

1 **Quantitative Mining of Compositional Heterogeneity in Cryo-**
2 **EM Datasets of Ribosome Assembly Intermediates**

3

4 Jessica N. Rabuck-Gibbons¹, Dmitry Lyumkis^{1,2}, James R. Williamson¹

5

6 1. Department of Integrative Structural and Computational Biology, Department of
7 Chemistry, and The Skaggs Institute for Chemical Biology,

8 The Scripps Research Institute, La Jolla, CA 92037, USA

9 2. Laboratory of Genetics and Helmsley Center for Genomic Medicine, The Salk

10 Institute for Biological Studies, La Jolla, CA 92037, USA

11

1 **Summary**

2
3 Macromolecular complexes are dynamic entities whose function is often intertwined with
4 their many structural configurations. Single particle cryo-electron microscopy (cryo-EM)
5 offers a unique opportunity to characterize macromolecular structural heterogeneity by
6 virtue of its ability to place distinct populations into different groups through
7 computational classification. However, current workflows are limited, and there is a
8 dearth of tools for surveying the heterogeneity landscape, quantitatively analyzing
9 heterogeneous particle populations after classification, deciding how many unique
10 classes are represented by the data, and accurately cross-comparing reconstructions.
11 Here, we develop a workflow that contains discovery and analysis modules to
12 quantitatively mine cryo-EM data for a set of structures with maximal diversity. This
13 workflow was applied to a dataset of *E. coli* 50S ribosome assembly intermediates,
14 which is characterized by significant structural heterogeneity. We identified new branch
15 points in the assembly process and characterized the interactions of an assembly factor
16 with immature intermediates. While the tools described here were developed for
17 ribosome assembly, they should be broadly applicable to the analysis of other
18 heterogeneous cryo-EM datasets.

19 20 **Keywords**

21
22 Cryo-electron microscopy; Single particle analysis; Ribosome biogenesis; Heterogeneity
23 analysis

1 Introduction

2

3 Cryo-electron microscopy (cryo-EM) is a rapidly evolving, powerful technology for
4 solving the structures of a wide variety of biological assemblies. The “resolution
5 revolution” in cryo-EM (Kühlbrandt, 2014), caused in part by advances in direct electron
6 detectors and improved data acquisition and analysis workflows, has led to high-
7 resolution structural insights into a wide variety of biological processes performed by
8 macromolecular assemblies (Fernandez-Leiro and Scheres, 2016). There have been
9 steady, but consistent improvements to achievable resolution, and the collective tools
10 are now enabling structure determination at true atomic resolution (Bartesaghi *et al.*,
11 2015; Tan *et al.*, 2018; Nakane *et al.*, 2020; Yip *et al.*, 2020; Zhang *et al.*, 2020). There
12 have also been numerous advances in workflows for analyzing structurally
13 heterogeneous particle populations, and data processing software now routinely include
14 strategies for handling distributions of structures that arise from compositional or
15 conformational changes in the macromolecular species of interest (Elmlund and
16 Elmlund, 2012; Gao *et al.*, 2004; Klaholz, 2015; Liao, Hashem and Frank, 2015; Nakane
17 *et al.*, 2018; Scheres, 2016; Spahn and Penczek, 2009; Wang *et al.*, 2013; White *et al.*,
18 2017; Zhong *et al.*, 2021; Grant, Rohou and Grigorieff, 2018; Lyumkis *et al.*, 2013;
19 Punjani and Fleet, 2021b; Punjani and Fleet, 2021a). However, most current cryo-EM
20 workflows still focus on achieving the maximum possible resolution, which requires
21 selecting and averaging potentially heterogeneous subsets of the data in the interest of
22 increasing the particle count for the homogeneous regions of a map. This strategy
23 comes at the expense of either eliminating particle populations that do not conform to
24 the predominant species or neglecting dynamic and labile regions of reconstructed
25 maps, which are often of biological interest.

26

27 Another challenge in cryo-EM heterogeneity analysis is that there is no way to define
28 the number of distinct structures in a given dataset *a priori*. It is up to the researcher to
29 employ a classification strategy and to heuristically determine the number of distinct
30 classes. Furthermore, there is no set procedure to determine the threshold for
31 examining map features and differences between maps. Thresholds are often set in a
32 subjective manner in order to best display the features of interest in the maps, although
33 an approach was recently described where a voxel-based false discovery rate could be
34 determined to establish a noise threshold for contouring (Beckers, Jakobi and Sachse,
35 2019). Thus, determining the final number of classes in a dataset and quantitatively
36 comparing a set of maps in order to tell a concise biological story with statistical
37 significance remains a challenge.

38

39 The process of ribosome assembly provides a useful case study for mining and
40 quantitatively assessing structural heterogeneity in cryo-EM data. The bacterial 70S

1 ribosome is a complex macromolecular machine composed of three ribosomal RNAs
2 (rRNAs) and ~50 ribosomal proteins (r-proteins) that form a large 50S subunit and a
3 small 30S subunit. Ribosome assembly occurs within several minutes *in vivo*, and the
4 process includes transcription and translation of the rRNAs and r-proteins, folding of the
5 rRNA and r-proteins, and docking of the r-proteins on the rRNA scaffold. rRNA folding
6 events and proper r-protein binding are facilitated by ~100 ribosome assembly factors.
7 Given the efficiency and speed of the assembly process, structural intermediates are
8 difficult to isolate and purify. However, perturbations in ribosome assembly lead to the
9 accumulation of numerous structural intermediates, which collectively inform molecular
10 mechanisms of ribosome assembly (Shajani, Sykes and Williamson, 2011; Stokes *et*
11 *al.*, 2014; Sashital *et al.*, 2014; Sykes *et al.*, 2010; Jomaa *et al.*, 2014; Li *et al.*, 2013; Ni
12 *et al.*, 2016; Davis *et al.*, 2016; Rabuck-Gibbons *et al.*, 2020). The major parts of the
13 ribosome that are often present or missing in assembly intermediates are the central
14 protuberance (CP), the L7/12 and L1 stalks, and the base (Figure 1A).

15
16 We previously developed a genetic approach by which the amount of a given r-protein,
17 in our case bL17, could be titrated by the addition of the small molecule homoserine
18 lactone (HSL) (Davis *et al.*, 2016). Limiting the amount of bL17 induced a roadblock in
19 ribosome assembly, causing intermediates to accumulate. In the first work using the
20 bL17-lim strain (Davis *et al.*, 2016), we identified thirteen distinct structures that fell into
21 four main structural classes (Figure 1B). Here, we will continue to use the nomenclature
22 for the main classes used by Davis and Tan, et al. These categories, ordered least to
23 most mature, are the B class which is missing the base, CP, and both stalks, the C
24 class in which the base is formed, but the central protuberance (CP) is either misdocked
25 or altogether missing, the D class in which the base of the 50S ribosome is missing, and
26 the E class, which contains both the base and the CP, but has variability in the
27 presence or absence of the stalks. Some of the “missing” regions (primarily rRNA, but
28 they may also include r-proteins) described above are not present in the reconstructed
29 maps but are in fact present in the sample and within individual particle images,
30 meaning that they contribute to “biological noise”. This becomes relevant for some of
31 the decisions that need to be made in the data analysis workflow, as will be discussed
32 below. In previous work, the four main initial classes belonging to the 50S assembly
33 intermediates (B, C, D, E) were each further subdivided by an additional round of
34 subclassification, resulting in thirteen distinct structures. While several different
35 subclassification schemes were attempted at that time using heuristics to determine the
36 number of subclasses, no attempt was made to establish quantitative criteria by which
37 the subclassification or coverage of relevant classes would be complete. While classes
38 were identified belonging to the 30S and 70S (F class and A class in Davis *et al.*, 2016),
39 they are not explicitly described in our previous work or in the work described here.

40

1 As our goal is to define broad trends in ribosome assembly through various
2 perturbations, it is important to quantitatively assess differences between intermediates
3 that accumulate under various specific conditions and to organize them into a ribosome
4 assembly landscape (Davis *et al.*, 2016; Bernstein *et al.*, 2004; Harnpicharnchai *et al.*,
5 2001; Jomaa *et al.*, 2011; Loerke, Giesebrecht and Spahn, 2010; Nikolay *et al.*, 2018;
6 Razi, Guarné and Ortega, 2017; Uicker, Schaefer and Britton, 2006). To this end, we
7 developed a data processing framework to analyze cryo-EM datasets methodically and
8 quantitatively in order to assess the number of distinct structures, the significant
9 differences among them, and to place these structures into a biological context. When
10 we apply our complete workflow to a dataset of ribosome assembly intermediates from
11 bL17-lim, we discover a total of forty-one different structures that are identifiable based
12 on a defined set of cutoff parameters. These structures include several novel
13 intermediates, such as the most immature assembly intermediate observed to date, and
14 an independent pathway contingent on the binding of a ribosome assembly factor, as
15 well as late-stage assembly intermediates. Together, these are organized into a revised
16 assembly landscape for the 50S ribosomal subunit under bL17-lim conditions.

17

18 **Results and Discussion**

19

20 ***An overview of the heterogeneity processing workflow***

21 There are two main phases in the framework for systematic analysis of heterogeneous
22 ensembles of macromolecular conformations (Figure 2). The first phase is a discovery
23 phase, which begins with iterative rounds of hierarchical classification and sub-
24 classification using a defined set of thresholding parameters. The goal of this first phase
25 to uncover the broad spectrum of distinct classes in a cryo-EM dataset, starting with
26 traditional pre-processing and data cleaning steps (e.g. motion correction, particle
27 picking, CTF estimation, and initial 2D and 3D classification). The initial data cleaning
28 steps defined here are intended to be very lenient, such that the only particles removed
29 from the dataset are clear artifacts or molecular species that are not of interest. For
30 example, in the case of 50S ribosome assembly intermediate analysis, we remove
31 particles that are obvious 30S or 70S ribosomes and proteasomes from the stack, but
32 we do not remove any classes that could possibly be 50S assembly intermediates.
33 After the cleaning steps, an iterative subclassification strategy is used to parse out
34 molecular heterogeneity. After an initial round of classification, each class (class X) is
35 subjected to a $n=2$ subclassification, resulting in two potential subclasses, X1 and X2.
36 Both subclasses are processed and binarized, and then difference maps X1-X2 and X2-
37 X1 are calculated, to determine if there is more heterogeneity that can be mined from
38 each class X. If the difference volumes don't reach a chosen molecular weight or
39 resolution threshold, then the subclassification is rejected, and further subclassification

1 is terminated. If neither of these two criteria are reached, the binary subclassification
2 process is iteratively repeated until one of the convergence criteria are met.

3
4 The second phase is an analysis phase, which is intended to quantitatively define and
5 distinguish structural features between maps, and further, to establish the number of
6 structural states using a given set of quantitative cutoffs. During this hierarchical
7 difference analysis, the full matrix of difference maps is calculated, and the molecular
8 weights of the difference maps are used as a metric to cluster the classes, which can be
9 visualized as a particle dendrogram. A line can be drawn through the dendrogram at a
10 chosen molecular weight threshold, which identifies similar maps that can be combined.
11 Next, to qualitatively differentiate between classes, the resulting set of maps are
12 compared to a catalog of coarse-grained structural features that are calculated from a
13 reference structure, in this case the bacterial 50S ribosome. It is convenient to use
14 features such as rRNA helices and r-proteins, that may be present or absent in various
15 classes. The presence of these coarse-grained reference features is quantitatively
16 analyzed using hierarchical clustering to organize and visualize the patterns of variation
17 among the final set of particle classes. For our dataset of bacterial ribosome assembly
18 intermediates, these features are used to place the observed classes into a putative
19 assembly pathway, based on a principle of parsimonious folding and unfolding.

20
21 ***A divisive resolution-limited subclassification approach facilitates identifying***
22 ***novel species***

23 A major challenge in the analysis of heterogeneous datasets is the accurate
24 identification of a broad diversity of structural states. To address this, we developed a
25 classification strategy to mine an experimental cryo-EM dataset for distinct particle
26 populations. Classification and refinement of particle classes can be undertaken using
27 a variety of software packages, and we have adopted the latest version of FrealignX,
28 whose code base is also implemented within *cisTEM* (Grant, Rohou and Grigorieff,
29 2018; Lyumkis *et al.*, 2013). We note that most processing packages that are capable of
30 classifying single-particle cryo-EM data can be employed for this purpose (Scheres,
31 2016; Nakane *et al.*, 2018; Zhong *et al.*, 2021; Punjani and Fleet, 2021b; Punjani and
32 Fleet, 2021a).

33
34 In typical cryo-EM workflows, 3D classification is performed several times, with different
35 choices for the total number of classes (n). If n is too small, the resulting classes may
36 have averaged properties leading to loss of structural diversity but potentially higher
37 resolution in the homogeneous regions. If n is too large, the data is subdivided into
38 nearly identical classes, but each class is characterized by lower resolution, because
39 the particle count contributing to the class decreases. For the characterization of
40 intrinsically heterogeneous datasets such as those encountered during ribosome

1 assembly, the goal of 3D classification is to capture the full range of structural diversity,
2 as opposed to a select few well-resolved classes. Therefore, we developed an iterative
3 subclassification strategy to systematically mine the data and identify distinct structural
4 intermediates, including species that are rare and underpopulated.

5
6 With the knowledge that our test dataset harbored at least thirteen intermediates (Davis
7 *et al.*, 2016), we started with $n=10$ in order to evaluate parameters for subclassification.
8 The ten initial classes are shown in Figure 3A. While we expected that we would find
9 the previous B, C, D and E classes in the dataset, the B-class was not present, and
10 rather, multiple classes that are subtle variations of the E-class were present. This
11 exemplifies one of the pitfalls of classification that we term “hiding”, where subclasses
12 can be mixed, only to emerge at subsequent stages of subclassification. A survey of
13 various classification parameters within FrealignX revealed that lowering the
14 *res_high_class* parameter, which is the resolution of the data to be used for
15 classification, ameliorated class hiding and had a strong effect on the classes that
16 emerged. This parameter is typically set to just below the estimated resolution limit of
17 the data. However, by setting *res_high_class* to 20Å, the gross class heterogeneity
18 increased, and the expected B-class emerged (Figure 3B). The resolution threshold for
19 classification is frequently defaulted and determined automatically during classification,
20 but it may also be explicitly set by the user or limited to the resolution of the first Thon
21 ring (Scheres, 2012; Scheres, 2016; Scheres *et al.*, 2008). With the well-defined
22 ribosome assembly case study, we show that a lower resolution threshold during
23 classification helps to identify particle subsets that are substantially distinct from the
24 predominant species.

25
26 We also examined different iterative subclassification strategies, with various numbers
27 of classes used for each stage of subclassification. In order to test the success of these
28 strategies, we selected a final n of ~ 30 , which was chosen because it provided a
29 convenient number to evaluate a variety of subclassification schemes, and because it
30 was close to twice the final number of classes found in the original bL17-lim dataset
31 (Davis *et al.*, 2016). The five classification schemes (Figure 3C) tested were: (1) a
32 simple 1-round classification with $n=30$, (2) a 2-round hierarchical classification of 6
33 initial classes, each subdivided into 5 ($n_1=6, n_2=5$; total $n=30$), (3) a 3-round hierarchical
34 subclassification of 2 initial classes each subdivided into 3, with a second subdivision
35 into 5 ($n_1=2, n_2=3, n_3=5$; total $n=30$), (4) a 3-round hierarchical classification of 5 initial
36 classes subdivided into 3, then subdivided into 2 ($n_1=5, n_2=3, n_3=2$; total $n=30$), and (5) a
37 5-round hierarchical binary subclassification strategy, where 2 initial classes were
38 subdivided into 2 until $n=32$ was reached ($n_1=2, n_2=2, n_3=2, n_4=2, n_5=2$; total $n=32$).

