

1 Annotating CryoET Volumes: A Machine 2 Learning Challenge

3

4 Author List:

5 Ariana Peck^{1*}, Yue Yu^{1*}, Jonathan Schwartz^{1*}, Anchi Cheng¹, Utz Heinrich Ermel¹, Saugat
6 Kandel¹, Dari Kimanius¹, Elizabeth Montabana¹, Daniel Serwas¹, Hannah Siems¹, Feng Wang²,
7 Zhuowen Zhao¹, Shawn Zheng¹, Matthias Haury¹, David Agard^{1,2}, Clinton Potter¹, Bridget
8 Carragher¹, Kyle Harrington^{1†}, Mohammadreza Paraan^{1†}

9

10 Affiliation:

11 1. Chan Zuckerberg Institute for Advanced Biological Imaging (CZ Imaging Institute), 3400
12 Bridge Parkway, Redwood City CA 94065, USA

13 2. Department of Biochemistry & Biophysics and the Howard Hughes Medical Institute,
14 University of California, San Francisco, San Francisco, CA 94143, United States.

15 * These authors contributed equally to this work.

16 † Corresponding author.

17 Abstract

18 Cryo-electron tomography (cryoET) has emerged as a powerful structural biology tool for
19 understanding protein complexes in their native cellular environments. Presently, 3D volumes of
20 cellular environments can be acquired in the thousands in a few days where each volume
21 provides a rich and complex cellular landscape. Despite numerous innovations, localizing and
22 identifying the vast majority of protein species in these volumes remains prohibitively difficult.
23 Machine learning based methods provide an opportunity to automate the process of labeling
24 and annotating cryoET volumes. Due to current bottlenecks in the annotation process, and a
25 lack of large standardized datasets, training datasets for machine learning algorithms have been
26 scarce. Here, we present a defined “phantom” sample, along with “ground truth” annotations,
27 that will be the basis of a machine learning challenge to bring cryoET and ML experts together
28 and spur creativity to address this annotation problem. We have also set up a cryoET data
29 portal that provides additional diverse sets of annotated 3D volumes from cryoET experts
30 across the world for the machine learning challenge.

31

32 Introduction

33 Much of cell biology is still uncharted territory. The structural biology method of cryo electron
34 tomography (cryoET) is uniquely poised to expand our understanding of cellular function in
35 health and disease¹. Samples are cryopreserved rather than chemically fixed for cryoET, which
36 maintains structural integrity of cellular components, enabling us to peer inside cells and
37 visualize molecular interactions with up to subnanometer resolution. When this technique is
38 applied to lamellae^{2,3}—frozen slices of cells—it provides a detailed 3D view
39 (<https://tinyurl.com/yx2t9wb9>) of the cellular state at the moment of freezing⁴. In eukaryotic cells,
40 this view includes many thousands of different protein complexes (tens of nanometers in size)
41 and organelles (hundreds of nanometers in size)^{5–7}.

42
43 Like tomography methods used in medical imaging such as Computed Tomography (CT) and
44 Magnetic Resonance Imaging (MRI), cryoET also illuminates a region of interest (ROI) from
45 different orientations and collects corresponding 2D projection images (Fig. 1a). These 2D
46 images are then aligned and reconstructed to compute a 3D volume of the ROI (Fig. 1b), the
47 result of which is called a tomogram. This type of data provides both the fine sampling required
48 to obtain near-atomic structures of protein complexes (3–4 Ångstroms)^{8,9}, and also the field of
49 view (600–1000 nm) to visualize each protein complex in its biological context.

50
51 However, there are several challenges to obtaining high resolution maps of protein complexes
52 with cryoET. First, unlike CT and MRI, in cryoET the range of projection views is restricted to
53 less than +/- 60° (Extended Data Fig. 1a) which results in artifacts in the reconstructed volumes,
54 commonly known as the missing wedge artifact¹⁰. Second, because high energy electrons are
55 damaging to biological samples, the number of electrons that can be used is severely limited
56 and thus the resulting images and volumes have very low signal to noise ratios (SNRs) (Fig.
57 1c). As a result, we need to identify and average together thousands of identical entities of
58 different orientations from many 3D volumes in order to produce a high resolution map of a
59 protein complex of interest¹ (Fig. 1d). The process of identifying these individual entities of
60 interest in the noisy 3D volumes is referred to as annotating, picking, or labeling.

61
62 Annotating cryoET tomograms remains a significant bottleneck due to the low SNRs of the
63 projections and volumes, the structural complexity of cellular samples, the diversity and
64 heterogeneity of the molecules of interest, and the large numbers of molecules that are required
65 for high resolution maps. In most cases, annotation is the most time-intensive and laborious part
66 of cryoET data processing since it often relies heavily on manual input. Comprehensive labeling
67 is critical both to obtaining high-resolution structures of individual protein complexes and
68 understanding the organization of large-scale ultra-structures, so new methods to annotate
69 cellular tomograms at scale are urgently needed.

70
71 Machine learning (ML) algorithms are well-suited to overcome this annotation bottleneck^{11–15}. At
72 the time of writing, 15,732 tomograms are publicly available through the recently launched
73 cryoET Data Portal¹⁶ (cryoetdataportal.czscience.com). While machine learning has been
74 leveraged to provide membrane segmentations for all datasets in the portal, labeling particles is

75 a far more difficult task due to their diversity, lower contrast, and crowding. As a result only 5%
76 of tomograms in the portal have molecular annotations. Further, labeling strategies that can be
77 readily adapted to account for data characteristics that vary with acquisition parameters and
78 better capture the intrinsic heterogeneity of molecules would be advantageous compared to
79 traditional approaches. ML methods to date are effective for the specific cases they were
80 developed on, but typically do not generalize sufficiently to meet the diverse needs of the
81 cryoET community.
82

