ does not depend on the azimuthal angle of \hat{k} , which we can

14 therefore take to be zero in order to evaluate (B.12): $\hat{k} = \sin \theta \, \hat{x} + \cos \theta \, \hat{z}$. Instead of the

15 integration over the sphere given by the unit vector, \hat{n} , we need to perform the integral in (B.12)

16 over the great circle perpendicular to \hat{k} . Therefore, we can parametrize \hat{n} , in the integration in

17 (B.12) by

$$18 \quad \hat{n} = (-\cos \theta \sin \beta \, , \cos \beta \, , \sin \theta \sin \beta) \quad . \quad (\text{B.13})$$

19 Eq. (B.13) is a parametrization of all the unit vectors perpendicular to \vec{k} as described in the last

20 paragraph. By changing β , we can sweep out the unit vector given by (B.13): these are the locus

21 of normals to \hat{k} and outside the cone of half angle α . So from (B.12)

22

$$k \text{ sp}(\vec{k}) = \frac{1}{4 \pi \cos \alpha} \int_0^{2\pi} d\beta \Theta(|\sin \theta \sin \beta| < \cos \alpha) \quad . \quad (\text{B.14})$$

The criterion $\Theta(|\sin \theta \sin \beta| < \cos \alpha)$ in B.14 is a rewrite for the constraint of the projection directions, $|\hat{n} \cdot \hat{z}| < \cos \alpha$, from Eq. B.12. Continuing from Eq. B.14.

$$k \text{ sp}(\vec{k}) = \frac{1}{\pi \cos \alpha} \int_0^{\pi/2} d\beta \Theta(\sin \beta < \cos \alpha / \sin \theta) \quad . \quad (\text{B.15})$$

If $\cos \alpha > \sin \theta$, then the argument of the indicator function in (B.15) is always true. If not the upper limit of β in the integral must be reduced to $\sin^{-1}(\cos \alpha / \sin \theta)$. This leads to:

$$k \text{ sp}(\vec{k}, \theta) = \frac{1}{\pi} \frac{\sin^{-1}(\frac{\cos \alpha}{\sin \theta})}{\cos \alpha} \quad , \quad \left| \frac{\pi}{2} - \theta \right| < \alpha, \quad (\text{B.16})$$

$$\text{"side-like"} \quad \frac{1}{2} \frac{1}{\cos \alpha} \quad , \quad \left| \frac{\pi}{2} - \theta \right| \geq \alpha \quad ,$$

which is (3.12).

Finally

$$\text{sp}(\vec{k}) = \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} \delta(\hat{n} \cdot \vec{k}) / C_{N,\text{top}} \quad . \quad (\text{B.17})$$

Using (B.8) and (B.17) using the parallel argument to (B.1)-(B.7) together, we note that

$$C_{N,\text{side}} + C_{N,\text{top}} = \int^{|\hat{n} \cdot \hat{z}| < \cos \alpha} d\hat{n} + \int^{|\hat{n} \cdot \hat{z}| > \cos \alpha} d\hat{n} = 4 \pi \quad . \quad (\text{B.18})$$

So

$$C_{N,\text{top}} = 4 \pi - C_{N,\text{side}} = 4 \pi (1 - \cos \alpha) = 8 \pi \sin^2 \alpha / 2 \quad . \quad (\text{B.19})$$

1
2
3
4
5
6
7
8

The parallel derivation to (B.14) now becomes:

$$k \text{ sp}(\vec{k}) = \frac{1}{8\pi \sin^2 \alpha/2} \int_0^{2\pi} d\beta \Theta(|\sin \theta \sin \beta| > \cos \alpha) \quad . \quad (\text{B.20})$$

9 This is the integration around the locus of points normal to \hat{k} and inside the cone of half-
10 angle, α . However, $\sin \beta$ may be replaced by $\cos \beta$ by shift of origin, and an overall factor of 4
11 introduced due to the 4 equivalent quadrants:

12

$$k \text{ sp}(\vec{k}) = \frac{1}{2\pi \sin^2 \alpha/2} \int_0^{\pi/2} d\beta \Theta(\cos \beta > \cos \alpha / \sin \theta) \quad . \quad (\text{B.21})$$