39

1 With the sole exception of the simple single-round classification with $n=30$, all of these
2 divisive schemes yielded new classes not previously identified (Supplemental Figure 1,
3 indicated by *). Furthermore, the iterative divisive approaches produced the greatest
4 range of structural diversity and avoided grouping together dissimilar classes. This
5 observation is perhaps not unexpected, as it is well known that a divisive classification
6 approach avoids local minima within the search space and is more robust than
7 attempting to produce a final number of classes directly (Gray, 1984; Sorzano *et al.*,
8 2010). Qualitatively, a first round of classification where n_1 is on the order of the number
9 of major classes works well, followed by smaller subdivisions. As an example, a three
10 round subclassification scheme with $[n_1 = 5, n_2 = 3, n_3 = 2]$, for a total of 30 final classes,
11 identified the greatest number of new structures, as shown in Supplemental Figure 1.
12 For this reason, we proceeded with the $n_1 = 5, n_2 = 3, n_3 = 2$ approach for our work,
13 although we note that the optimal classification scheme will likely vary with the distinct
14 heterogeneity spectrum for each unique dataset. Given that the observed classes are
15 relatively independent of the details of the subclassification, we turned our attention to
16 the criteria for termination of subclassification.

17

18 ***Defining an endpoint for subclassification***

19 The determination of when subclassification is complete is a key question in cryo-EM
20 analysis. Many times, classification is considered finished if a specific region of interest
21 can be resolved to a satisfactory resolution, depending on what question(s) the user
22 wishes to address. However, this subjective approach may be insufficient for the
23 purpose of uncovering hidden features and discovering new structural states, especially
24 if there are multiple datasets to be compared. To guide the analysis of our bL17-lim
25 dataset, and to establish a protocol that can be used to analyze other data with
26 statistical significance, our goal was to establish metrics by which we could confidently
27 terminate the subclassification. We adopted a simple metric to determine the endpoint
28 of subclassification. For any given class at any stage of subclassification, a test
29 subclassification is performed with $n=2$. If the two resulting subclasses differ by less
30 than a chosen noise threshold, or by less than a chosen molecular weight threshold,
31 then subclassification is complete, and the subdivision is rejected. Conversely, if the
32 thresholds are exceeded, the subclassification is retained, and the two resulting classes
33 are iteratively subjected to additional subclassification until the termination thresholds
34 are met (Figure 2).

35

36 There are at least two types of noise that need to be considered in the difference
37 analysis that are used to conclude subclassification. First, there is the intrinsic noise
38 floor in the map that arises from averaging noisy image data during the reconstruction
39 process. Second, there is biological noise, which can be broadly attributed to
40 conformational and compositional heterogeneity, resulting in density above the intrinsic

1 noise floor that cannot be interpreted in terms of a structure or slight shifts of well-
2 defined elements that may or may not be significant (Supplemental Figure 2). For
3 example, in the case of ribosome assembly, there are portions of rRNA that are present
4 in the sample, but do not resolve to a reasonable structure (Davis *et al.*, 2016). To
5 characterize a diverse set of classes, the goal is to identify significant differences that
6 exceed chosen thresholds for these noise components.

7
8 A three-step process was developed to remedy the above challenges, based on the
9 estimation of the real space noise in a given map. First, a low-pass filter is used to
10 reduce high-frequency information in the map (low-pass filter threshold, Table 1)
11 Clearly, this is inadvisable if high resolution is the goal for the experiment, but for
12 heterogeneity analysis, resolution is secondary to differentiating between broader
13 conformational and compositional differences. Second, it is important that the soft
14 spherical mask typically applied during classification is removed, and the standard
15 deviation of the unmasked map (σ_{map}) is calculated using standard cryo-EM analysis
16 programs. While the signal from the macromolecular object is included in this
17 calculation, that contribution to the standard deviation is negligible if the box size is
18 sufficiently large, so that voxels containing true signal represents 1-2% of the total map
19 volume. Effectively, σ_{map} provides a crude estimate of the intrinsic map noise. There
20 are several other ways to calculate a noise threshold, most recently the program
21 developed by Beckers *et al.* (Beckers, Jakobi and Sachse, 2019) which uses a false
22 discovery rate (FDR) to determine the threshold used for visualization and analysis, or
23 one can use the noise sampled from the periphery of the map. The values of $3\sigma_{\text{map}}$ are
24 highly correlated to the contour levels based on FDR as shown in Figure S3 but the
25 $3\sigma_{\text{map}}$ threshold generally exceeds the FDR threshold, and is thus more conservative.
26 Due to the prevalence of unresolved features in the ribosome data, we have used $3\sigma_{\text{map}}$
27 as a convenient threshold to eliminate noise. Third, each map is then binarized using a
28 $3\sigma_{\text{map}}$ threshold such that intensities greater than $3\sigma_{\text{map}}$ were set to 1, and intensities
29 less than $3\sigma_{\text{map}}$ were set to 0 (binarization threshold, Table 1). Other thresholds could
30 be devised and implemented, as long as they are applied consistently across classes.
31 These thresholded, binarized maps are used for the remainder of the analysis. Using
32 these maps is advantageous because the “noise” from flexible regions is removed from
33 the map, and there is a clear boundary of which parts of a structure are analyzed.
34 Further, binarization facilitates coarse-grained analysis and eliminates the need for
35 scaling.

36
37 To define the endpoint to classification, the above filtering and binarization steps are
38 applied after a test $n=2$ subclassification of a given class X into class X1 and X2. If
39 either class X1 or class X2 do not have a resolvable map, as defined by the resolution
40 limit (r-limit, Table 1), then the classification process is terminated. If the differences

1 between class X1 and class X2 are less than the volume limit (v-limit, Table 1), the
2 subclassification is terminated. However, if the differences between class X1 and X2 are
3 greater than the v-limit, then the subclassification is retained, and classes X1 and X2
4 are in turn further subdivided into 2 classes. This process then repeats on all classes
5 until either the r-limit or the v-limit are reached. This set of limits provides a consistent
6 and quantitative basis for iterative subclassification.

7

8 ***Segmented difference analysis between map identifies the exact number of*** 9 ***structural states***

10 Having discovered the structural variants in the data, we then asked how the different
11 maps compare to one another and where/what are the major differences. To address
12 this question, we developed a strategy to quantitatively assess similarities between the
13 classes. While the classification approach in the discovery phase is designed to
14 terminate once the structural features were no longer distinguishable using the r-limit or
15 v-limit, this procedure does not guarantee that individual structures within the collective
16 set of reconstructions are all distinct from one another. More specifically, a situation can
17 arise where two similar classes emerge (from hiding) in different branches of the
18 subclassification tree.

19

20 In the first step, difference maps are calculated between all of the binarized, thresholded
21 maps. Such difference maps are useful to identify regions of density that are distinct
22 between classes, and in our case, provide both qualitative and quantitative insight into
23 structural relationships between distinct assembly intermediates. Two specific examples
24 for distinct “D-classes” are shown in Figure 4A-B. The first two columns display two
25 distinct maps arising from some point during classification. The raw difference maps are
26 shown in the third column (map1-map2, red; map2-map1, blue). The approximate
27 molecular weight of these differences is also indicated. These difference maps are then
28 segmented to remove “dust” that may arise from minor conformational or compositional
29 variations between maps. This dust cannot be interpreted in biological terms at the
30 target resolution but may add up to a significant molecular weight (Table 1
31 segmentation threshold, Figure 4). Such difference maps can be computed for all
32 pairwise combinations of reconstructions arising from the classification procedure.

33

34 The pairwise difference maps are useful for both qualitative and quantitative
35 downstream analyses. To parse through structural differences, define an accurate final
36 number of *unique* structural variants in the data, and combine particles contributing to
37 similar maps, we employed a simple hierarchical clustering approach based on the
38 positive/negative molecular weight differences between structures. Based on the
39 clustering, it is possible to pare down the maps and combine particles from similar
40 reconstructions, even if they arise from different starting points in the classification

1 (Figure 5A). At this stage, two classes can be combined if the molecular weight
2 differences between the two classes are less than a given threshold. Since the
3 branchpoints of the dendrogram provide a measure of *molecular weight* differences
4 between maps, they can serve as a guide for analyzing the similarity between classes
5 overall based on the nodes of the dendrogram (Figure 5B). In the example in Figure 5B,
6 the dendrogram reveals that the leftmost structure is distinct from the other two and
7 needs to be treated independently, whereas the latter two can be combined into a single
8 class. Thus, although there are 42 distinct structures in Figure 5A, after hierarchical
9 clustering analysis and the subsequent merging of similar maps, there are 41 distinct
10 structures that will go forward in the analysis pathway. Collectively, these procedures
11 enable us to identify the exact number of structural states within the data, given the
12 limitations associated with identifying novel classes in the discovery phase and
13 according to the established criteria in the analysis phase, defined above.

14

15 ***Defining relationships between distinct structures***

16 An important step in analyzing differences between classes discovered within the above
17 procedures for heterogenous cryo-EM data analysis is to define *where* differences
18 between two maps are located. If a model (e.g. an atomic model or a cryo-EM structure)
19 exists as a reference, and if the reconstructed maps differ primarily by compositional
20 variation, then it is straightforward to use the model for interpreting the collective set of
21 maps (Davis *et al.*, 2016) in an “occupancy analysis.” In the case of bacterial ribosome
22 assembly, we have a well-defined reference model (Figure 6A). This reference structure
23 is broken into its individual r-RNA and r-protein parts, yielding theoretical cryo-EM
24 densities for each component (Figure 6B). Such individual densities can then be directly
25 compared to densities arising from experimental cryo-EM classification. It is important
26 that the reference densities are generated at (approximately) the same resolution as the
27 experimental densities arising from hierarchical clustering and difference analyses. The
28 theoretical maps are then binarized, which enables comparing the theoretical maps to
29 the binarized experimental maps arising from subclassification. Each binarized class
30 (Figure 6C) is then compared to each theoretical feature map by counting overlapping
31 voxels and normalizing to the theoretical volume, to define the fractional occupancy of
32 the selected feature in the map that can be completely present (Figure 6D), partially
33 present due to partial flexibility or a misdocked figure (Figure 6E), or completely missing
34 (Figure 6F). The complete set of fractional occupancies are given as an n by m matrix of
35 values between 0 and 1, where n describes the set of classes and m defines the
36 number of features.

37

38 The resulting fractional occupancies can be visualized as a heat map and subjected to
39 hierarchical clustering to organize the classes and features (Figure 6G). Clustering
40 along the feature (x-axis) groups elements (in this case, r-proteins and rRNAs), and

1 clustering along the map (y-axis) groups the maps according to their occupancy. As
2 expected, the B, C, D, and E, maps cluster well together. The occupancy matrix
3 facilitates the visualization of large blocks of structural features that co-vary across the
4 particle classes, providing cooperative folding blocks (Figure 6H) (Davis *et al.*, 2016).
5 This procedure enables a quantitative comparison of distinct sets of maps that differ by
6 compositional variants. We note that this procedure is not currently compatible with
7 conformational variability or density that is not represented in the reference. However, if
8 there are multiple reference models that differ by discrete conformational changes, the
9 current protocol can be extended to competitively compare occupancies against
10 different reference models.

11

12 ***Ordering structures in a ribosome assembly pathway***

13 In the final step that is relevant to defining an assembly process, we developed a
14 module that uses molecular weight differences to place ribosome assembly
15 intermediates into a pathway. In this analysis, a “folding” matrix is calculated from the
16 molecular weight difference that would need to be added to a given map to create a
17 second map, and the “unfolding” matrix is calculated from the molecular weight that
18 would need to be subtracted from one map to create a second. Each element of the
19 folding/unfolding matrix can be considered as the driving force/barrier for a structural
20 transition between two classes. By postulating that folding proceeds by incremental
21 assembly, with minimal unfolding, a parsimonious transition graph can be constructed
22 with allowed passages between classes based on simple criteria – there is a molecular
23 weight cutoff unfolding transitions, and there is a limit set to the number of transitions
24 emanating from each class. Large unfolding events are unlikely, given the large
25 number of states that are close in molecular weight, but small unfolding events must be
26 permitted to allow for structural rearrangements required to transition between classes.
27 Finally, it is likely that structural transitions proceed from a finite manifold of close
28 intermediates. The folding and unfolding matrices can be used to construct a directed
29 graph of allowed transitions using these criteria, as shown in Figure 7.

30

31 ***Analysis of bL17-lim data using the quantitative heterogeneity mining protocol***

32 Our quantitative mining protocol was developed using data collected from newly purified
33 assembly intermediates from the previously characterized bL17-limitation strain (Davis
34 *et al.*, 2016). We collected new cryo-EM data (Supplementary Table 1 and subjected it
35 to our workflow. In the discovery phase, we employed an updated high-resolution limit
36 for refinement ($res_high_class=20\text{\AA}$) and an initial $n5>n3>n2$ hierarchical classification
37 scheme, followed by additional rounds of binary subdivision. All maps were binarized
38 according to the $3\sigma_{map}$ threshold determined individually for each map. To determine if
39 subclassification was complete, we selected a v-limit of 1.5 kDa. The rationale for this
40 choice is that 1.5 kDa represents the size of the smallest RNA helix present in the

1 bacterial ribosome and therefore corresponds to the smallest feature that we would like
2 to capture in the data. For our purposes, smaller features can be assumed to be either
3 biological and/or experimental noise. After iterative subclassification, the total number of
4 classes is 42. The similarity between all of the maps was then analyzed by the
5 hierarchical clustering analysis as described above, and at this stage, pairs of classes
6 were combined with a 10 kDa difference threshold (Figure 5A, dotted red line). This
7 cutoff was chosen because it is close to the average molecular weight of all proteins
8 and rRNA features, and we wished to reduce the complexity of our data. We found one
9 pair of structures that were similar to one another according to our established criteria
10 for biological significance, and the particles belonging to these classes were accordingly
11 combined (Figure 5A). Thus, using this protocol, a total of forty-one ribosome assembly
12 intermediates were identified using quantitative metrics and similarity analysis, with
13 minimal heuristic intervention.

14
15 The classes were then subjected to an occupancy analysis to view the sets of
16 cooperative folding blocks across the different classes. For the bL17-lim dataset here,
17 the maps are compared to the reference crystal structure (PDB 4ybb) (Figure 6A). The
18 reference 4ybb structure is filtered to 10Å and segmented into volumes corresponding
19 to individual r-proteins and rRNA helices, resulting in 139 theoretical map segments
20 (Figure 6B). The fractional analysis revealed five major structural blocks (Figure 6G,H)
21 The largest, block I (red) is composed of structural elements that are largely present in
22 all of the classes. These elements are found on the back of the ribosome and represent
23 the structural core that can form without bL17. Block II (green) represents the central
24 protuberance, which is fully formed in the D and E classes but is either missing or
25 misdocked in the B and C classes. Block III (yellow) maps to the base of the ribosome
26 and the L1 stalk. These features are mostly present in the C and E classes but are
27 missing in the B and D classes. These two blocks represent parallel pathways in
28 assembly (Davis *et al.*, 2016), as it is unlikely that the base of the ribosome would be
29 unfolded or disordered in order to form the central protuberance, and vice versa. Block
30 IV (blue) represents density that is specific to the base of the L7/12 stalk and is mostly
31 present in the D and some of the E classes. Finally, block V (purple) represents density
32 that is mostly missing in all maps, and is composed of h68, bL9, the L7/12 stalk, and the
33 top of the L1 stalk. These represent features that are among the last of the ribosome to
34 fold (like h68) or are flexible elements (the stalks). bL9 is a special case, as the
35 conformation in the crystal structure is an artifact due to crystallization; in cryo-EM
36 structures, bL9 wraps around to the interface between the 30S and 50S subunits and is
37 often flexible. These central blocks are very similar to the ones that we discovered
38 previously (Davis *et al.*, 2016), but this updated occupancy matrix will allow us to
39 compare the blocks that arise from other depletion or deletion strains in order to explore
40 the cooperative block-like behavior or ribosome assembly in future work.