83 Previous machine learning challenges, both in cryoEM and cryoET, have been valuable in
84 benchmarking several algorithms using simulated datasets^{17–19} or limited real-world datasets²⁰.
85 Here, for the first time, we have organized a machine learning challenge
86 (<https://cryoetdataportal.czscience.com/competition>) based on a real-world diverse sample and
87 a large cryoET dataset to spur innovation in this domain. Contestants will be tasked with
88 developing ML algorithms that can robustly perform multi-class particle labeling on hundreds of
89 unseen tomograms after being trained on a limited set of annotated tomograms from the same
90 dataset, as well as any other data — either synthetic or experimental — of the competitors'
91 choice. To encourage generalizability, we have chosen five target particles that have diverse
92 shapes and collectively span nearly an order of magnitude in molecular weight (Fig. 1d). Four
93 target particles, Virus-like Particles²¹ (VLPs), Thyroglobulin (THG)²², Beta-galactosidase²³, and
94 Apoferritin²⁴, were mixed with cellular lysate, which naturally contains the fifth particle, 80S
95 ribosomes²⁵ in abundance (Fig. 2a). The lysate also naturally includes various non-target
96 particles, such as nucleosomes, filaments, proteasomes, and membrane-bound proteins, along
97 with structural elements like membranes, which frequently confuse picking algorithms (Extended
98 Data Fig. 2). We generated reference annotations through an elaborate and rigorous workflow,
99 as described below, that drove the development of several new tools to aid particle picking but
100 also underscored the need for a more streamlined and fully automated solution to tomogram
101 annotation. These “ground truth” labels will be used to score participants’ results and will be
102 released on the CryoET Data Portal as a resource to benchmark future algorithm development
103 after the contest ends. We anticipate that this challenge will deliver novel algorithms to
104 streamline tomogram annotation, provide a benchmark dataset and a standardized pipeline to
105 guide the continuous improvement of ML models, and engage the ML and cryoET communities
106 to jointly tackle challenges facing the field.

107 Creating a “phantom” sample

108 Prior to setting up the machine learning challenge, we carried out a survey of interested parties,
109 which indicated that there was limited enthusiasm for basing a challenge around synthetic
110 datasets nor for the several well annotated existing ribosome datasets. Inspired by the work of
111 Ishengulova et al²⁷, we opted to create a sample from a mixture of cell lysate and known
112 proteins. We refer to this created sample as a “phantom”, based on nomenclature used in the
113 biomedical research community²⁸ to denote objects used as stand-ins for tissues to evaluate
114 systems and methods for imaging. The primary reason for using this phantom rather than
115 cellular material is the inherent complexity and time-consuming nature of preparing and
116 annotating cell-based samples. Working with intact cells introduces a multitude of variables and

117 intricate structures that are difficult to label exhaustively or accurately. This makes it
118 exceptionally challenging to create “ground truth” datasets suitable for training machine learning
119 models and scoring contestants’ results. Therefore, we developed this phantom to mimic the
120 cellular environment to the extent that it provides a relatively large set of annotations for diverse
121 protein structures. For users who are interested in tuning machine learning algorithms to the
122 crowded *in situ* datasets, the cryoET data portal (cryoetdataportal.czscience.com), discussed
123 below, provides multiple annotated datasets of protein complexes *in situ*.

124 Our phantom sample comprises several key components designed to simulate a reduced level
125 of the diversity and complexity of a cellular tomogram (Fig. 2). Firstly, we utilized HEK293T cell
126 lysate^{29,30} that serves as a foundational component, offering a realistic backdrop of cellular
127 material, including common elements, such as ribosomes and membranes (Fig. 2a). To further
128 diversify the sample, we mixed in five commercially available proteins, apo ferritin, thyroglobulin
129 (THG), Beta-galactosidase, Beta-amylase, and Human Serum Albumin (HSA), as well as virus-
130 like particles (VLPs) (see Acknowledgements). These components collectively provide a range
131 of molecular weights, shapes, and therefore complexity (Table 1) for testing and developing
132 machine learning algorithms for protein complex annotation.

133 To produce the phantom sample, we used functionalized electron microscopy grids^{31,32}
134 optimized for cell lysates. The surface of the grid was covered with a layer of graphene oxide
135 (Fig. 2b) and then functionalized with GFP-nanobodies that bind to a GFP-Tag attached to the
136 LAMP1 C-terminus found on lysosome surfaces in HEK293T cells (Extended Data Fig. 1b). This
137 design selectively enriches lysosomes on the grids while other larger organelles such as
138 mitochondria are washed away. However, because of the “gentle” washing step, abundant
139 cellular protein complexes such as ribosomes and nucleosomes remain on the grid. The binding
140 of lysosomes onto the graphene oxide layer creates a spacer effect (Fig. 2c, 2d) that facilitates
141 a consistent sample thickness (~150–250 nm) similar to that typically seen in cellular samples.
142 Sample thickness is a critical factor in cryoET as samples that are too thick have reduced
143 contrast, while samples that are too thin provide very limited cellular volumes. Using our sample
144 preparation method we were able to achieve two features that are otherwise improbable: 1)
145 multiple layers of particles stacked on top of each other (Fig. 2d, 2e, Extended Data Fig. 3)
146 similar to *in situ* samples, 2) exclusion of most of the proteins of interest from the air-water
147 interface (Fig. 2d, Extended Data Fig. 4) which can cause protein denaturation or induce
148 preferred orientations. In these phantom samples, we observed that the air-water interface is
149 consistently populated by very small protein densities (Fig. 2c, Extended Data Fig. 4). We
150 hypothesize these densities to mostly belong to the added HSA and Beta-amylase, as well as
151 small proteins in the lysate. Among the larger proteins, a small subset of ribosomes partially
152 makes contact with the air-water interface. Crucially, preventing preferential orientation allows
153 for a more accurate representation of molecules (Extended Data Fig. 5), which is essential for
154 training robust machine learning models.

155 Our phantom sample captures many of the essential elements of cellular cryoET data, including
156 inherent noise and contrast, a range of shapes and sizes of protein complexes, and the
157 challenges associated with distinguishing species with similar-looking shapes. For example, this
158 similarity arises when the small dimension of a larger protein has similar features to the large

159 dimension of a smaller protein, so that an algorithm interprets the similar views as belonging to
160 the same class of proteins (compare the 2D class views of THG, Beta-galactosidase, and Beta-
161 amylase from Extended Data Fig. 5). The phantom sample also reflects the heterogeneity seen
162 in ribosomes from cells in various states and includes membranes and membrane-bound
163 proteins, although these are not the primary focus of the current challenge. We also observe low
164 copy numbers of additional cellular structures, such as nucleosomes, proteasomes, and
165 luminal proteins which are scattered sparsely through the tomograms, increasing the
166 complexity of the sample (Extended Data Fig. 2).

167 It is important to note however, that our phantom sample does not capture the crowdedness of
168 real cellular environments. While it includes a variety of components, and mimics many aspects
169 of cellular complexity, the density and spatial constraints typical within a living cell are not
170 replicated. While this limitation is acknowledged, the phantom sample nevertheless serves as
171 an initial well annotated dataset for developing, training, and advancing specialized machine
172 learning models for cryoET data analysis, with the anticipation that this will lead to algorithms
173 capable of labeling protein complexes in realistic cellular environments.