14 If $\cos \alpha > \sin \theta$, then the condition of the indicator function cannot be fulfilled, and the left-
15 hand side = 0. Otherwise

16

$$k \text{ sp}(\vec{k}) = \frac{1}{2\pi \sin^2 \alpha/2} \int_0^{\pi/2} d\beta \Theta(\beta < \arccos(\cos \alpha / \sin \theta)) \quad . \quad (\text{B.22})$$

18

19 So

$$k \text{ sp}(\vec{k}) = \frac{\arccos(\cos \alpha / \sin \theta)}{2\pi \sin^2 \alpha/2} \quad \text{for} \quad \left| \frac{\pi}{2} - \theta \right| \leq \alpha \quad , \quad (\text{B.23})$$

$$\text{"top-like"} = 0 \quad \text{for} \quad \left| \frac{\pi}{2} - \theta \right| > \alpha \quad . \quad (\text{B.24})$$

22

23

1 This is (3.18). Thus, the sampling is zero in directions close to along the z-axis, for the top like
2 cases.

3

4

5

1 References

- 2 1. Bartesaghi, A., et al., *Atomic Resolution Cryo-EM Structure of beta-Galactosidase.*
3 Structure, 2018. **26**(6): p. 848-856 e3.
- 4 2. Tan, Y.Z., et al., *Sub-2 Å Ewald curvature corrected structure of an AAV2 capsid variant.*
5 Nat Commun, 2018. **9**(1): p. 3628.
- 6 3. Zivanov, J., et al., *New tools for automated high-resolution cryo-EM structure*
7 *determination in RELION-3.* Elife, 2018. **7**.
- 8 4. Wyrick, J., et al., *Tomography of a Probe Potential Using Atomic Sensors on Graphene.*
9 ACS Nano, 2016. **10**(12): p. 10698-10705.
- 10 5. Fernandez-Leiro, R. and S.H. Scheres, *Unravelling biological macromolecules with cryo-*
11 *electron microscopy.* Nature, 2016. **537**(7620): p. 339-46.
- 12 6. Noble, A.J., et al., *Routine single particle CryoEM sample and grid characterization by*
13 *tomography.* Elife, 2018. **7**.
- 14 7. Noble, A.J., et al., *Reducing effects of particle adsorption to the air-water interface in*
15 *cryo-EM.* Nat Methods, 2018. **15**(10): p. 793-795.
- 16 8. D'Imprima, E., et al., *Protein denaturation at the air-water interface and how to prevent*
17 *it.* Elife, 2019. **8**.
- 18 9. Russo, C.J. and L.A. Passmore, *Progress towards an optimal specimen support for*
19 *electron cryomicroscopy.* Curr Opin Struct Biol, 2016. **37**: p. 81-9.
- 20 10. Grigorieff, N., *Three-dimensional structure of bovine NADH:ubiquinone oxidoreductase*
21 *(complex I) at 2.2 Å in ice.* J Mol Biol, 1998. **277**(5): p. 1033-46.
- 22 11. Penczek, P.A., *Three-dimensional spectral signal-to-noise ratio for a class of*
23 *reconstruction algorithms.* J Struct Biol, 2002. **138**(1-2): p. 34-46.
- 24 12. Diebold, C.A., et al., *Conical Fourier shell correlation applied to electron tomograms.* J
25 Struct Biol, 2015. **190**(2): p. 215-23.
- 26 13. Dudkina, N.V., et al., *Interaction of complexes I, III, and IV within the bovine respirasome*
27 *by single particle cryoelectron tomography.* Proc Natl Acad Sci U S A, 2011. **108**(37): p.
28 15196-200.
- 29 14. Dang, S., et al., *Cryo-EM structures of the TMEM16A calcium-activated chloride channel.*
30 Nature, 2017. **552**(7685): p. 426-429.
- 31 15. Tan, Y.Z., et al., *Addressing preferred specimen orientation in single-particle cryo-EM*
32 *through tilting.* Nat Methods, 2017. **14**(8): p. 793-796.
- 33 16. Naydenova, K. and C.J. Russo, *Measuring the effects of particle orientation to improve*
34 *the efficiency of electron cryomicroscopy.* Nat Commun, 2017. **8**(1): p. 629.
- 35 17. Lyumkis, D., *Challenges and opportunities in cryo-EM single-particle analysis.* J Biol
36 Chem, 2019. **294**(13): p. 5181-5197.
- 37 18. Penczek, P.A., *Resolution measures in molecular electron microscopy.* Methods Enzymol,
38 2010. **482**: p. 73-100.
- 39 19. Rosenthal, P.B. and R. Henderson, *Optimal determination of particle orientation,*
40 *absolute hand, and contrast loss in single-particle electron cryomicroscopy.* J Mol Biol,
41 2003. **333**(4): p. 721-45.
- 42 20. Sorzano, C.O., et al., *A review of resolution measures and related aspects in 3D Electron*
43 *Microscopy.* Prog Biophys Mol Biol, 2017. **124**: p. 1-30.