1
2 The ordering module was used to calculate an initial pathway in the absence of bL17-
3 lim, which was modified by hand, as elements like the misdocked central protuberance
4 and non-native structural elements can have large effects on molecular weight
5 differences but may arise earlier in the order of assembly. We found the same initial
6 super classes as previously reported (B, C, D and E classes). While the classes we
7 found were similar to the initial bL17 data (Supplemental Figure 4), the new classes
8 enabled refinement of our bL17-lim ribosome assembly pathway. First, we found a
9 YjgA-dependent pathway through the assembly process (Figure 7, classes denoted by
10 *). YjgA was only bound if the central protuberance and the L1 stalks were present. We
11 also discovered three potential parallel processes in the C class where the earliest
12 event could either be the completion of the L1 stalk, the partial docking of the central
13 protuberance, or the formation of the base of the L7/12 stalk. We did not previously
14 observe the formation of the base of the L7/12 structure in the assembly pathway for
15 any class. We also found an immature B class (Figure 7, structure 1) and an immature
16 D class where the base was missing, but the L7/12 stalk was absent or present (Figure
17 7, structures 2 and 3), which were not present in the original set of 13 structures (Davis
18 *et al.*, 2016). In particular, the immature B class represents the least mature pre-50S
19 intermediate identified to date. We also identified several structures that seem to be
20 transition points between the two classes (Figure 7, structures 6 and 7), and we observe
21 formation of density at the base of the structures, which is lacking in other D classes
22 and is present in other E classes. These new discoveries inform a better understanding
23 of ribosome assembly in the context of bL17 limitation, and the data analysis process
24 will allow us to quantitatively assess cryo-EM data from other limitation strains and
25 ribosome assembly defects.

26

27 **Conclusions**

28

29 Heterogeneity analysis in cryo-EM provides exciting opportunities to discover new
30 biology, but current workflows suffer from numerous challenges. The work here
31 addresses three challenges that researchers face in the analysis of cryo-EM data, as
32 exemplified using a case study of ribosome assembly intermediates: establishing a
33 divisive approach to classification with well-defined endpoints to discover novel states, a
34 comprehensive difference analysis between distinct structures, and the application of
35 well-defined criteria (thresholds) for limiting classification. The application of specific
36 thresholds and limits (Table 1) has been critical to the success of analyzing ribosome
37 assembly intermediate data. The implementation of this workflow has allowed us to
38 identify an additional 28 ribosome assembly intermediates (counting the 41 assembly
39 intermediates after merging similar classes), which include an independent pathway for
40 the assembly factor YjgA and the earliest intermediate discovered to date in the

1 ribosome assembly process. The discovery and analysis modules of this workflow
2 provide a powerful analysis for quantitatively interrogating heterogeneous cryo-EM data
3 for complex biological processes.

4

5 **Acknowledgements**

6

7 Molecular graphics and analyses were performed with the USCF Chimera package
8 (supported by NIH P41 GM103311). This work was supported by grants from the NIH
9 DP5-OD021396 and U54 AI150472 (to D.L.) R35-GM136412 (to JRW).

10

11 **Author Contributions**

12 Jessica N. Rabuck-Gibbons: Conceptualization, Investigation, Methodology, Software,
13 Formal Analysis, Data Curation, Writing – Original Draft, Writing – Review & Editing,
14 Visualization.

15

16 Dmitry Lyumkis: Conceptualization, Investigation, Writing – Review & Editing,
17 Resources.

18

19 James R. Williamson: Conceptualization, Methodology, Software, Data Curation,
20 Writing – Review & Editing, Visualization, Supervision, Project Administration, Funding
21 Acquisition.

22

23 **Declaration of Interests**

24 The authors declare no competing interests.

25

26 **References**

27

28 Baldwin, P. R. and Lyumkis, D. (2020) 'Non-uniformity of projection distributions
29 attenuates resolution in Cryo-EM', *Progress in biophysics and molecular biology*, 150,
30 pp. 160-183.

31 Baldwin, P. R. and Lyumkis, D. (2021) 'Tools for visualizing and analyzing Fourier space
32 sampling in Cryo-EM', *Progress in Biophysics and Molecular Biology*, 160, pp. 53-65.

33 Bartesaghi, A., Merk, A., Banerjee, S., Matthies, D., Wu, X., Milne, J. L. and
34 Subramaniam, S. (2015) '2.2 Å resolution cryo-EM structure of β -galactosidase in
35 complex with a cell-permeant inhibitor', *Science*, 348(6239), pp. 1147-1151.

36 Beckers, M., Jakobi, A. J. and Sachse, C. (2019) 'Thresholding of cryo-EM density
37 maps by false discovery rate control', *IUCrJ*, 6(1), pp. 18-33.

38 Bernstein, K. A., Gallagher, J. E., Mitchell, B. M., Granneman, S. and Baserga, S. J.
39 (2004) 'The small-subunit processome is a ribosome assembly intermediate', *Eukaryotic
40 cell*, 3(6), pp. 1619-1626.

- 1 Davis, J. H., Tan, Y. Z., Carragher, B., Potter, C. S., Lyumkis, D. and Williamson, J. R.
2 (2016) 'Modular assembly of the bacterial large ribosomal subunit', *Cell*, 167(6), pp.
3 1610-1622. e15.
- 4 Elmlund, D. and Elmlund, H. (2012) 'SIMPLE: Software for ab initio reconstruction of
5 heterogeneous single-particles', *Journal of structural biology*, 180(3), pp. 420-427.
- 6 Fernandez-Leiro, R. and Scheres, S. H. (2016) 'Unravelling biological macromolecules
7 with cryo-electron microscopy', *Nature*, 537(7620), pp. 339-346.
- 8 Gao, H., Valle, M., Ehrenberg, M. and Frank, J. (2004) 'Dynamics of EF-G interaction
9 with the ribosome explored by classification of a heterogeneous cryo-EM dataset',
10 *Journal of structural biology*, 147(3), pp. 283-290.
- 11 Grant, T., Rohou, A. and Grigorieff, N. (2018) 'cisTEM, user-friendly software for single-
12 particle image processing', *elife*, 7, pp. e35383.
- 13 Gray, R. (1984) 'Vector quantization', *IEEE Assp Magazine*, 1(2), pp. 4-29.
- 14 Harnpicharnchai, P., Jakovljevic, J., Horsey, E., Miles, T., Roman, J., Rout, M.,
15 Meagher, D., Imai, B., Guo, Y. and Brame, C. J. (2001) 'Composition and functional
16 characterization of yeast 66S ribosome assembly intermediates', *Molecular cell*, 8(3),
17 pp. 505-515.
- 18 Jomaa, A., Jain, N., Davis, J. H., Williamson, J. R., Britton, R. A. and Ortega, J. (2014)
19 'Functional domains of the 50S subunit mature late in the assembly process', *Nucleic
20 Acids Research*, 42(5), pp. 3419-3435.
- 21 Jomaa, A., Stewart, G., Martín-Benito, J., Zielke, R., Campbell, T. L., Maddock, J. R.,
22 Brown, E. D. and Ortega, J. (2011) 'Understanding ribosome assembly: the structure of
23 in vivo assembled immature 30S subunits revealed by cryo-electron microscopy', *Rna*,
24 17(4), pp. 697-709.
- 25 Klaholz, B. P. (2015) 'Structure sorting of multiple macromolecular states in
26 heterogeneous cryo-EM samples by 3D multivariate statistical analysis', *Open Journal
27 of Statistics*, 5(07), pp. 820.
- 28 Kühlbrandt, W. (2014) 'The resolution revolution', *Science*, 343(6178), pp. 1443-1444.
- 29 Li, N., Chen, Y., Guo, Q., Zhang, Y., Yuan, Y., Ma, C., Deng, H., Lei, J. and Gao, N.
30 (2013) 'Cryo-EM structures of the late-stage assembly intermediates of the bacterial
31 50S ribosomal subunit', *Nucleic acids research*, 41(14), pp. 7073-7083.
- 32 Liao, H. Y., Hashem, Y. and Frank, J. (2015) 'Efficient estimation of three-dimensional
33 covariance and its application in the analysis of heterogeneous samples in cryo-electron
34 microscopy', *Structure*, 23(6), pp. 1129-1137.
- 35 Loerke, J., Giesebrecht, J. and Spahn, C. M. (2010) 'Multiparticle cryo-EM of
36 ribosomes', *Methods in enzymology*: Elsevier, pp. 161-177.
- 37 Ludtke, S. J. (2016) 'Single-particle refinement and variability analysis in EMAN2. 1',
38 *Methods in enzymology*: Elsevier, pp. 159-189.
- 39 Lyumkis, D., Brilot, A. F., Theobald, D. L. and Grigorieff, N. (2013) 'Likelihood-based
40 classification of cryo-EM images using FREALIGN', *Journal of structural biology*, 183(3),
41 pp. 377-388.
- 42 Nakane, T., Kimanius, D., Lindahl, E. and Scheres, S. H. (2018) 'Characterisation of
43 molecular motions in cryo-EM single-particle data by multi-body refinement in RELION',
44 *Elife*, 7, pp. e36861.

- 1 Nakane, T., Kotecha, A., Sente, A., McMullan, G., Masiulis, S., Brown, P. M., Grigoras,
2 I. T., Malinauskaite, L., Malinauskas, T. and Miehling, J. (2020) 'Single-particle cryo-EM
3 at atomic resolution', *Nature*, 587(7832), pp. 152-156.
- 4 Ni, X., Davis, J. H., Jain, N., Razi, A., Benlekbir, S., McArthur, A. G., Rubinstein, J. L.,
5 Britton, R. A., Williamson, J. R. and Ortega, J. (2016) 'YphC and YsxC GTPases assist
6 the maturation of the central protuberance, GTPase associated region and functional
7 core of the 50S ribosomal subunit', *Nucleic acids research*, 44(17), pp. 8442-8455.
- 8 Nikolay, R., Hilal, T., Qin, B., Mielke, T., Bürger, J., Loeke, J., Textoris-Taube, K.,
9 Nierhaus, K. H. and Spahn, C. M. (2018) 'Structural visualization of the formation and
10 activation of the 50S ribosomal subunit during in vitro reconstitution', *Molecular cell*,
11 70(5), pp. 881-893. e3.
- 12 Pettersen, E. F., Goddard, T. D., Huang, C. C., Couch, G. S., Greenblatt, D. M., Meng,
13 E. C. and Ferrin, T. E. (2004) 'UCSF Chimera—a visualization system for exploratory
14 research and analysis', *Journal of computational chemistry*, 25(13), pp. 1605-1612.
- 15 Punjani, A. and Fleet, D. J. (2021a) '3D Flexible Refinement: Structure and Motion of
16 Flexible Proteins from Cryo-EM', *bioRxiv*.
- 17 Punjani, A. and Fleet, D. J. (2021b) '3D Variability Analysis: Resolving continuous
18 flexibility and discrete heterogeneity from single particle cryo-EM', *Journal of Structural
19 Biology*, 213(2), pp. 107702.
- 20 Rabuck-Gibbons, J. N., Popova, A. M., Greene, E. M., Cervantes, C. F., Lyumkis, D.
21 and Williamson, J. R. (2020) 'SrmB rescues trapped ribosome assembly intermediates',
22 *Journal of molecular biology*, 432(4), pp. 978-990.
- 23 Razi, A., Guarné, A. and Ortega, J. (2017) 'The cryo-EM structure of YjeQ bound to the
24 30S subunit suggests a fidelity checkpoint function for this protein in ribosome
25 assembly', *Proceedings of the National Academy of Sciences*, 114(17), pp. E3396-
26 E3403.
- 27 Russo, C. J. and Passmore, L. A. (2014) 'Ultrastable gold substrates for electron
28 cryomicroscopy', *Science*, 346(6215), pp. 1377-1380.
- 29 Sashital, D. G., Greeman, C. A., Lyumkis, D., Potter, C. S., Carragher, B. and
30 Williamson, J. R. (2014) 'A combined quantitative mass spectrometry and electron
31 microscopy analysis of ribosomal 30S subunit assembly in *E. coli*', *Elife*, 3, pp. e04491.
- 32 Scheres, S. H. (2012) 'RELION: implementation of a Bayesian approach to cryo-EM
33 structure determination', *Journal of structural biology*, 180(3), pp. 519-530.
- 34 Scheres, S. H. (2016) 'Processing of structurally heterogeneous cryo-EM data in
35 RELION', *Methods in enzymology*: Elsevier, pp. 125-157.
- 36 Scheres, S. H., Núñez-Ramírez, R., Sorzano, C. O., Carazo, J. M. and Marabini, R.
37 (2008) 'Image processing for electron microscopy single-particle analysis using XMIPP',
38 *Nature protocols*, 3(6), pp. 977-990.
- 39 Shajani, Z., Sykes, M. T. and Williamson, J. R. (2011) 'Assembly of bacterial
40 ribosomes', *Annual review of biochemistry*, 80, pp. 501-526.
- 41 Sorzano, C., Bilbao-Castro, J., Shkolnisky, Y., Alcorlo, M., Melero, R., Caffarena-
42 Fernández, G., Li, M., Xu, G., Marabini, R. and Carazo, J. (2010) 'A clustering approach
43 to multireference alignment of single-particle projections in electron microscopy', *Journal
44 of structural biology*, 171(2), pp. 197-206.

- 1 Spahn, C. M. and Penczek, P. A. (2009) 'Exploring conformational modes of
2 macromolecular assemblies by multiparticle cryo-EM', *Current opinion in structural
3 biology*, 19(5), pp. 623-631.
- 4 Stokes, J. M., Davis, J. H., Mangat, C. S., Williamson, J. R. and Brown, E. D. (2014)
5 'Discovery of a small molecule that inhibits bacterial ribosome biogenesis', *Elife*, 3, pp.
6 e03574.
- 7 Sykes, M. T., Shajani, Z., Sperling, E., Beck, A. H. and Williamson, J. R. (2010)
8 'Quantitative proteomic analysis of ribosome assembly and turnover in vivo', *Journal of
9 molecular biology*, 403(3), pp. 331-345.
- 10 Tan, Y. Z., Aiyer, S., Mietzsch, M., Hull, J. A., McKenna, R., Grieger, J., Samulski, R. J.,
11 Baker, T. S., Agbandje-McKenna, M. and Lyumkis, D. (2018) 'Sub-2 Å Ewald curvature
12 corrected structure of an AAV2 capsid variant', *Nature communications*, 9(1), pp. 1-11.
- 13 Tan, Y. Z., Baldwin, P. R., Davis, J. H., Williamson, J. R., Potter, C. S., Carragher, B.
14 and Lyumkis, D. (2017a) 'Addressing preferred specimen orientation in single-particle
15 cryo-EM through tilting', *Nature methods*, 14(8), pp. 793-796.
- 16 Tan, Y. Z., Baldwin, P. R., Davis, J. H., Williamson, J. R., Potter, C. S., Carragher, B.
17 and Lyumkis, D. (2017b) 'Addressing preferred specimen orientation in single-particle
18 cryo-EM through tilting', *Nature methods*, 14(8), pp. 793.
- 19 Uicker, W. C., Schaefer, L. and Britton, R. A. (2006) 'The essential GTPase RbgA
20 (YlqF) is required for 50S ribosome assembly in *Bacillus subtilis*', *Molecular
21 microbiology*, 59(2), pp. 528-540.
- 22 Wang, Q., Matsui, T., Domitrovic, T., Zheng, Y., Doerschuk, P. C. and Johnson, J. E.
23 (2013) 'Dynamics in cryo EM reconstructions visualized with maximum-likelihood
24 derived variance maps', *Journal of structural biology*, 181(3), pp. 195-206.
- 25 White, H., Ignatiou, A., Clare, D. and Orlova, E. (2017) 'Structural study of
26 heterogeneous biological samples by cryoelectron microscopy and image processing',
27 *BioMed research international*, 2017.
- 28 Wolfram Research, I. (2020) *Mathematica*. Version 12.2 edn. Champaign, Illinois:
29 Wolfram Research, Inc.
- 30 Yip, K. M., Fischer, N., Paknia, E., Chari, A. and Stark, H. (2020) 'Atomic-resolution
31 protein structure determination by cryo-EM', *Nature*, 587(7832), pp. 157-161.
- 32 Zhang, K., Pintilie, G. D., Li, S., Schmid, M. F. and Chiu, W. (2020) 'Resolving individual
33 atoms of protein complex by cryo-electron microscopy', *Cell research*, 30(12), pp. 1136-
34 1139.
- 35 Zhong, E. D., Bepler, T., Berger, B. and Davis, J. H. (2021) 'CryoDRGN: reconstruction
36 of heterogeneous cryo-EM structures using neural networks', *Nature Methods*, 18(2),
37 pp. 176-185.