174 Establishing “ground truth” labels

175 Ground truth labels, here, are defined as the x,y,z coordinates of the center of each particle in
176 its tomogram with the origin of the coordinate system being the corner of the tomogram.
177 Generating “ground truth” labels required multiple months of work by a large team of people and
178 spurred the development of several new methods for particle picking, visualization, and curation
179 (Fig. 3, Extended Data Fig. 6). This effort included manual annotation, the use of template
180 matching and machine learning algorithms, and the careful curation of the selected particles
181 using a variety of approaches. Below we will first briefly outline the in-house tools that were
182 developed as part of these efforts; these will be described in more depth in forthcoming
183 publications (Table S1). We have already repurposed several of these tools for other projects at
184 our institute so we anticipate these new methods will prove valuable beyond this challenge.
185

186 *DenoisET* implements the *Noise2Noise* algorithm³⁷ for denoising cryoET tomograms. This ML
187 algorithm learns to denoise imaging data from training on paired noisy measurements and is
188 thus suitable for methods like cryoET in which the clean signal is difficult to realistically simulate
189 and cannot be measured. However, because cryoET data are collected as a series of frames of
190 the same underlying field-of-view but with different realizations of the stochastic noise process,
191 training data can be readily generated by reconstructing paired tomograms after splitting the
192 acquired frames into half-sets. Our implementation relies on a similar U-Net architecture as
193 used in Topaz-Denoise³⁸ and leverages the contrast transfer function-deconvolved tomograms
194 produced by AreTomo3²⁶ to improve contrast enhancement. The increased SNR provided by
195 denoising proved critical for our manual annotation efforts and benefited our machine learning
196 labeling workflows as well (Fig. 1C).
197

198 *Copick* is a storage-agnostic and server-less platform designed specifically for cryoET datasets.
199 This package permits efficient access of tomograms and annotations both programmatically, to

200 assist algorithm development, and visually in ChimeraX³⁵ and napari³⁹ for inspection and
201 manual labeling. Specifically, tomograms are stored in the multi-scale OME-Zarr format⁴⁰ to
202 enable rapid and parallel data loading from any filesystem at different resolutions. Annotations
203 are associated with a particle class and other unique identifiers and can be overlaid in an
204 editable format on static copies of reference labels. Collectively, these design choices enabled a
205 large team to manually curate our initial picks in parallel and simplified methods development by
206 providing a unifying framework for data access and storage. We leveraged this framework to
207 easily transfer annotations among all of the in-house tools described here.
208

209 *Slab-picking* was developed as an alternative to 3D template matching to generate initial
210 candidate picks (Fig. 3A). For this approach, tomograms were divided along the z-axis into
211 slabs of uniform thickness and projected along that axis to increase the SNR. These projections
212 were then uploaded into CryoSPARC⁴¹ as mock micrographs. Conventional 2D approaches for
213 particle picking such as blob-picking and template matching were applied to locate candidate
214 particles in each micrograph, and 2D classification was performed to remove false positives by
215 manually deselecting classes judged to be incorrect. The locations were then mapped back to
216 their corresponding tomograms, and the z-height of each particle was refined based on local
217 intensity statistics. Due to time constraints, we were limited to generating candidate picks for all
218 particle types from non-overlapping 300 Å thick slabs. However, we expect that using a sliding
219 window approach to retain particles positioned at the slab boundaries and adjusting the slab
220 height to match the particle size would yield more hits.
221

222 *ArtiaX*³⁶ is a ChimeraX³⁵ plug-in that provides a toolbox for visualizing, selecting, and editing
223 particle picks in tomograms (Fig. 3B). This package was extended in several ways to facilitate
224 and accelerate large-scale manual curation efforts. First, the package was rendered
225 interoperable with the copick framework to enable annotators to rapidly render tomograms at
226 different locations and at distinct voxel spacings, quickly switch between tomograms of interest,
227 and independently curate the same set of candidate picks in parallel. Second, shortcut keys
228 were added to scan candidate picks by recentering the field-of-view on each particle and rapidly
229 delete false positives. Finally, a new feature provides orthoslice views through the tomogram to
230 aid identification and disambiguate particles that are indistinguishable in projection.
231

232 *Copicklive* provides an interactive web viewer to track picking statistics and tomogram curation
233 in real-time (Fig. 3B). These features proved useful for maximizing annotation coverage across
234 the full set of tomograms and reducing duplicate efforts during a week-long manual picking
235 marathon that included 42 individual participants who collectively annotated 147 tomograms and
236 over 29,000 particle annotations. In addition, the interface provides tools to facilitate particle
237 rejection and class reassignment based on 2D projections of subvolumes centered around
238 candidate picks. Because copick is the foundation of this software package, the results of
239 manual curation are automatically saved in a standardized format that can be easily accessed
240 by other software.
241

242 *DeepFindET* is an adaptation of the *DeepFinder* package¹², a CNN-based algorithm to
243 simultaneously label multiple particle types in cellular cryoET data. While initial results from

244 using this package out-of-the-box were promising, the following extensions significantly
245 improved performance on our phantom dataset. First, the model was switched from a U-Net to
246 Residual U-Net architecture for more stable training. Second, additional geometric-based
247 augmentations were added to diversify the training data. Third, the original semi-supervised
248 clustering approach was replaced with a size-based selection scheme to reduce bias and
249 increase confidence in the particle labels during inference. The model was trained on synthetic
250 data containing a mixture of target particles and a limited set of phantom tomograms labeled by
251 expert annotators (Fig. 3C).

252
253 *CellCanvas* is a flexible tool for building geometric models of cellular architecture, with an
254 associated napari plug-in to enable interactive painting-based segmentation and model
255 refinement (Fig. 3C). *CellCanvas*' segmentation capabilities were used to generate an
256 intentionally over-picked set of candidate particles to minimize false negatives. The first step of
257 this pipeline generated a voxel-wise embedding for each tomogram in the phantom dataset
258 using a Swin UNETR model, which was pre-trained on medical imaging data⁴² and fine-tuned
259 on four synthetic tomograms generated using PolNet⁴³. Clusters in embedding space were
260 transformed into segmentation masks by interactively training an XGBoost classifier on manual
261 annotations. These segmentation masks were in turn converted to particle picks. While we
262 found the combination of pre-computed embeddings and quick-to-train interactive models highly
263 effective for labeling the phantom data, *CellCanvas* was intentionally built with a plug-and-play
264 design that enables replacing the embedding model and classifier with other models for
265 increased flexibility.