- 1 21. Jensen, G.J., *Alignment error envelopes for single particle analysis*. J Struct Biol, 2001.
2 **133**(2-3): p. 143-55.
- 3 22. Penczek, P.A., *Image restoration in cryo-electron microscopy*. Methods Enzymol, 2010.
4 **482**: p. 35-72.
- 5 23. Heymann, J.B., *Single-particle reconstruction statistics: a diagnostic tool in solving*
6 *biomolecular structures by cryo-EM*. Acta Crystallogr F Struct Biol Commun, 2019. **75**(Pt
7 1): p. 33-44.
- 8 24. Bracewell, R.N., *Strip Integration in Radio Astronomy*. Austrian Journal of Physics, 1956.
9 **9**: p. 198.
- 10 25. Stagg, S.M., et al., *ResLog plots as an empirical metric of the quality of cryo-EM*
11 *reconstructions*. J Struct Biol, 2014. **185**(3): p. 418-26.
- 12 26. Passos, D.O. and D. Lyumkis, *Single-particle cryoEM analysis at near-atomic resolution*
13 *from several thousand asymmetric subunits*. J Struct Biol, 2015. **192**(2): p. 235-44.
- 14 27. Saad, A., et al., *Fourier amplitude decay of electron cryomicroscopic images of single*
15 *particles and effects on structure determination*. J Struct Biol, 2001. **133**(1): p. 32-42.
- 16 28. Stagg, S.M., et al., *A test-bed for optimizing high-resolution single particle*
17 *reconstructions*. J Struct Biol, 2008. **163**(1): p. 29-39.
- 18 29. *Guinier Plot*. 2019.
- 19 30. Baldwin, P.R. and P.A. Penczek, *Estimating alignment errors in sets of 2-D images*. J
20 Struct Biol, 2005. **150**(2): p. 211-25.
- 21 31. Milotti, E., *1/f noise: a pedagogical review*. 2013.
- 22 32. Lyumkis, D., et al., *Likelihood-based classification of cryo-EM images using FREALIGN*. J
23 Struct Biol, 2013. **183**(3): p. 377-388.
- 24 33. Voss, N.R., et al., *A toolbox for ab initio 3-D reconstructions in single-particle electron*
25 *microscopy*. J Struct Biol, 2010. **169**(3): p. 389-98.
- 26 34. Zhang, C., et al., *Analysis of discrete local variability and structural covariance in*
27 *macromolecular assemblies using Cryo-EM and focused classification*. Ultramicroscopy,
28 2018.
- 29 35. Baxter, W.T., et al., *Determination of signal-to-noise ratios and spectral SNRs in cryo-EM*
30 *low-dose imaging of molecules*. J Struct Biol, 2009. **166**(2): p. 126-32.
- 31 36. Crowther, R.A., et al., *Three dimensional reconstructions of spherical viruses by fourier*
32 *synthesis from electron micrographs*. Nature, 1970. **226**(5244): p. 421-5.
- 33 37. Henderson, R., *The potential and limitations of neutrons, electrons and X-rays for atomic*
34 *resolution microscopy of unstained biological molecules*. Q Rev Biophys, 1995. **28**(2): p.
35 171-93.
36