38

39 **Materials and Methods**

40

41 *Cell Growth and Isolation of Ribosomal Particles*

42 Cells were grown and ribosomal particles were isolated as in (Davis *et al.*, 2016).
43 Briefly, strain JD321 was grown in M9 media (48mM Na₂HPO₄, 22mM KH₂PO₄,
44 8.5mM NaCl, 10mM MgCl₂, 10mM MgSO₄, 5.6mM glucose, 50mM Na₃*EDTA, 25mM
45 CaCl₂, 50mM FeCl₃, 0.5mM ZnSO₄, 0.5mM CuSO₄, 0.5mM MnSO₄, 0.5mM CoCl₂,

1 0.04mM d-biotin, 0.02mM folic acid, 0.08mM vitamin B1, 0.11mM calcium pantothenate,
2 0.4nM vitamin B12, 0.2mM nicotinamide, 0.07mM riboflavin, and 7.6mM
3 (14NH₄)₂SO₄] with tetracycline (10 mg/mL), chloramphenicol (35 mg/mL), and limiting
4 conditions HSL (0.1 nM) and harvested at OD=0.5. Cells were lysed in Buffer A (20mM
5 Tris-HCl, 100mM NH₄Cl, 10mM MgCl₂, 0.5mM EDTA, 6mM β-mercaptoethanol; pH
6 7.5) by a mini bead beater, and the clarified lysate was fractionated on a 10-40% w/v
7 sucrose gradient (50mM Tris-HCl, 100mM NH₄Cl, 10mM MgCl₂, 0.5mM EDTA, 6mM β-
8 mercaptoethanol; pH 7.5).

9

10 *Electron Microscopy Data Collection*

11 Fractions containing the ribosomal intermediates were spin-concentrated with a 100
12 kDa MW filter (Amicon) and buffer exchanged into Buffer A. 3 μl of this sample was
13 added to a plasma cleaned (Gatan, Solarus) 1.2mm hole, 1.3mm spacing holey gold
14 grids (Russo and Passmore, 2014). Grids were manually frozen in liquid ethane, and
15 single particle data was collected using Legion on a Titan Krios microscope (FEI) with
16 a K2 summit direct detector (Gatan) in super-resolution mode (pixel size of 0.66Å at
17 22,500 magnification). A dose rate of ~5.8e⁻⁷/pix/sec was collected across 50 frames
18 with a fluence of 33-35e⁻⁷/Å² at a tilt of -20° to compensate for preferred orientation (Tan
19 *et al.*, 2017b).

20

21 *FrealignX Classifications*

22 After conversion from Relion to FrealignX parameters, global refinements were
23 performed in FrealignX, and all occupancies were randomized across the parameter
24 files. A final value of 20Å was selected for *res_high_class*, and after every 10 cycles of
25 classification/refinement, all classes were aligned to a C class scaffold using custom
26 scripts for a 3D alignment with Chimera (Pettersen *et al.*, 2004) while running FrealignX.
27 For each classification step, 50 refinement/classification cycles were performed. After
28 initial classification, each class was selected in a parameter file for subsequent rounds
29 of classification using the merge_classes.exe in cisTEM (Grant, Rohou and Grigorieff,
30 2018) and custom scripts. The occupancies were randomized across the parameter
31 files, and the same cycle of 50 cycles of refinement/classification interspersed with 3D
32 alignment with Chimera every 10 cycles. FSC curves and Euler plots were generated by
33 FrealignX and cisTEM (Grant, Rohou and Grigorieff, 2018), and 3DFSC plots were
34 calculated by the 3DFSC server (Tan *et al.*, 2017a). The SCF was calculated according
35 to the process in (Baldwin and Lyumkis, 2021; Baldwin and Lyumkis, 2020). The
36 3DFSCs and all maps shown were visualized in Chimera (Pettersen *et al.*, 2004), and
37 the details for each map are indicated in Table S1.

38

39 *Calculation of σ values*

1 For analysis, each map was first filtered to 10Å. To calculate σ which was used as a
2 measure of noise, each map was unmasked by expanding the *outer_radius* in FrealignX
3 so that the spherical particle mask would be larger than the box size. The Fourier
4 folding of signal along the edges of the box was negligible. Relion 2.1 was used to
5 calculate the σ value using the *relion_image_handler* command. Relion 2.1 was then
6 used to create binarized maps using the *relion_image_handler* command, and the
7 binarization threshold was set to 3σ .

8

9 *Hierarchical clustering analysis*

10 Thresholded, binarized maps were given as input to a custom Mathematica script
11 (Wolfram Research, 2020). The Mathematica script calculated the segmented
12 difference maps between all maps and calculated the molecular weights of the
13 differences maps (in kilodaltons) using Equation 1 (Ludtke, 2016):

$$MW = n_{pixels} * pixelsize^3 * \rho / 1000$$

14 Density ρ is 0.81 daltons/Å³. The MW difference matrix was clustered using the
15 Euclidean distance metric and Ward's linkage and displayed in a dendrogram. Similar
16 maps were averaged together after hierarchical clustering analysis using EMAN2
17 ((Ludtke, 2016).

18

19 *Occupancy Analysis*

20 The thresholded and binarized maps were given as input, and the reference map from
21 the *E. coli* 50S subunit crystal structure (PDB ID 4YBB) was segmented into 139
22 elements comprised of individual ribosomal proteins and rRNA helices according to the
23 23S secondary structure. Theoretical densities for each r-protein and rRNA helix were
24 calculated for each element at 10Å using the *pdb2mrc* command from EMAN. Prior to
25 binarization, voxels that had overlapping theoretical density from two structural
26 elements, were assigned to the smaller of the two theoretical volumes so that each pair
27 of volumes is nonoverlapping. Each voxel density was binarized to either 0 or 1 using a
28 threshold of 0.016, which is the threshold that gave the approximately correct molecular
29 weight for individual r-proteins and rRNAs helices. The relative volumes in the binarized
30 experimental and reference maps were calculated, which gave a fractional occupancy
31 between 0 and 1 for each element. The occupancy values were clustered across the
32 rows (classes) and columns (rRNA/protein elements) using an unsupervised
33 hierarchical clustering using the Euclidean distance metric and Ward's linkage method,
34 as implemented in Mathematica.

35

36 *Parsimonious folding/unfolding matrices.* A pathway diagram was constructed by using
37 the $n \times n$ molecular weight difference matrices, \mathbf{M}_f and \mathbf{M}_u , from a set of n structures.
38 Each difference map ($M_i - M_j$) has negative elements corresponding to folding that occurs
39 in the transition from class i to class j , and positive elements that correspond to

1 unfolding that occurs in the transition from class i to class j . The volume changes for
2 folding and unfolding form the elements of \mathbf{M}_f or \mathbf{M}_u , noting that $\mathbf{M}_u = \mathbf{M}_f^T$. The
3 matrices \mathbf{M}_f and \mathbf{M}_u are used to construct a directed graph \mathbf{G} , comprised of the set of
4 vertices v_i , and a set of directed edges, e_{ij} , representing the allowed transitions
5 between classes. The set of edges is initialized as the set of e_{ij} where $M_{u,i,j} > M_{f,i,j}$, such
6 that only net folding transitions are allowed. The set of edges is pruned using two
7 global parameters: θ_{unf} as a maximum threshold for unfolding, and n_{branch} , as a limit on
8 the number of transitions emanating from a single class. The unfolding threshold limits
9 unreasonable structural rearrangements, while the branching threshold limits transitions
10 to a small set of the closest transitions. Edges are eliminated if the unfolding exceeds
11 the threshold such that $M_{u,i,j} > \theta_{unf}$, unless elimination of the edge results in a
12 disconnected graph \mathbf{G} . Next, for each vertex v_i , the set of remaining edges e_{ik}
13 emanating from v_i , are sorted into the order based on the $M_{f,i,k}$, retaining at most the
14 n_{branch} edges, again, unless deleting the edge would result in a disconnected graph \mathbf{G} .
15 The resulting transition graph \mathbf{G} should have one or more *source* vertices (classes) that
16 are the earliest classes in the assembly pathway, and one or more *sink* vertices that are
17 the most mature classes in the pathway. Tuning of the parameters θ_{unf} and n_{branch} ,
18 adjusts the connectivity and degree of branching of the resulting graph. The graph
19 vertices are annotated with thumbnails of the map, followed by manual layout of the
20 graph into a sensible order in Adobe Illustrator. The values of θ_{unf} and n_{branch} used to
21 generate the graph in Figure 7 were 390 kDa and 3, respectively.

22

23 *Data Deposition and Software Availability*

24 Mathematica scripts and example parameter files, where needed, will be available upon
25 request. All maps are deposited at EMPIAR as noted in the Key Resources Table.

26

27 **Figure Titles and Legends**

28

29 Figure 1. Description of the bacterial large ribosomal subunit and prior assembly
30 intermediates identified by cryo-EM. (A) PDB ID 4YBB labeled with prominent features
31 identifiable on the large ribosomal subunit, including the central protuberance (CP),
32 base, L1 stalk, and L7/12 stalks. These terms are used throughout the paper. (B)
33 Primary classes identified within the original bL17-lim dataset (Davis *et al.*, 2016). From
34 left to right: B class (red), C class (yellow), D class (green), and the E class (blue).

35

36 Figure 2. Workflow for cryo-EM heterogeneity analysis.

37

38 Figure 3. A divisive resolution-limited subclassification approach facilitates identifying
39 rare structural variants. (A) FrealignX classification with *res_high_class* parameter set to
40 Nyquist (5.24Å). (B) FrealignX classification with the *res_high_class* parameter set to

1 20Å. Using a lower resolution cutoff leads to the identification of a broader range of
2 classes. (C) Results of the five different classification schemes. The colors (A,B)
3 correspond to the classes found in (C).
4

5 Figure 4. Segmented difference analysis helps to define molecular weight differences
6 between map pairs. (A) Example where two maps would have been considered
7 different before segmentation, but are not different after segmentation. (B) Example
8 where two maps are different both before and after segmentation. Numbers indicate
9 positive (Map1-Map2) and negative (Map2-Map1) molecular weight differences.
10

11 Figure 5. Hierarchical clustering is used to combine similar maps under a given
12 threshold. (A) Hierarchical clustering analysis of the maps that result after the terminal
13 subclassification (total n=42). The red dashed line indicates the 10kDa MWCO used to
14 combine similar maps at this step, and the red stars indicate maps that are combined
15 after this analysis. After combining similar maps, the final number of classes is thus 41.
16 (B) Close-up example of two combined maps in (A). The leftmost structure is distinct
17 from the other two by ~135 kDa and needs to be treated independently, whereas the
18 two rightmost structures can be combined into a single class.
19

20 Figure 6. Results of occupancy analysis on the full dataset mapped onto the ribosomal
21 scaffold (A) Reference crystal structure 4ybb. (B) Binarized maps of the individual
22 proteins and rRNA helices created by segmenting the crystal structure into 139
23 individual helices and proteins, and calculating theoretical 10A maps in Chimera. (C) An
24 example of a binarized experimental map arising from sub-classification. The pixels
25 from the binarized experimental map that are located in the theoretical binarized map
26 are counted and normalized to an occupancy value of 0-1. (D) Example of an E class
27 (blue) where the occupancy of an rRNA helix (h82, salmon) is fully occupied, with the
28 corresponding occupancy block underneath. (E) Partial occupancy example of h82
29 (salmon) with a C class (yellow). (F) Example where rRNA (h82, salmon) is missing in
30 the experimental data (B class, red). G. Occupancy analysis plot, where the individual
31 proteins and helices are shown on the x-axis, the experimental maps are on the y-axis,
32 and the normalized occupancy values are shown from white (0) to dark blue (1).
33 Hierarchical clustering of both structure elements and experimental maps was
34 performed on the occupancy matrix using a squared Euclidean distance metric and
35 Ward's linkage. (H) Occupancy analysis blocks mapped back to the reference structure
36 4YBB, and the numbering system is the same as in (G).
37

38 Figure 7. Revised ribosome assembly map from bL17-lim. (Assembly pathway drawn by
39 analyzing the folding and unfolding molecular weight matrices and revised by hand).
40

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Supplemental Information Titles and Legends

Supplemental Figure 1. Results of the five tested classification schemes grouped by class. The structures are colored by classification scheme. Unique classes are shown by an asterisk (*), and classes that are similar are underscored by red brackets. Clustering of (A) the B classes, (B) the C classes, (C) the D classes, and (D) the E classes resulting from the tested classification schemes. Any 70S or “junk” classes that result from the subclassifications are omitted for clarity.

Supplemental Figure 2. Example of ambiguous density and features for B class particles. From left to right: (B class filtered to 5Å and shown at $3\sigma_{\text{map}}$, $2\sigma_{\text{map}}$, and $1.5\sigma_{\text{map}}$. In particular, at the $2\sigma_{\text{map}}$ threshold, noise above background is visible proximal to the main particle that is likely due to disordered rRNA (black arrows).

Supplemental Figure 3. Comparison of the the $3\sigma_{\text{map}}$ threshold that is used in our current analysis versus the confidence map FDR threshold (Beckers, 2019). The black line represents $y=x$, and the red and black dots represent thresholds at 1% and 0.01% FDR, respectively. The measures are highly correlated, and the $3\sigma_{\text{map}}$ threshold is generally more conservative than either FDR threshold.

Supplemental Figure 4. Hierarchical clustering analysis of the original bL17-lim data (orange) together with the structures solved by the new data processing workflow (blue). The red dotted line indicates the 10.0 kDa cutoff applied to determine similarity between classes. The original classes typically have counterparts within the new data (red underlined structures), but the new workflow is able to identify many more structural intermediates.