266
267 *Minislab curation* was developed to facilitate particle curation which is the filtering of picked
268 particles. In this approach, subvolumes centered on individual particles were extracted from the
269 tomograms, and the intensity was integrated along the z-axis to generate per-particle 2D
270 projections (Fig. 3D). These “minislabs” were then tiled to generate mock micrographs for
271 further processing in CryoSPARC⁴¹. The particles were cleaned up and curated by various
272 combinations of performing manual picking, 2D classification, running *ab initio* reconstruction
273 with multiple classes, and applying Topaz¹⁴, a CNN-based particle picker that was retrained for
274 each particle class. The selected particle projections were mapped back to their positions in the
275 tomograms. While minislab curation cannot distinguish between particles that appear similar in
276 projection, curation of 2D projections was more efficient and benefited from the increased
277 signal-to-noise compared to 3D subvolumes.

278
279 The tools described above, and other published software packages were stitched together into
280 several workflows to generate “ground truth” labels (Fig. 3, Extended Data Fig. 6). Initial picks
281 were generated either using PyTom^{33,34} 3D template matching (VLP, ribosome) on undenoised
282 tomograms or 2D slab-picking (apo ferritin, Beta-galactosidase, thyroglobulin) on denoised
283 tomograms. Preliminary picks for a subset of tomograms were manually curated in ChimeraX³⁵,
284 leveraging new features in ArtiaX³⁶ for facile annotation and copick for rapid file access and
285 exchange. These picks were curated by expert annotators for the next step of model training. To
286 generate a more refined set of picks across the full phantom dataset with less false negatives or
287 missed particles, DeepFindET and the *CellCanvas* classifier were trained on a combination of

288 synthetic tomograms generated by PolNet⁴³ and real data curated by the expert annotators. The
289 PolNet training set contained the five particles of interest in addition to membranes and a sixth
290 particle, Beta-amylase. Minislabs for the predicted picks were extracted from denoised
291 tomograms for final rounds of curation. For VLP and apoferritin, curation was achieved by
292 applying a Topaz model trained separately for each particle. For thyroglobulin and Beta-
293 galactosidase, 2D classification was performed and the retained picks from both CellCanvas
294 and DeepFindET were merged into a single set. After removing duplicates, the picks underwent
295 a round of manual curation followed by multi-class *ab initio* reconstruction to reject any
296 remaining false positives. For the ribosome, both Topaz and iterative rounds of 2D classification
297 were used to eliminate contamination. Given the intrinsic heterogeneity of ribosomes, we
298 performed final rounds of 2D classification and multi-class *ab initio* reconstruction to provide a
299 confidence score for each retained particle. In the case of particles curated by Topaz, any true
300 positives rejected by the model were rescued through manual curation to minimize false
301 negatives. This long, complex, and convoluted approach required to generate accurate
302 annotations highlights the need for more streamlined algorithms that generalize well across
303 datasets and diverse target particles and thus the urgent need for this ML challenge.
304

305 While we refer to our final curated set as “ground truth”, we acknowledge the following caveats.
306 Based on our rigorous curation efforts, we are confident that the less challenging particle labels
307 contain few false positives but are less confident about the more challenging particles. We also
308 acknowledge that our ground truth labels very likely missed some true positives. These caveats
309 are due to the fact that even a well-trained expert finds it exceptionally difficult to be certain in
310 labeling the smaller particles and it would take years for experts to fully annotate this large
311 dataset. We also note that we opted for a more inclusive set of particles than would be included
312 in the typical downstream processing task of subtomogram averaging. Thus, our final labeled
313 particles are conformationally heterogeneous and include partial particles at the tomogram
314 boundaries if their identities seem clear. In the case of the VLPs, we retained not only the
315 abundant 29 nm diameter icosahedral form but also other spherical particles and tubular
316 species. Our intention in making these choices was to focus the challenge on the specific task
317 of labeling without constraining the labels to the most homogeneous subset of particles in the
318 data. The final ground truth set includes thousands of picks for each of the THG, Beta-
319 galactosidase, and VLP particles, and tens of thousands of picks for each of the 80S ribosome
320 and apoferritin particles (Table 1, Extended Data Fig. 5).

321 The Machine Learning Challenge

322 The cryoET field has benefited from an increase in machine learning contributions, but particle
323 picking, a problem that is well-suited for machine learning methods, remains a major open
324 problem. We are hosting an ML challenge (<https://cryoetdataportal.czscience.com/competition>)
325 to push the limits of particle picking across particle sizes and draw attention from the broader
326 ML community. We deliberately selected a small training subset of the full dataset to
327 recapitulate the common situation of a cryoET researcher who can only afford enough effort to
328 annotate a handful of tomograms, yet needs annotations for hundreds or thousands of
329 tomograms for high-resolution subtomogram averaging. We have designed evaluation metrics

330 that span a range of particle sizes, incentivizing models that can perform well for small particles.
331 We note that due to the challenges of annotating particles, even by well-trained human experts,
332 there is the possibility that a fraction of the particles have not been labeled in the ground truth.
333 These false negatives will be explored after the competition is over. We will use a variety of 2D
334 and 3D classification methods to assess the particles provided by the top 10 submissions to
335 determine if they include a larger fraction of true positives than the ground truth we provide.
336 This analysis will be shared with the community and published as part of a larger assessment of
337 the ML challenge results, the winning algorithms, and the lessons learned.
338

339 *Data description:* The 3D volumes provided to the participants on Kaggle- for training and
340 testing- are denoised tomograms (Table 2). The three other tomogram types (weighted back
341 projection, CTF-corrected, Isonet-corrected) are available on Kaggle only for training and not
342 testing. All the tomogram types are available to the participants on the cryoET data portal (as
343 well as raw data such as frames and tilt series). The ground truth labels for these tomograms
344 are x,y,z coordinates of the center of each labeled particle. The origin for this coordinate system
345 is the lower bottom left corner of each tomogram (which would be increasing positively along the
346 x, y, and z axes of the 3D image array). These ground truth labels belong to 6 distinct classes
347 which are evaluated individually. Therefore, participants are required to submit their labels in the
348 same x,y,z coordinate system per class.
349

350 *Evaluation metric:* The evaluation metric for the cryoET particle picking challenge uses a
351 distance threshold-based point matching approach, where matches are determined based on a
352 multiple of the particle radius. The metric employs microaveraging, averaging matches across
353 all picks rather than per tomogram. To mitigate the penalty for false positives, an F-score with a
354 beta value of 4 is utilized (e.g. 1 false negative is penalized the same as 16 false positives),
355 emphasizing precision and reducing the penalty for false positives given the annotation
356 uncertainty for small particles. Additionally, the metric incorporates weighted scoring for different
357 particle groups, assigning higher importance to hard-to-pick particles (Fig. 4). A weight of 1 is
358 assigned to the “easy” to pick particles: apoferritin, ribosome, and virus-like particles (VLPs),
359 while a weight of 2 is assigned to the more difficult particles to pick: thyroglobulin (THG) and
360 Beta-galactosidase. A weight of 0 is assigned to Beta-amylase, because it’s not part of the
361 challenge, although annotations are provided for training purposes.
362

363 *Dataset split:* We curated a selection of 492 good-quality tomograms (Extended Data Fig. 7)
364 from the phantom dataset and divided them into three subsets: training, testing, and validation.
365 To ensure that these subsets had similar distributions across all particle species, we employed
366 the Kolmogorov-Smirnov test to assess the similarity between distributions. The 492 tomograms
367 were then grouped into 16 bins based on this analysis, and tomograms randomly sampled from
368 these bins were used to construct the three subsets. We experimented with various data split
369 ratios and observed that DeepFindET was able to perform well with smaller training sets, which
370 is the most common situation in CryoET. Consequently, we opted to submit only seven
371 tomograms to be used for training, 25% of the remaining tomograms are used for public
372 evaluation and 75% are used for private evaluation and final scoring of the submitted models
373

374 *Establishing metrics parameters with synthetic datasets:* To ensure that the F-beta scores
375 effectively distinguish performance differences among models (especially models that will
376 perform better than DeepFindET), we determined the metric's beta value using a set of results
377 with known ranking quality. This set was generated by mixing varying percentages (in 10%
378 increments) of picks from a fine-tuned DeepFindET model and ground truth data. By comparing
379 their F-beta scores, we observed that datasets containing 80% or more ground-truth results
380 yielded mock datasets that score better than DeepFindET. Based on this insight, we added an
381 additional 20 datasets ranging from 80% to 100% (in 2% increments) ground truth with the
382 remaining picks coming from DeepFindET and used them to compare the impact of different
383 beta values. For this comparison we aimed to find the smallest beta value that ensured a
384 consistent top-5 ranking. Through this analysis, we determined that setting beta=4 produced a
385 more consistent ranking.
386

387 *Example notebooks:* To reduce the onboarding time for competitors, an extensive set of
388 example notebooks is being provided. These include: DeepFindET, TomoTwin, and 3D U-Net
389 models. For each model type (except TomoTwin which uses a pretrained model), training and
390 inference notebooks are provided. These notebooks also leverage the copick library for
391 handling cryoET datasets, for which we provide PyTorch Datasets and utility functions to
392 simplify the creation of data loaders, metadata tracking, and model performance analysis.

393 Enhancing training infrastructure and data using the 394 cryoET data portal

395 All training data described here (7 tomograms) as well as additional published training data are
396 available on the CZ CryoET Data Portal (CZCDP¹⁶, cryoetdataportal.czscience.com) under
397 dataset IDs 10440 (experimental data) and 10441 (PolNet⁴³-simulations). The CZ cryoET Data
398 Portal is an open platform specifically designed to facilitate algorithm development, including by
399 researchers outside the cryoET field. It expands on the data archiving and sharing efforts made
400 by the Electron Microscopy Public Image Archive (EMPIAR)⁴⁴, and the Electron Microscopy
401 Data Bank⁴⁵, by providing a large collection of curated and standardized datasets and
402 annotations under public domain (Creative Commons 0) license, as well as a growing set of
403 open source tools geared towards machine learning development and cryoET visualization.
404

405 Independent of domain expertise, ML Challenge participants can familiarize themselves with the
406 training data, as well as other available cryoET datasets, by browsing the Portal's web interface
407 and visualizing the provided ground truth annotations in their web browser using the integrated
408 Neuroglancer-based viewer. Outside of the browser, datasets can be queried using a GraphQL-
409 API, accessed using a python-based API client and visualized using a plugin to the ND-image
410 viewer napari³⁹. In contrast to existing archives, all datasets are accessible with a consistent
411 layout and metadata schema in public cloud storage, and directly link raw data, tilt series,
412 tomograms and annotations. All processed image data are stored in the cloud-ready next
413 generation file format OME-Zarr⁴⁰, allowing participants to train particle picking networks on
414 datasets exceeding available local storage by streaming tomograms whole or in part, and

415 removing dependencies on domain-specific image formats (such as MRC). PyTorch⁴⁶ and
416 Tensorflow⁴⁷ infrastructure relying on these features is available to participants through the
417 direct integration of CZCDP datasets with the copick API.
418
419 We expect that the considerable volume (15,732 tomograms) and diversity (44 species) of the
420 data already hosted on the CryoET Data Portal will lend itself to the development of robust ML
421 algorithms that generalize well across samples. In addition, the growing collection of datasets
422 should serve as further motivation to ML Challenge participants, as any tools and algorithms
423 relying on CZCDP infrastructure and/or copick will be immediately applicable to the entire
424 current collection, as well as any future data deposited to the portal, with little or no need for
425 modification.
426
427 Taking into account the available datasets and infrastructure, we envision direct ways through
428 which the cryoET Data Portal could be leveraged to improve model training for this challenge.
429 For example, the synthetic data used to train CellCanvas and DeepFindET have been made
430 available (portal dataset ID: 10441), and can serve as convenient test sets for models and
431 infrastructure. In addition, although these synthetic tomograms lack the complexity of
432 experimental data, we observed that training on them significantly improved prediction accuracy
433 across all particle classes for both of the models we used. The deposited annotations may also
434 prove a valuable resource for model refinement. In particular, the 200,843 instances of non-
435 mitochondrial ribosome labels across 447 non-phantom tomograms could supplement the
436 limited training tomograms from the phantom dataset. Although this is the only target particle
437 that has been annotated in other Portal datasets, including molecular annotations for non-target
438 particles could potentially improve a model's ability to distinguish among different molecular
439 species. In a similar vein, the membrane segmentations available for all Portal datasets could
440 be used as negative training examples to prevent algorithms from mistaking these high contrast
441 features for target particles. We anticipate that the specific ways in which contestants
442 incorporate data and annotations from the CryoET Data Portal into their workflows will be
443 informative for future algorithm development.
444
445 A final distinguishing feature of the CZ cryoET data portal is the evolving nature of its corpus.
446 Annotations generated as part of the Challenge will be added to the cryoET data portal with
447 attribution to the authors and direct links to their tools and methods. This provides a simple path
448 for participants to present their results interactively to a broad audience and allows for easy
449 reproducibility and comparison of any methods developed as part of the Challenge.