Figure 1

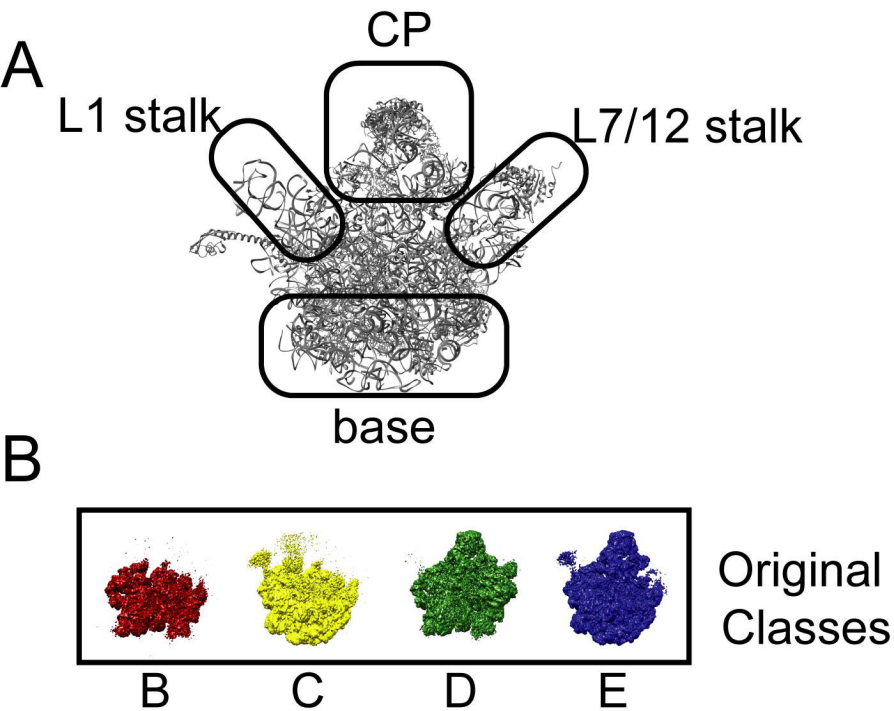


Figure 2

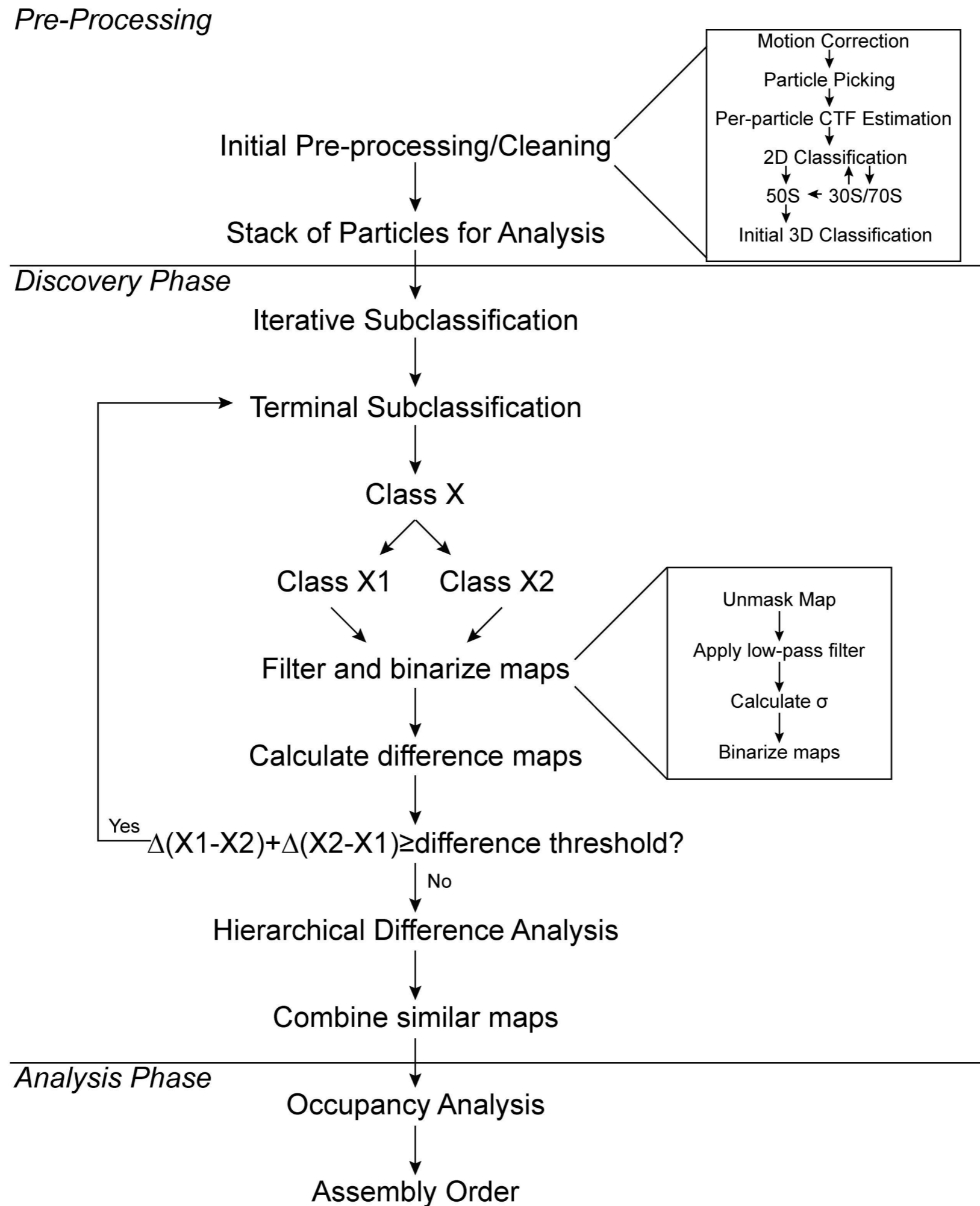
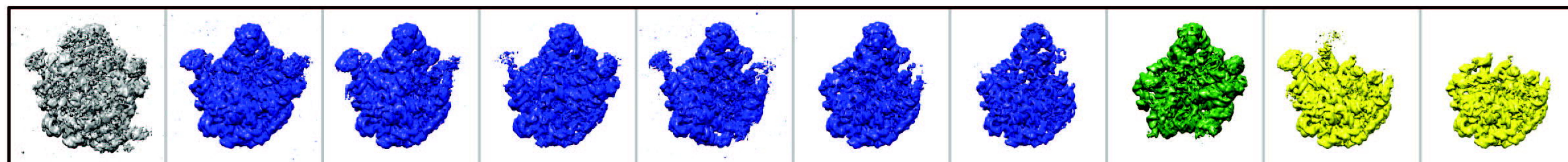


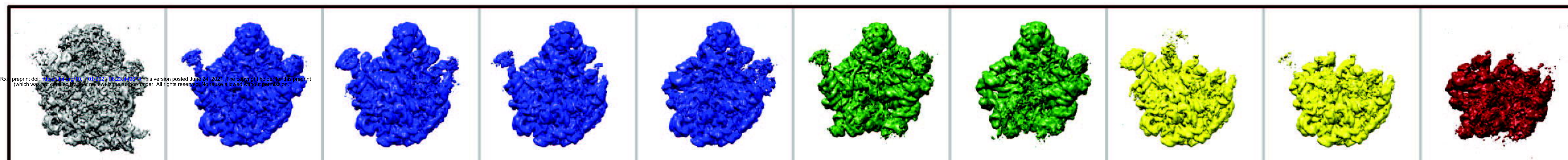
Figure 3

A



FrealignX
res_high_class=5.24Å

B



FrealignX
res_high_class=20Å

C

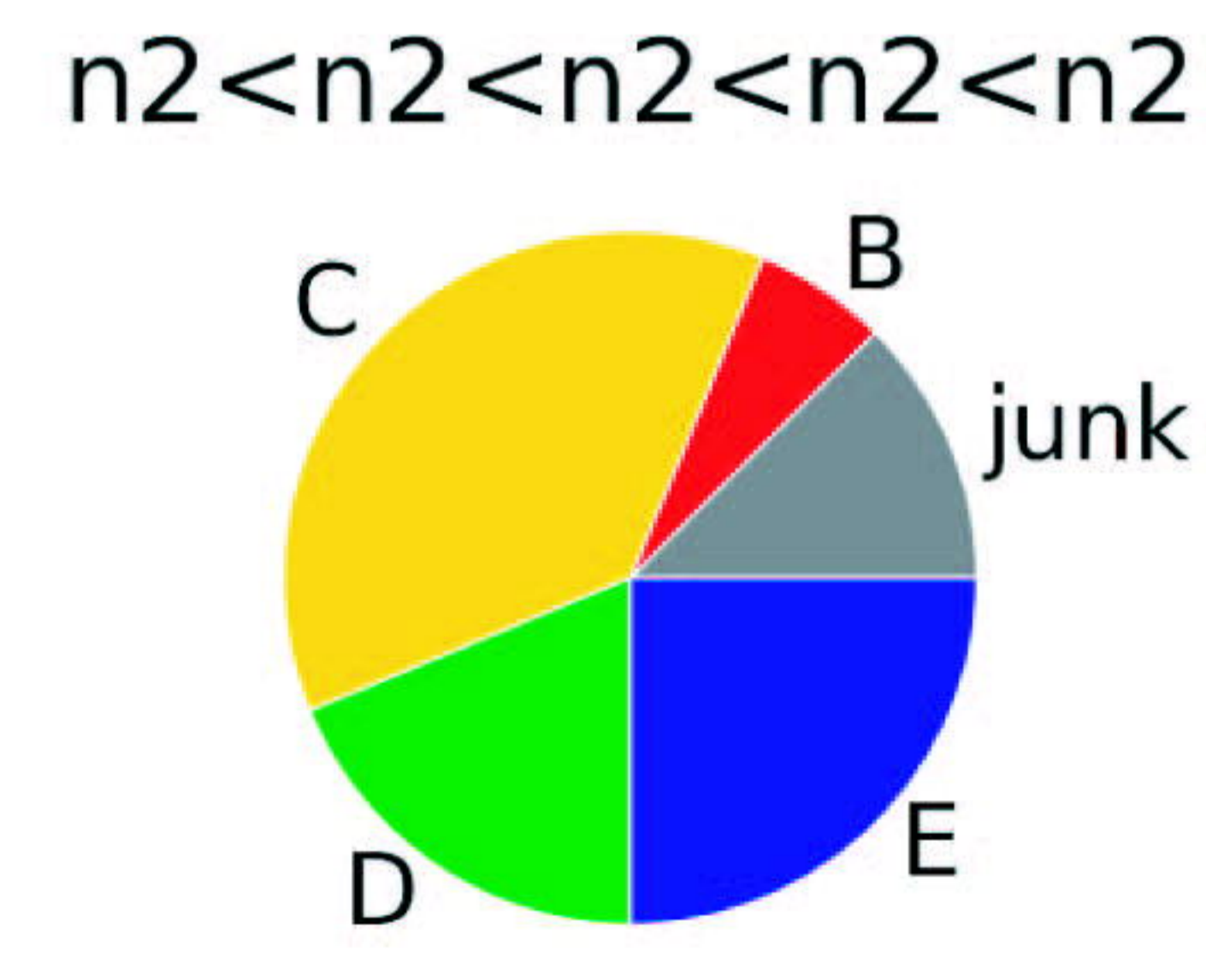
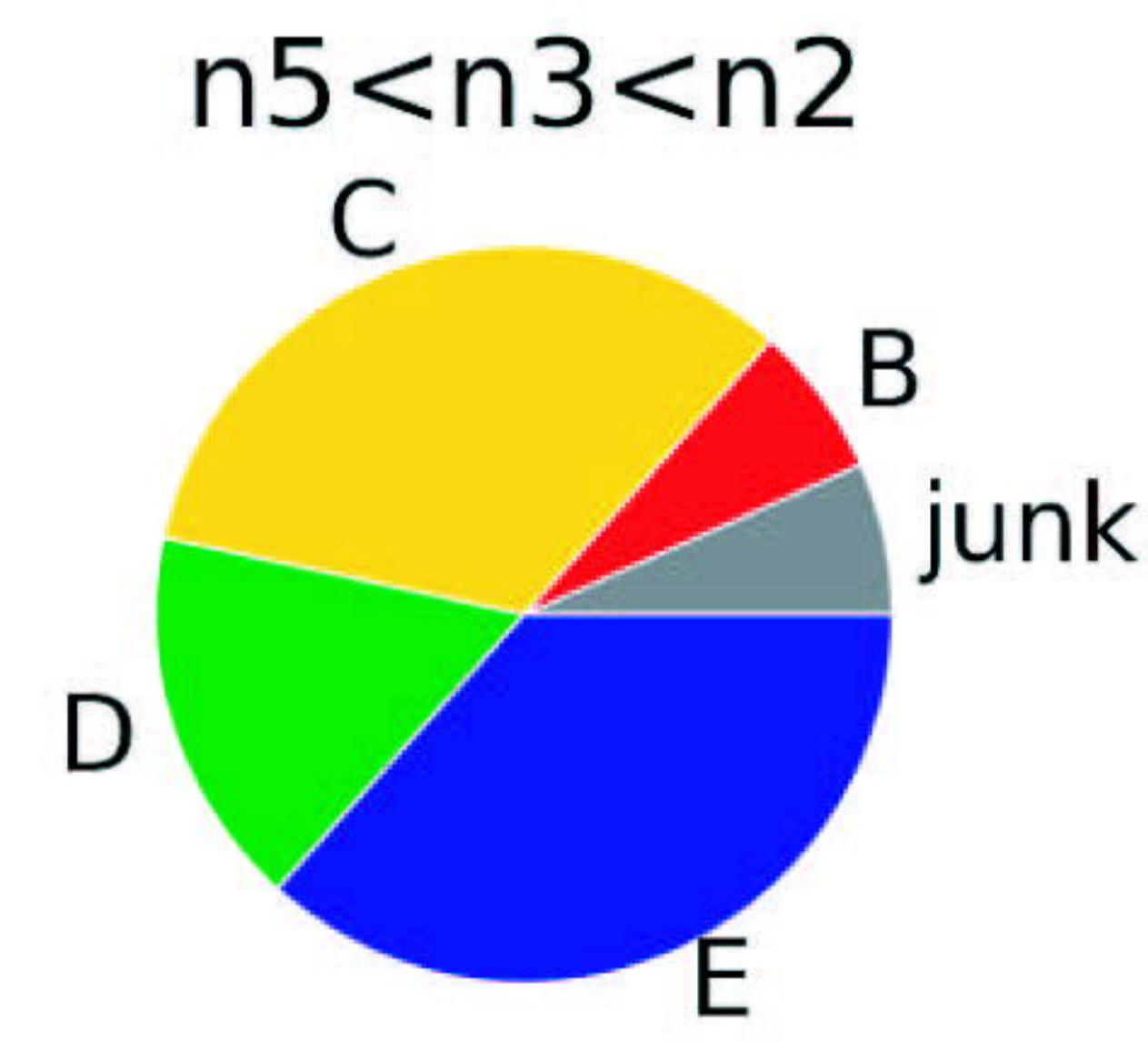
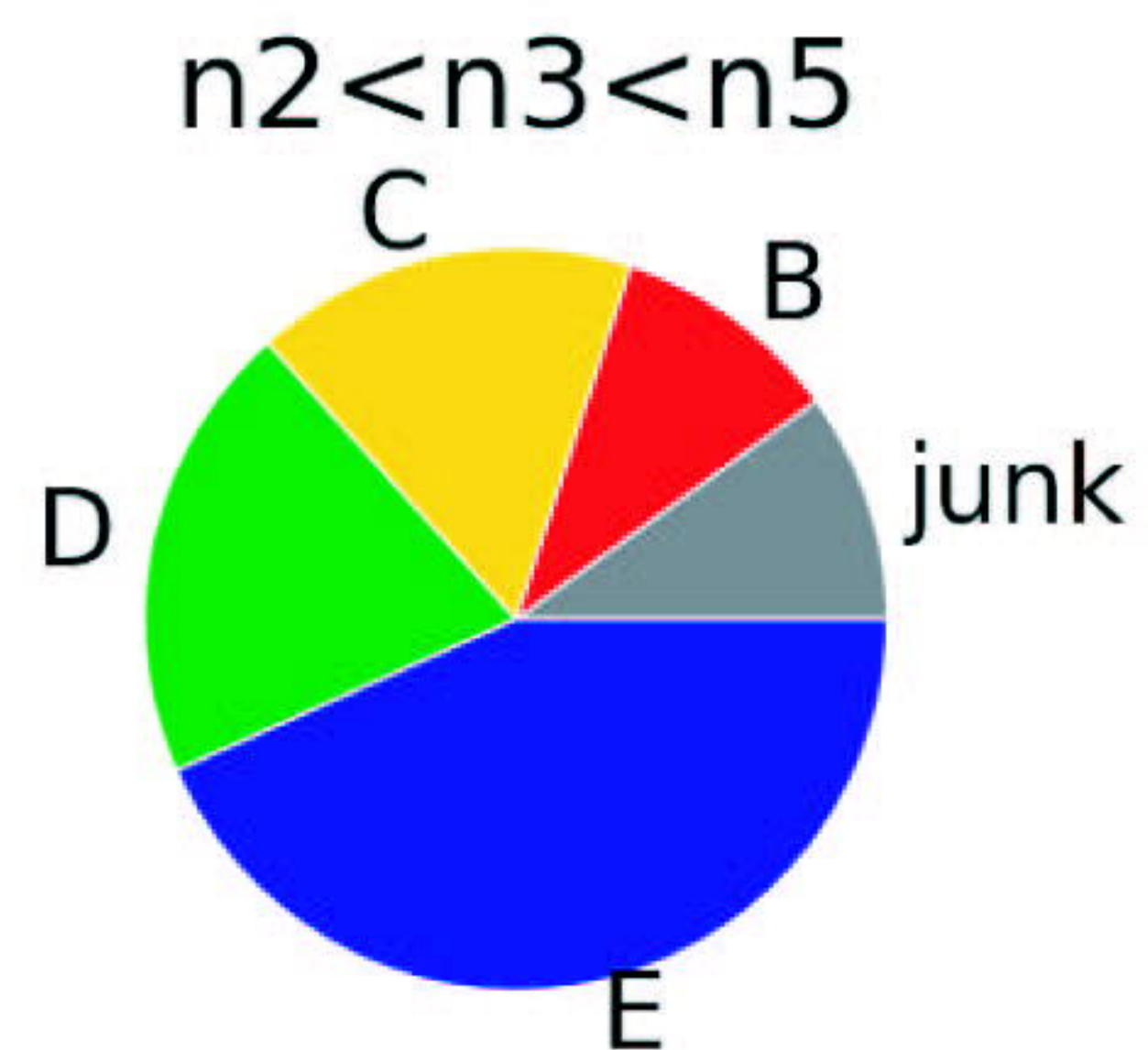
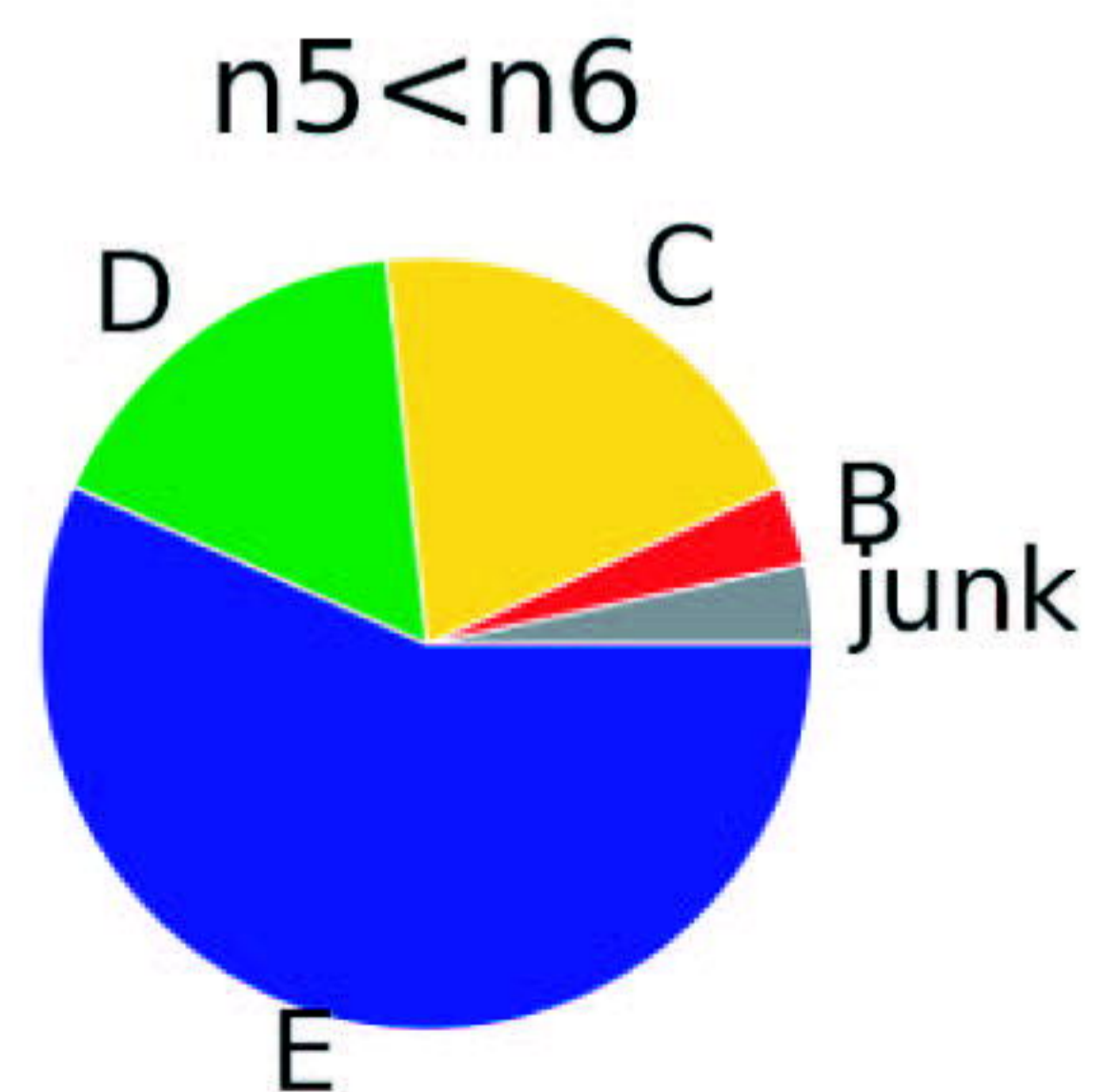
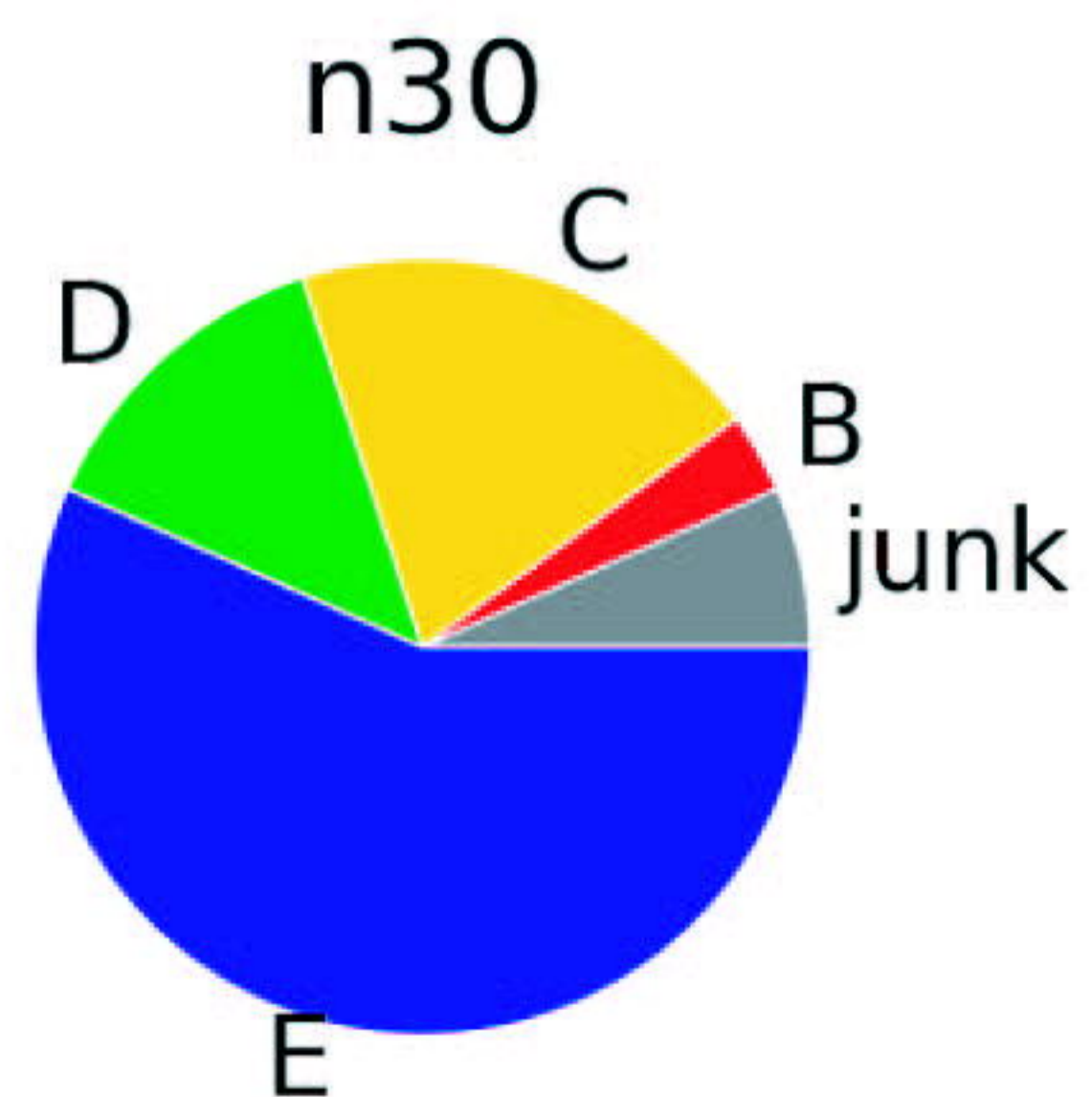
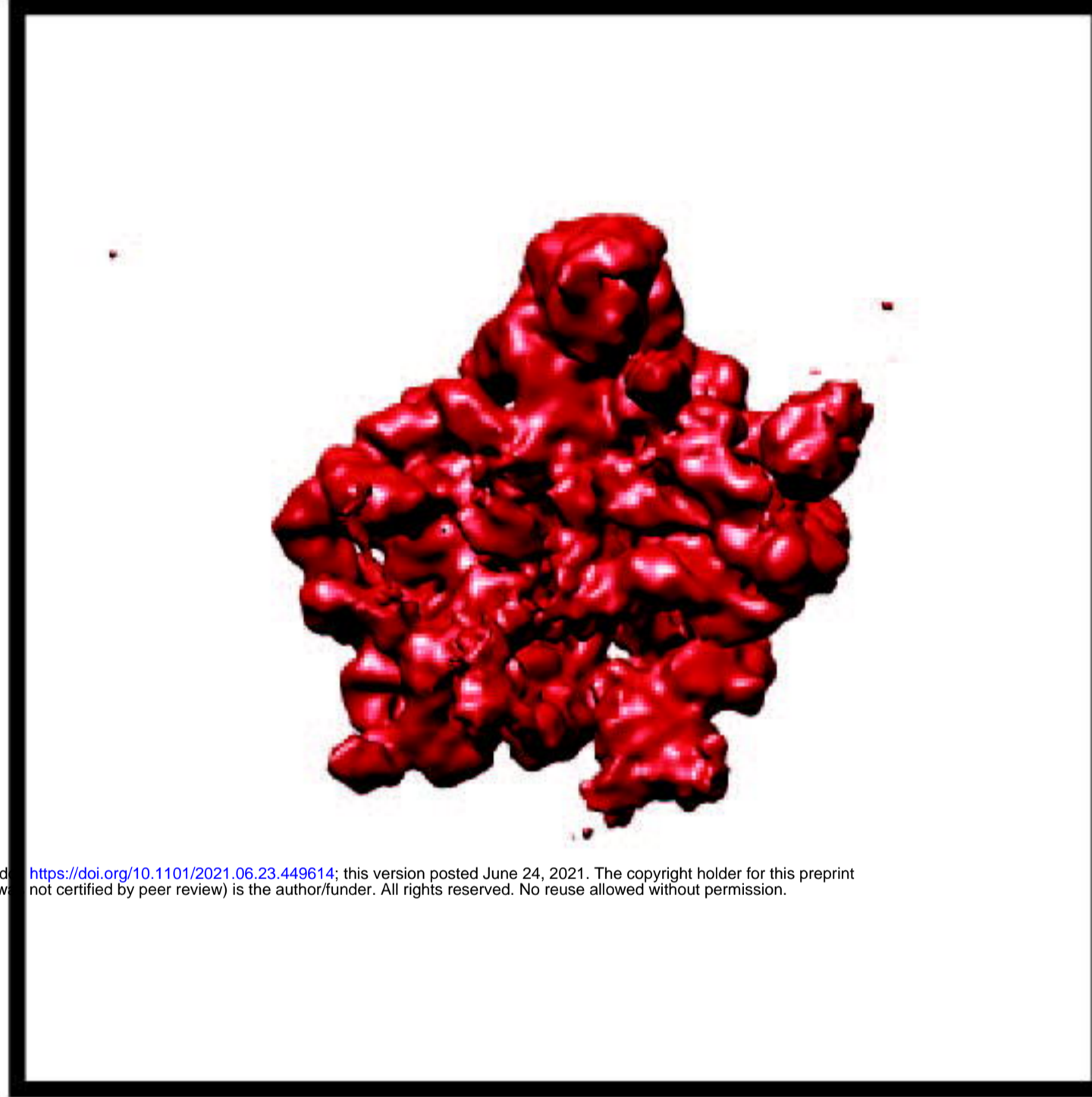