450 Discussion

451 Challenges like the one described here have catalyzed the development of state-of-the-art
452 methodologies in structural biology. A particularly instructive example is the biennial Critical
453 Assessment of protein Structure Prediction (CASP) competition that benchmarks protein
454 structure modeling algorithms against experimentally-determined structures⁴⁹. Compared to the
455 incremental progress achieved during CASP's first decades, the introduction of DeepMind's
456 AlphaFold in 2018 led to a major breakthrough in prediction accuracy⁵⁰. This leap in

457 performance was so profound that CASP has since shifted its focus to more challenging protein
458 modeling problems, while AlphaFold has increasingly been integrated into other structural
459 biology workflows^{51,52}. This advance critically depended not only on innovations in deep learning
460 research but also on the public availability of a rich corpus of experimental results^{50,53}.

461
462 Until recently, the conditions to hold a Machine Learning Challenge had not been met in the field
463 of cryoET. However, several developments have converged to make such a competition not
464 only possible but potentially game-changing for the field. As with protein structure modeling, a
465 wealth of experimental data is now available across multiple public data banks. EMPIAR⁴⁴ and
466 EMDB⁴⁵ have led the charge to make cryoET data openly accessible. Expanding on these
467 efforts, we launched the CryoET Data Portal to prioritize data standardization, tomogram
468 annotation and segmentation, including by non-cryoET practitioners so that domain-specific
469 knowledge is not a barrier to entry. The recent acceleration of data acquisition and processing is
470 poised to vastly expand the quantity of public data and it has also enabled us to generate a
471 benchmark dataset that includes hundreds of high-quality tomograms. Importantly, this phantom
472 dataset contains artifacts characteristic of experimental tomograms and is large enough to
473 robustly score annotation algorithms in blind predictions against withheld data. As the existing
474 labeling approaches tried in this project did not work out-of-the-box on our phantom tomograms,
475 rigorously annotating the full dataset required months of effort and stitching together in-house
476 tools and existing software into a convoluted workflow that underscored the need for a more
477 streamlined solution. The result of this effort is high-quality annotations across nearly five
478 hundred experimental tomograms for six distinct particle classes whose range of shapes and
479 sizes will identify deep learning algorithms most likely to generalize to the diverse targets of
480 interest to the cellular biology community.

481
482 We expect that our Machine Learning Challenge will serve as the foundation for a series of
483 contests designed to spur innovation in cryoET methods development, including from experts
484 outside the field, and provide a standardized platform to critically evaluate state-of-the-art
485 annotation algorithms. The results from these contests will in turn inform the design of future
486 benchmark datasets that best showcase the capabilities and identify the limitations of current
487 methods, offering the field a valuable resource to continually track its progress. Similar to CASP,
488 advances will enable us to expand the scope to more difficult challenges such as annotation in
489 crowded cellular environments and labeling molecular features along high-contrast membranes.
490 We hope that these challenges over time will deliver transformative ML methods that solve the
491 annotation bottleneck in cryoET and can be readily scaled across all available data to deepen
492 our insights into how molecules are structured and organized in their native context of the cell.
493

494 **Tables**

Sample	MW	Symmetry/copies	Average counts per tomogram	ID
Ribosome (80S; human)	4.3 Md	Monomer 1 x	80	EMD-3883
VLP (PP7)	3.4 Md	Icosahedral 60 x	6	EMD-41917
Thyroglobulin (bovine)	660 Kd	Homodimer 4 x	13	EMD-24181
Beta-galactosidase	540 Kd	D2 4 x	6	EMD-0153
Apo ferritin	450 Kd	Octahedral 24 x	47	EMD-41923
Beta-amylase (Sweet potato)	268 Kd	D2 4 x	5	EMD-30405
Albumin (Human)	66 Kd	Monomer 1 x	NA	EMD-43090

495
496
497
498
499
500

Table 1. Purified proteins (from different sources) present in the phantom dataset. All proteins except for Human Serum Albumin (HSA) and Beta-amylase are targets in the machine learning challenge. The listed EMD entries are for the maps used in Figures 1 and 2, and some initial picking with PyTom template matching.

Tomogram Type	Reconstruction workflow	Description	Contrast
Weighted Back Projection	AreTomo3	Real space WBP	Low
CTF-Deconvolved	AreTomo3	Local CTF correction of tilt series+ real space WBP	Medium
Denoised	AreTomo3 +DenoisET	Local CTF correction of tilt series+ real space WBP+tomogram denoising	High

Isonet-Corrected	AreTomo3 +DenoisET +Isonet	Local CTF correction of tilt series+ real space WBP+tomogram denoising+missing wedge estimation and correction	High
------------------	----------------------------------	--	------

501

502 **Table 2. Different tomogram types available in the challenge.** All of these tomograms have
503 the same tilt series alignment, therefore, the particle coordinates are the same across all types.

504

Methods

505

Sample preparation

506 The phantom dataset was generated from a combination of cell lysates, commercially available
507 purified protein, and purified protein from collaborators. The cell lysates were from a HEK293T
508 cell line generated at Chan Zuckerberg BioHub, San Francisco by Manuel Leonneti's group with
509 a knock-in GFP-Tag on the C-terminus of the lysosomal house-keeping protein LAMP1. The
510 homogenization protocol was developed at our collaborator's lab³⁰ and adapted by us for
511 cryoET sample preparation. Briefly, the cells were lysed using a hypotonic homogenization
512 buffer (25mM Tris*HCl pH 7.5, 50mM sucrose, 0.2mM EGTA, 0.5mM MgCl₂) and shearing
513 forces generated using a 23G syringe. To protect the organelle membranes the lysate was
514 immediately mixed with a sucrose buffer (2.5M sucrose, 0.2mM EGTA, 0.5mM MgCl₂) to re-
515 equilibrate the homogenates to an isotonic osmolarity. The nuclear fraction was separated by
516 centrifuging the lysate at 1000 x g for 10 minutes. The lysosomes with the LAMP1-GFP tag from
517 this lysate were then purified on-grid using functionalized electron microscopy grids^{31,32}. This
518 technology was adapted for organelles at the CZ Imaging Institute and University of California,
519 San Francisco in a joint collaboration with David Agard's lab. The lysosomes were captured on-
520 grid using GFP nanobodies. These GFP nanobodies are attached to the maleimide groups
521 covering the surface of the grid using a kck linker. All lysis steps were at 4°C and the lysate was
522 kept on ice until plunge freezing.

523

524 Plunge freezing was done using a Leica GP2, a Whitman #1 blotting paper, and liquid ethane at
525 -180°C as a cryogen for fast freezing. The functionalized grid was loaded into the chamber
526 which was set at 4°C and 95% humidity. The lysate was pipetted up and down (20 strokes) in
527 10 ul volumes in 10 repeated rounds, resulting in an overall application of 100 ul of lysate to the
528 grid surface. The grid surface was then washed with PBS two times by pipetting. Before adding
529 the purified protein species most of the buffer volume left on the grid was pipetted away leaving
530 the grid and the lysosomes attached to it just hydrated enough. This was done so that the
531 dilution factor for the 6 purified protein species would be limited. Each of the purified proteins
532 were added one at a time in 1 µl volumes as follows: 1) THG (from bovine thyroid, Sigma-
533 Aldrich T9145) at a concentration of ~17.8 mg/mL (A280=19.35), 2) Apoferritin (from equine
534 spleen, Sigma-Aldrich 178440) at a concentration of ~5 mg/mL (A280= 4.73), 3) Beta-
535 galactosidase (from *E. Coli*, Sigma-Aldrich G5635) at a concentration of ~6 mg/mL
536 (A280=13.79), 4) Beta-amylase (from sweet potato, Sigma-Aldrich A8781) at a concentration of
537 ~5 mg/ml (A280=4.75), 5) Human Serum Albumin (Sigma-Aldrich A3782) at a concentration of

538 ~50 mg/mL (A₂₈₀= 20.15), 6) Virus-like Particles (from collaborators at NYSBC) at a
539 concentration of ~7.5 mg/mL (A₂₈₀=26.91). The target concentration for most of the species
540 was aimed at 5 μ M after the 6-fold dilution due to mixing. These are the per-species errors
541 between the 5 μ M goal and the final values as a result of challenges in volume handling at high
542 concentrations: 1) THG: 15%, 2) Apoferritin: 33%, 3) Beta-galactosidase: 27%, 4) Beta-
543 amylase: 16%. HSA concentration was kept high because it served as a background protein.
544 VLP concentration could not be pushed further because of difficulties with volume handling at
545 high concentrations. That said, all proteins are present in sufficient amounts in all the
546 tomograms. After all the purified proteins were added to the grid, back-side blotting was done
547 for 6 seconds and the grid was plunged into liquid ethane and then stored in liquid nitrogen.

548 Data collection

549 All 1,089 tilt series were collected on one grid on a Krios G4 equipped with an X-FEG electron
550 gun, a Falcon 4i direct electron detector, and the SelectrisX energy filter. The pixel size was set
551 to 1.54 Å/pix, and the total dose to 62.93 e⁻/Å² linearly spread over 31 tilt images spanning a
552 range of -45° to +45° in 3° increments. The software used for data collection was TFS Tomo 5.
553 Utilizing the beam-image shift data collection feature, 9 targets were imaged at each stage
554 position. The movies were saved in the EER format.
555

556 For a rapid quality check of the sample, 2D data were also collected on the same grid for single
557 particle analysis. The virus-like particles were refined to a resolution of 3.94 Å. Briefly, 124 movie
558 frames were collected and all the down-stream processing was done in CryoSparc⁴¹. To
559 generate templates for VLPs, 12 micrographs were manually picked. Template matching
560 resulted in 3,844 particles which were filtered down to 606 particles after 2D classification. Ab-
561 initio model generation and the homogenous refinement resulted in a 3.93 Å map of an
562 icosahedral VLP (Extended Data Fig. 8).

563 Data processing

564 Motion correction, tilt-series alignment, and tomogram reconstruction were all performed using
565 AreTomo3²⁶. Specifically, raw frames were partitioned into non-overlapping groups of 2000
566 frames each, and batches of 10 frames within each group were integrated to generate rendered
567 frames. Every two and four rendered frames were summed for measuring global and local
568 motion, respectively. Motions measured on group sums were then interpolated to individual
569 rendered frames for more accurate correction of rapid motion. For local motion estimates, these
570 integrated frames were further subdivided into 5x5 patches. Both gain and motion corrected
571 frames are summed to generate corrected tilt images. CTF parameters were estimated for each
572 tilt image. These tilt images were then aligned into tilt-series, with global alignment followed by
573 4x4 patch-based local alignment. Tomograms were reconstructed using weighted-back
574 projection, either with (for denoising) or without (for 3D template matching) applying a local CTF
575 deconvolution and correction to the tilt-series. Even and odd pairs of CTF-deconvolved
576 tomograms were produced to enable denoising. Paired tomograms were generated by splitting
577 the motion-corrected frames into even and odd sets and then applying the alignment

578 parameters determined for the full tilt-series to generate paired tomograms. Tomograms were
579 denoised using denoisET, an in-house implementation of *Noise2Noise*³⁷. This algorithm relies
580 on a U-Net architecture composed of five downsampling and upsampling blocks each, with a
581 kernel size of 3 pixels applied during each convolutional layer. This model was trained for 10
582 epochs on subvolumes extracted from 43 pairs of even and odd CTF-deconvolved tomograms
583 before being applied to full tomograms to denoise the full dataset.