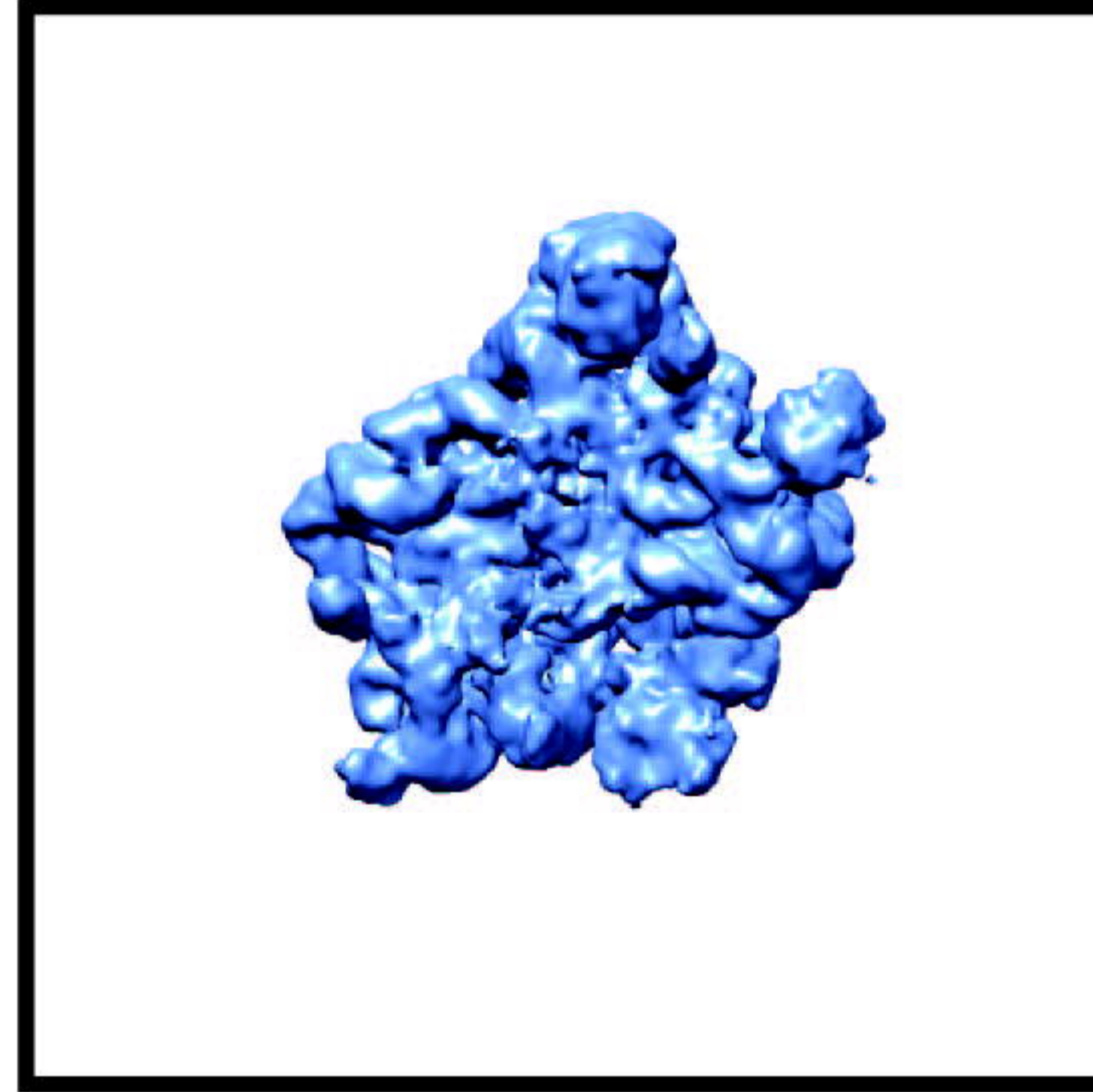
Figure 4

A

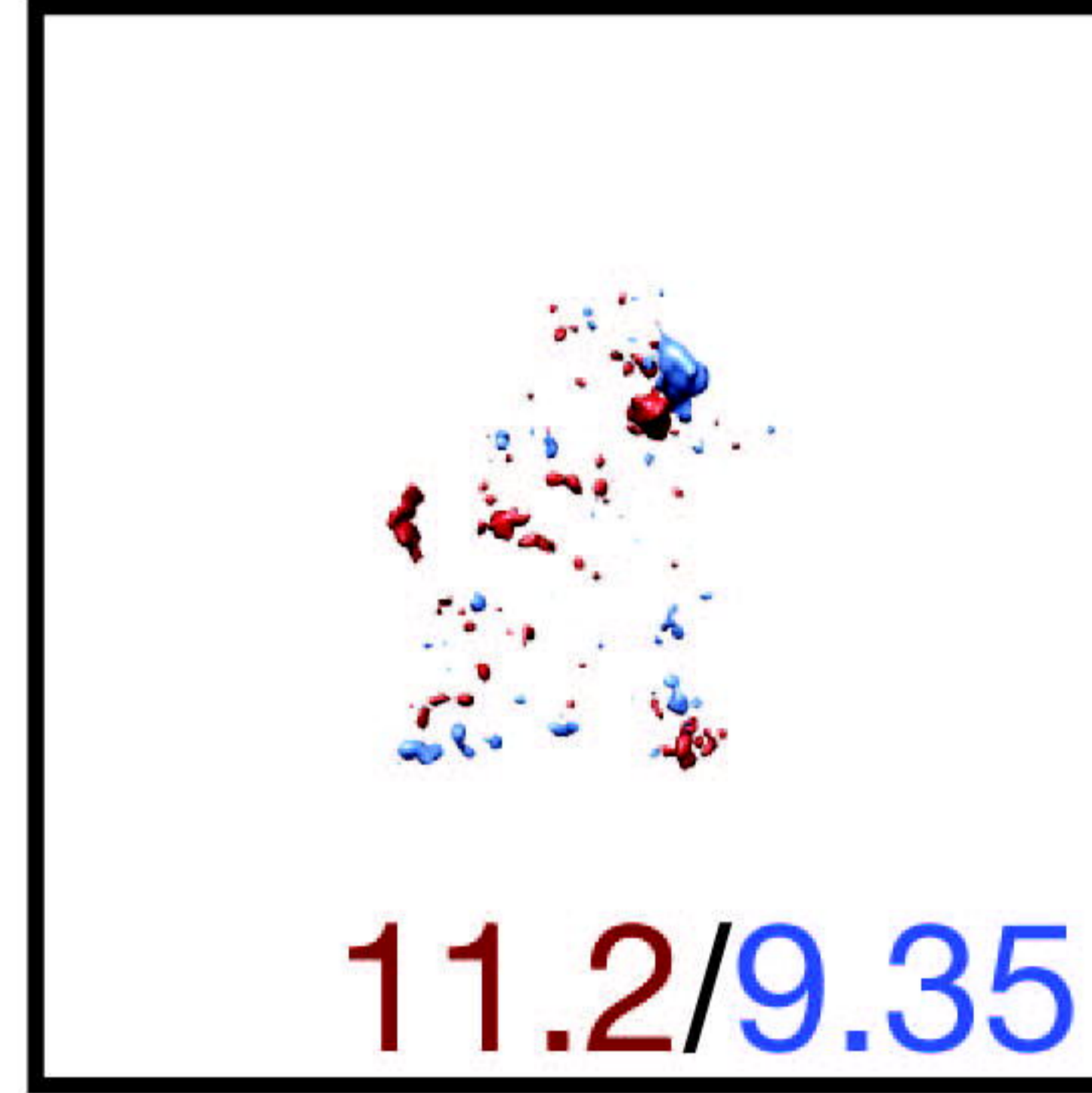
Map1



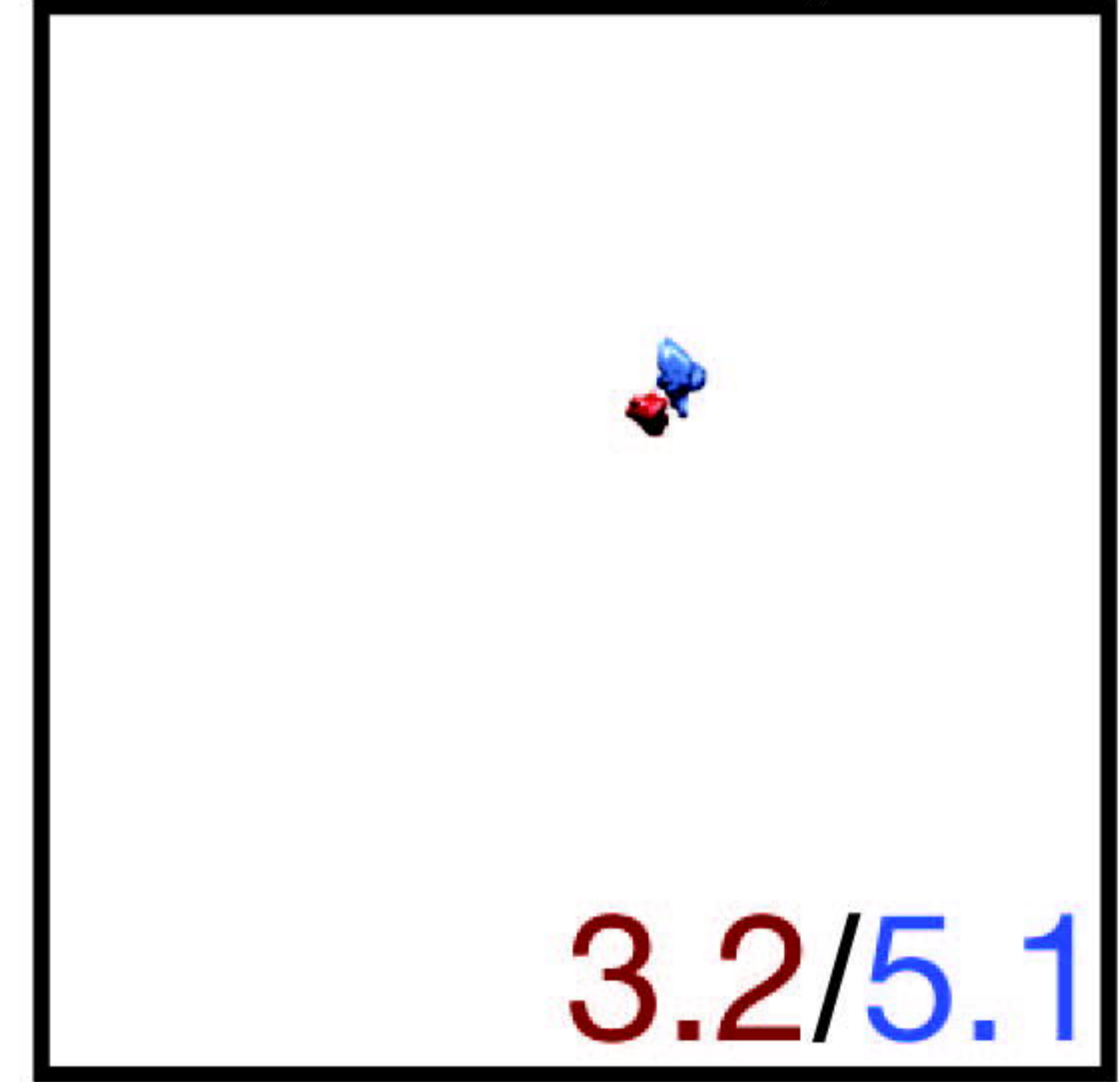
Map2



Difference (kDa)



Segmented
Difference (kDa)



bioRxiv preprint doi: <https://doi.org/10.1101/2021.06.23.449614>; this version posted June 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

B

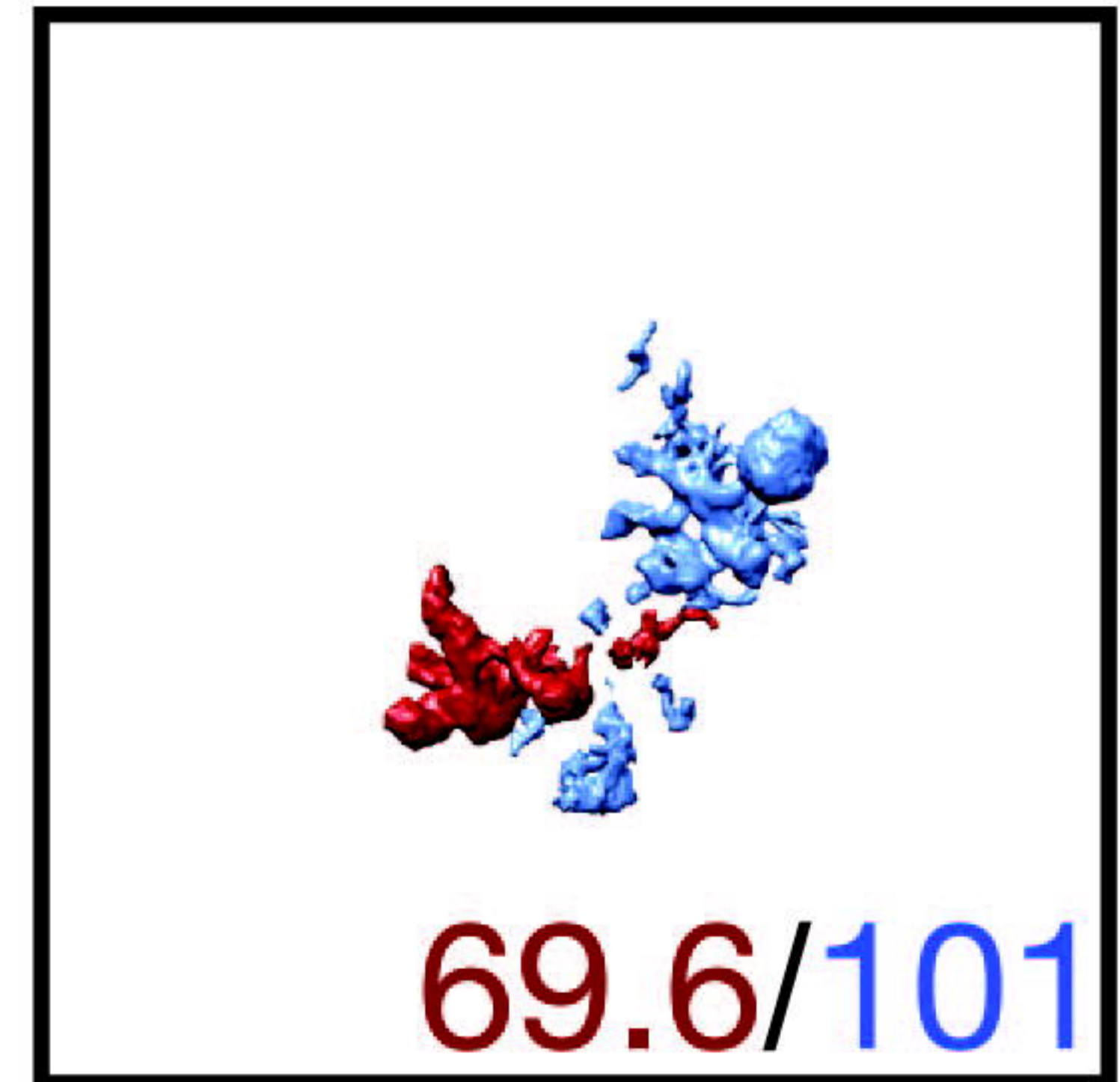
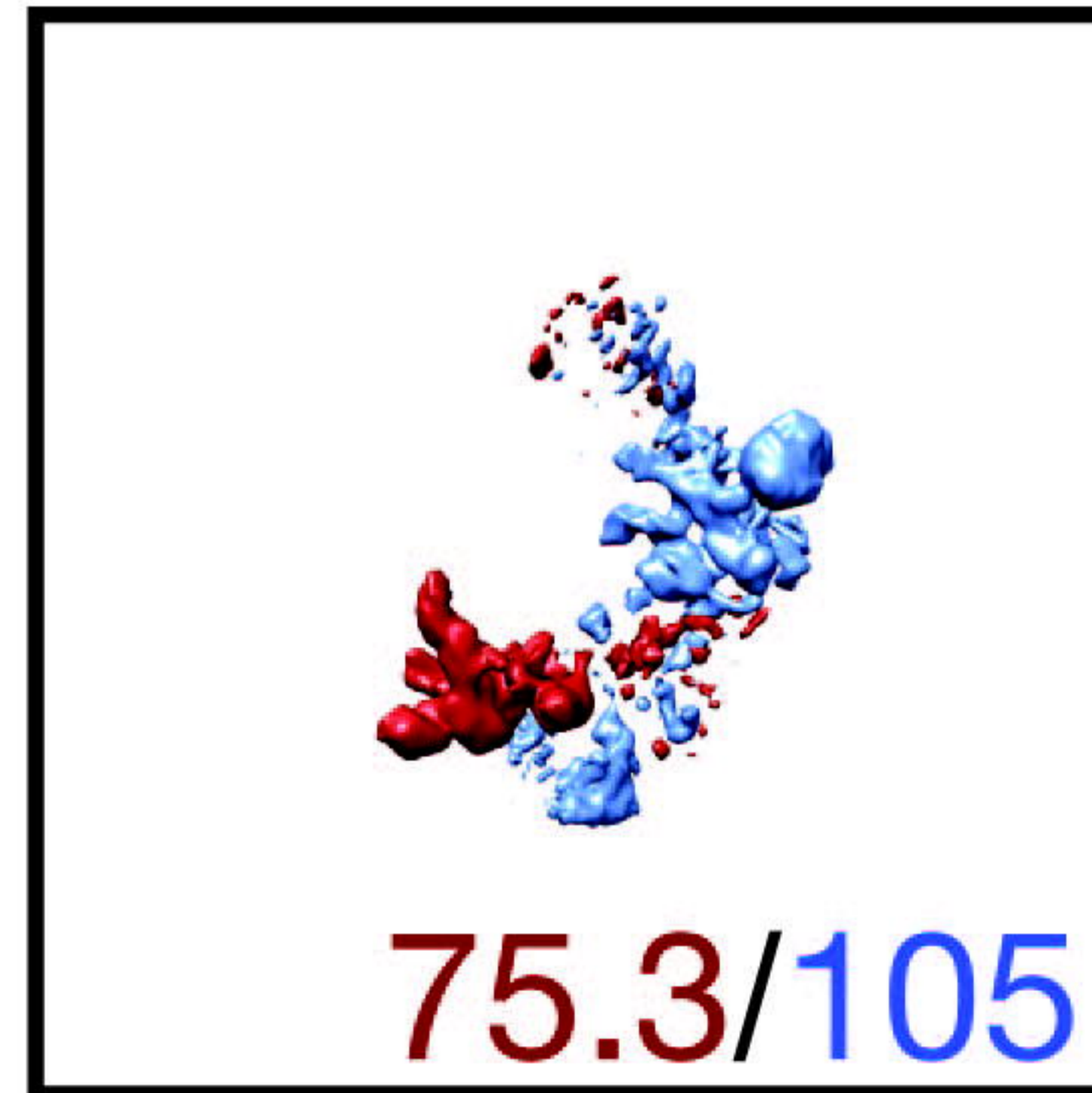
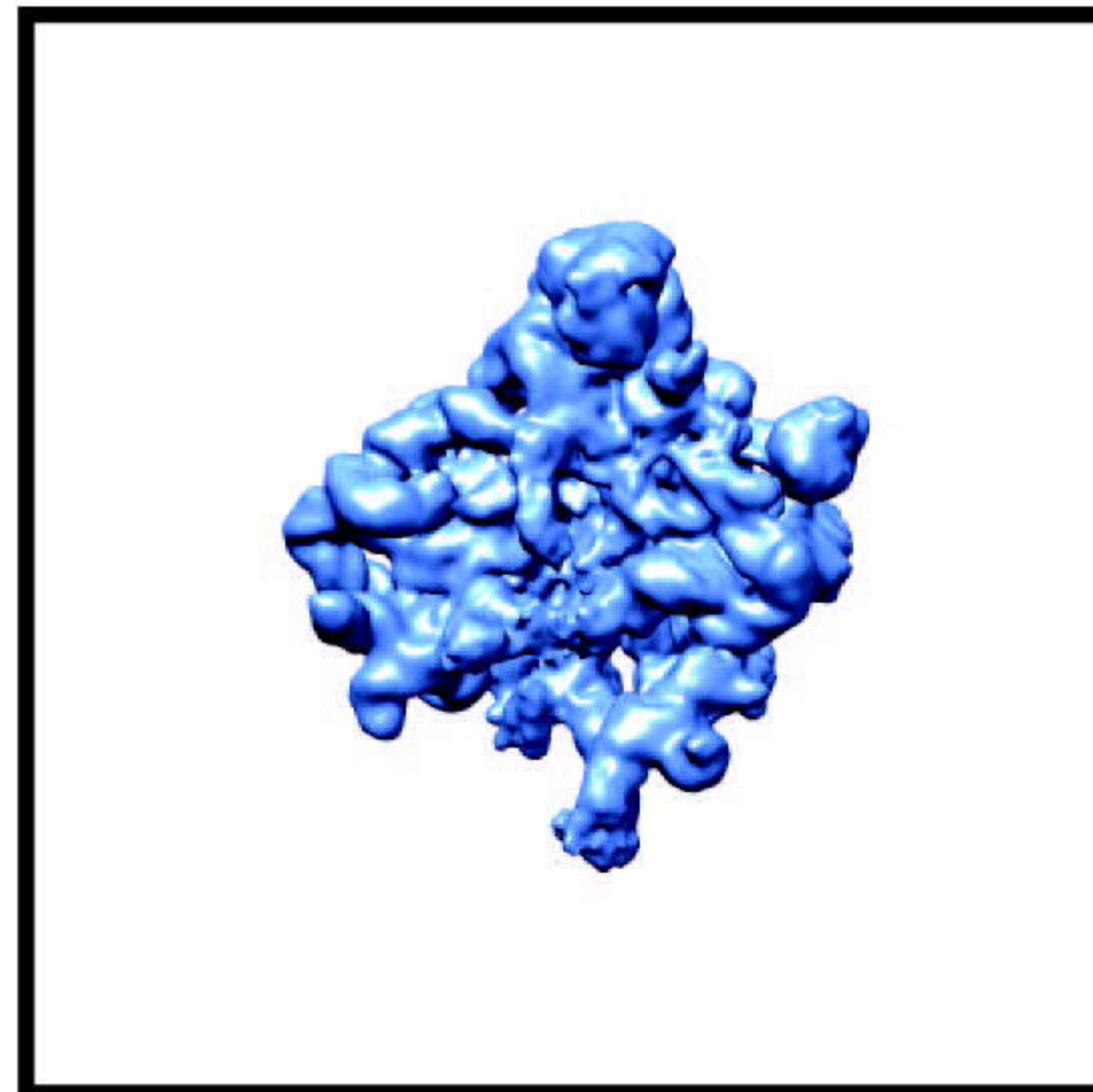
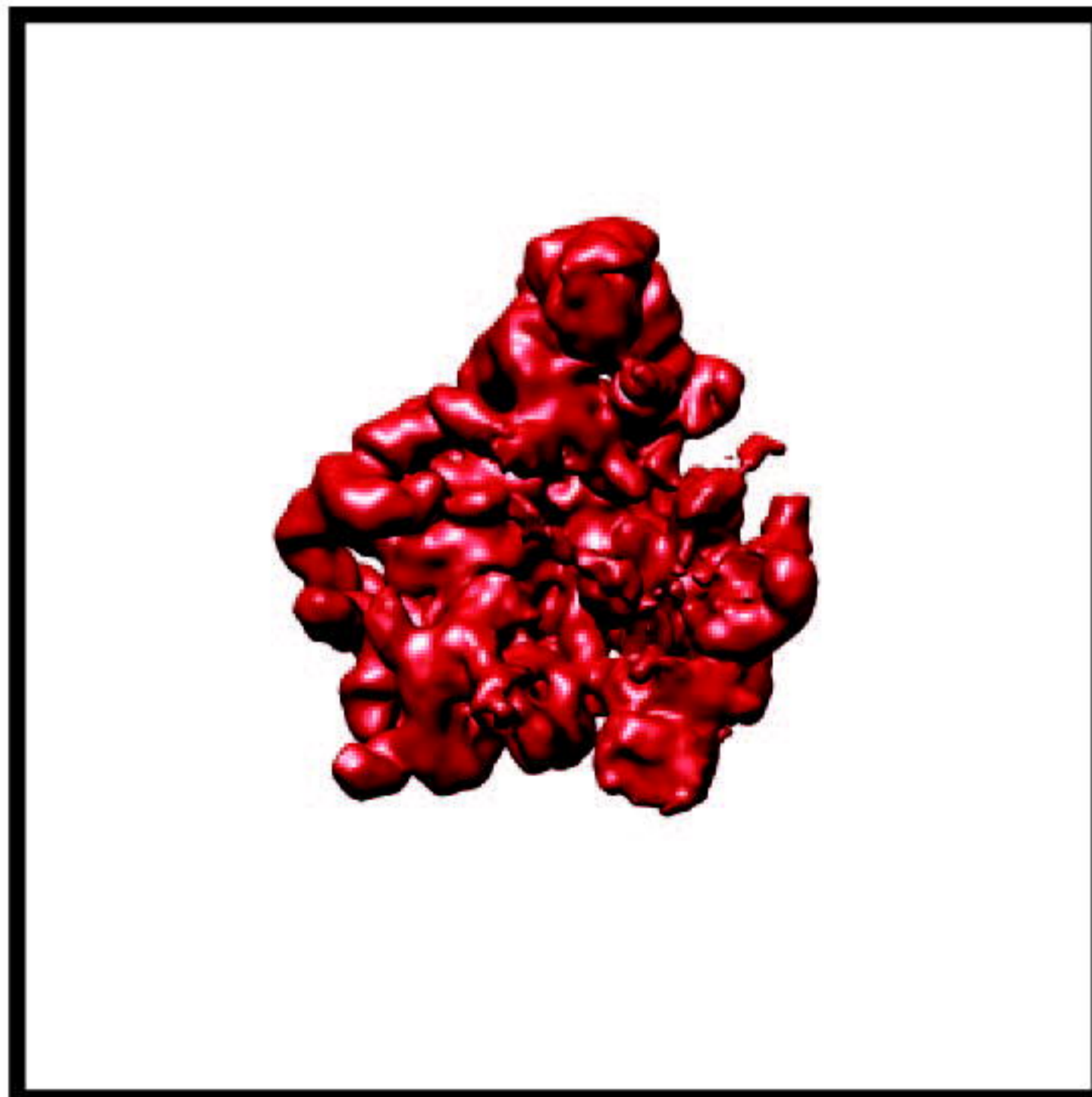
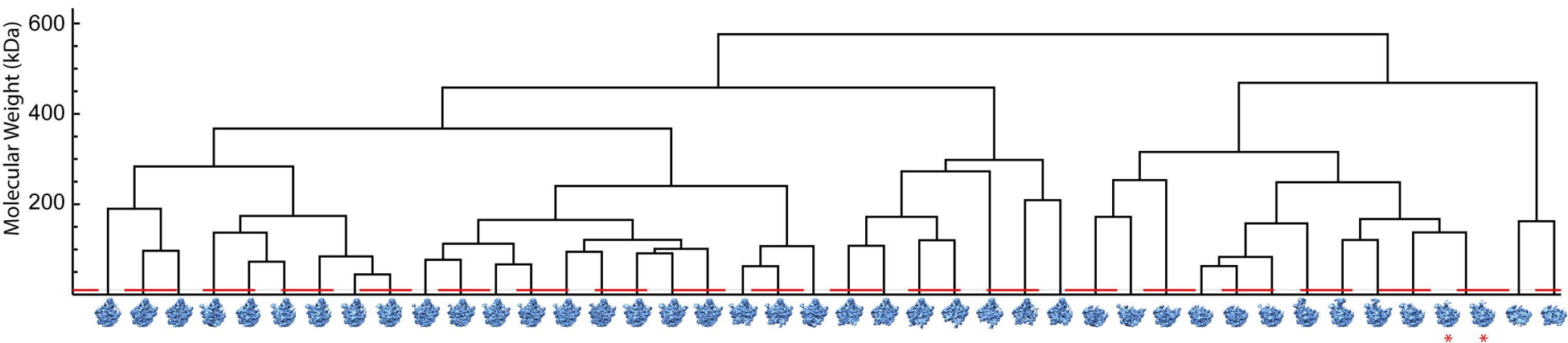


Figure 5

A



B

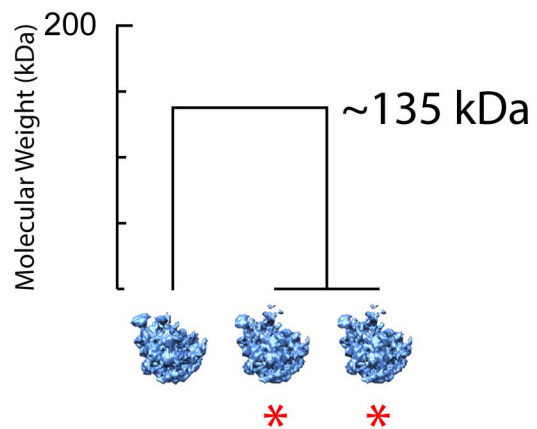
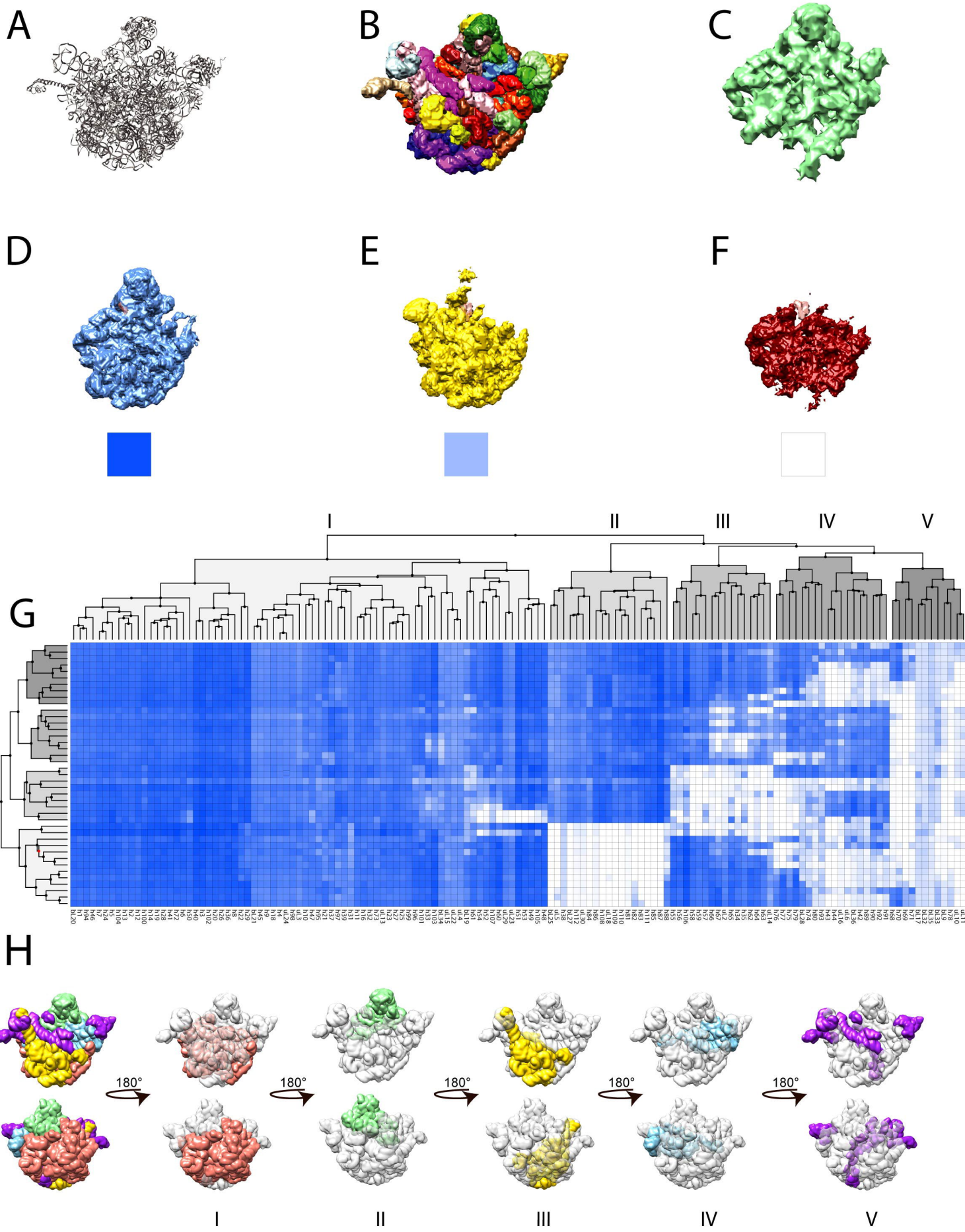
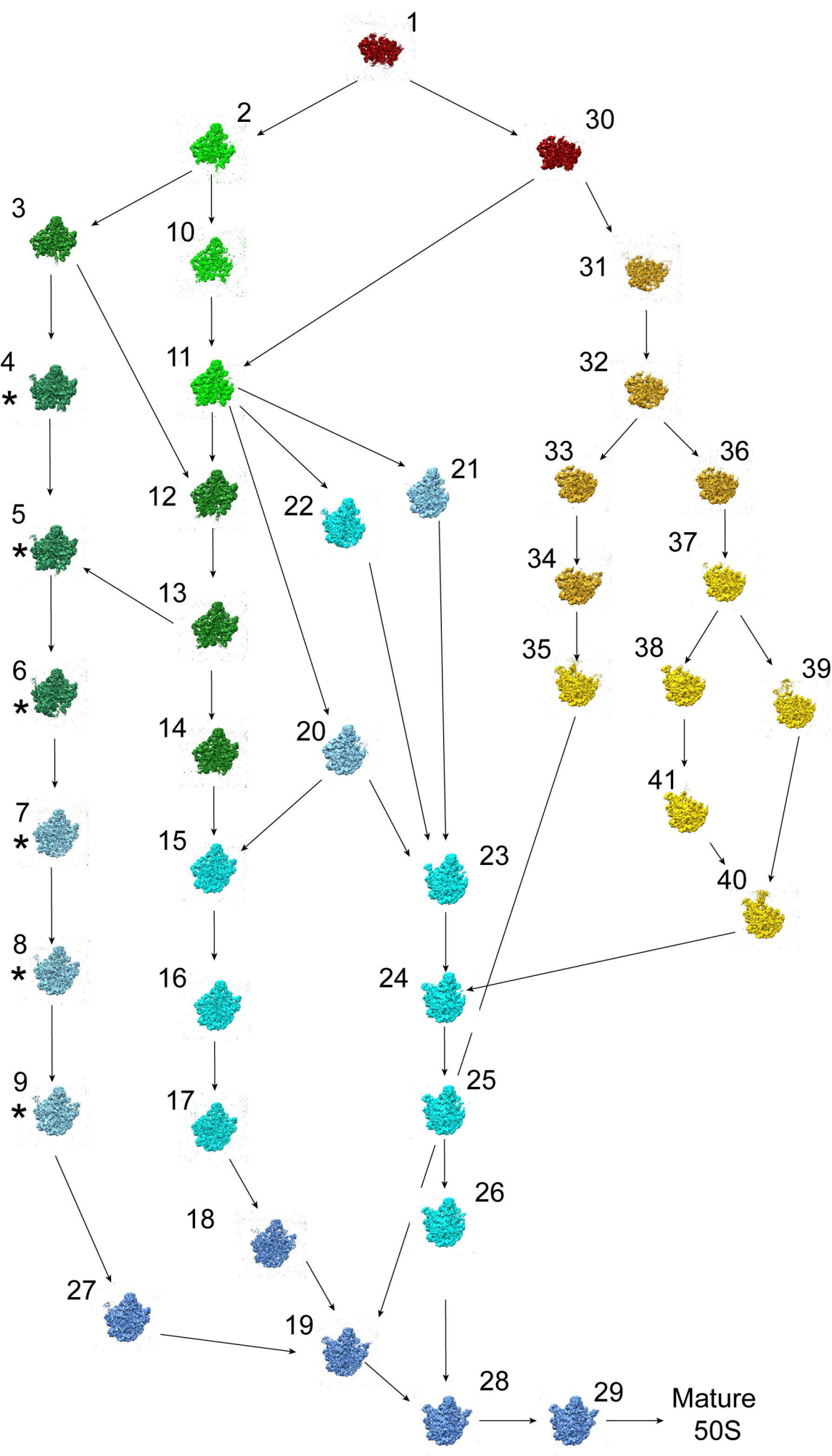


Figure 6





	Threshold/limit	Description	Value used in this paper
limits	r-limit	The minimum resolution necessary for a map	10Å
	v-limit	volume limit: molecular weight difference limit for terminal subdivision	1.5 kDa
Thresholds	low pass filter threshold	used to normalize resolution between maps and to focus on lower-resolution differences between maps	10Å
	binarization threshold	threshold at which maps are binarized; pixel values below this limit are set to 0, values above this limit are set to 1	$3\sigma_{\text{map}}$
	segmentation threshold	defines the volume of dust to be removed from difference maps	1.5 kDa
	difference threshold	defines the lower limit for acceptable differences between maps	10 kDa