584 Generating ground truth

585 Two approaches were taken to generate initial picks. For the 80S ribosome and VLP,
586 PyTom's^{33,34} 3D template matching algorithm was applied to the CTF-uncorrected and
587 undenoised tomograms with a 10 Å pixel size. Templates were generated internally in PyTom
588 from published maps (EMDB IDs: 3883 and 41917). For THG, apoferritin, and Beta-
589 galactosidase, slab-picking was used to identify candidate particles. Specifically, denoised
590 tomograms with a 5 Å pixel size were divided along the z-axis into non-overlapping 300 Å thick
591 slabs, which were then projected along that axis to generate mock micrographs for further
592 processing in CryoSPARC⁴¹. Blob picking was performed across a range of particle diameters,
593 followed by iterative rounds of 2D classification to separate the different particle classes from
594 each other and filter contamination based on visual inspection. Candidate particles were
595 mapped back to their tomograms, and the depth of each particle was refined by locating the z-
596 coordinate with the maximum integrated intensity in a subvolume centered on the particle.
597 Particle coordinates from PyTom and slab-picking were stored in copick, and candidate picks for
598 147 of the 492 tomograms were manually curated during a week-long pickathon marathon using
599 the ChimeraX³⁵ copick plug-in and copicklive to track curation statistics.
600 A more refined and expansive set of picks was generated using two deep learning approaches.
601 One of these, DeepFindET (an adaptation of DeepFinder¹²), uses a Residual U-Net architecture
602 with three 3D convolutional layers and a receptive field size of 680 Å to predict segmentation
603 masks from annotated tomograms. The other, CellCanvas, involves a multi-step pipeline to
604 segment tomograms. In the first step, a Swin UNETR model pre-trained on computed
605 tomography data⁴² generates a voxel-wise embedding for each tomogram. In the second step,
606 an interactively trained XGBoost classifier transforms clusters in embedding space into multi-
607 class segmentation masks so that each voxel of the tomogram is assigned a probability for each
608 class label. For initial training data, synthetic tomograms were generated in PolNet⁴³ with a
609 thickness of 180 nm and the five target particles, Beta-amylase, and membranes randomly
610 distributed throughout the sample volume. An additive Gaussian noise model was applied, but
611 neither the contrast transfer function or sample motion were modeled. The DeepFindET model
612 was trained on 24 of these synthetic tomograms and then fine-tuned on five manually annotated
613 tomograms from the phantom dataset and curated picks from the initial round of template
614 matching and slab picking. For CellCanvas, the pre-trained Swin UNETR model was fine-tuned
615 on four synthetic tomograms, while the classifier model was iteratively trained on manually
616 annotated phantom tomograms. The segmentation masks predicted by DeepFindET were
617 converted to individual picks using a size-based threshold, with the cut-off set to two thirds of
618 each particle's maximum dimension. In the case of CellCanvas, segmentation masks were
619 grouped by the Watershed algorithm into similar regions, and the centroid of each region was

620 labeled as the particle with the highest local probability density. For subsequent curation efforts,
621 predictions for both models were pooled for THG and Beta-galactosidase, while the predictions
622 for VLP, 80S ribosome, and apoferritin came exclusively from DeepFindET. We chose not to
623 use the predictions from CellCanvas for these particles because the large number of candidate
624 particles and particle miscentering proved computationally prohibitive to our downstream
625 curation pipeline. Beta-amylase picks came from 2D classes of it that were generated during the
626 classification of THG and Beta-galactosidase picks.
627 For the final rounds of curation, subvolumes centered on each particle were extracted from the
628 denoised tomograms and projected along the z-axis to generate 2D per-particle projections.
629 These “minislabs” were then tiled in 2D and loaded into CryoSPARC⁴¹ as mock micrographs,
630 each containing 240 candidate particles. A single particle pick was positioned in the center of
631 each tile. In the case of the VLPs, 80S ribosome, and apoferritin, we manually annotated 2, 14,
632 and 19 micrographs per class respectively and used these positive labels to train a Topaz
633 model¹⁴ for each particle. The trained Topaz models were then applied to the remaining
634 micrographs. To increase the accuracy of our labels, we inspected minislabs of both the
635 rejected particles to rescue any true positives and the accepted particles to eliminate any false
636 positives during a final round of manual curation. In the case of the 80S ribosome, 2D
637 classification followed by *ab initio* reconstruction with two classes was performed to assign a
638 confidence score to each particle. Particles belonging to the manually-selected high-quality
639 classes and retained in the correct reconstruction class were given a score of 1. Particles that
640 segregated into the low-quality classes from 2D classification were assigned a score of 0, while
641 the remaining particles that were selected during 2D classification but did not contribute to the
642 correct reconstruction were assigned a score of 0.5. Only the score=1 ribosomes were retained
643 for the challenge. In the case of THG and Beta-galactosidase, minislab micrographs were
644 separately generated from the DeepFindET and CellCanvas picks, and 2D classification was
645 performed on each set. After selecting high-quality classes based on visual inspection, the
646 particles from each ML model were merged into a single set. Duplicates were removed based
647 on a distance threshold of 185 Å and 125 Å for THG and Beta-galactosidase, respectively. The
648 merged sets were manually curated, followed by *ab initio* reconstruction with two classes.
649 Particles that contributed to the correct reconstruction class were retained. Finally, we
650 generated minislabs from any positions that belonged to more than one particle class and
651 manually assigned one label based on visual inspection to remove these inter-class duplicates.
652 The final particle lists from these curation efforts were considered the ground truth annotations.

653 References

- 654 1. Young, L. N. & Villa, E. Bringing Structure to Cell Biology with Cryo-Electron Tomography.
655 *Annu. Rev. Biophys.* **52**, 573–595 (2023).
- 656 2. Rigort, A. & Plitzko, J. M. Cryo-focused-ion-beam applications in structural biology. *Arch.*
657 *Biochem. Biophys.* **581**, 122–130 (2015).
- 658 3. Villa, E., Schaffer, M., Plitzko, J. M. & Baumeister, W. Opening windows into the cell:

- 659 focused-ion-beam milling for cryo-electron tomography. *Curr. Opin. Struct. Biol.* **23**, 771–777
660 (2013).
- 661 4. Khavnekar, S. *et al.* Towards the Visual Proteomics of *C. reinhardtii* using High-throughput
662 Collaborative *in situ* Cryo-ET. *Microsc. Microanal.* **29**, 961–963 (2023).
- 663 5. Mahamid, J. *et al.* Visualizing the molecular sociology at the HeLa cell nuclear periphery.
664 *Science* **351**, 969–972 (2016).
- 665 6. Barad, B. A., Medina, M., Fuentes, D., Wiseman, R. L. & Grotjahn, D. A. Quantifying
666 organellar ultrastructure in cryo-electron tomography using a surface morphometrics pipeline.
667 *J. Cell Biol.* **222**, (2023).
- 668 7. Wu, G.-H. *et al.* CryoET reveals organelle phenotypes in huntington disease patient iPSC-
669 derived and mouse primary neurons. *Nat. Commun.* **14**, 692 (2023).
- 670 8. Xue, L. *et al.* Visualizing translation dynamics at atomic detail inside a bacterial cell. *Nature*
671 **610**, 205–211 (2022).
- 672 9. Tegunov, D., Xue, L., Dienemann, C., Cramer, P. & Mahamid, J. Multi-particle cryo-EM
673 refinement with M visualizes ribosome-antibiotic complex at 3.5 Å in cells. *Nat. Methods* **18**,
674 186–193 (2021).
- 675 10. Van Veen, D. *et al.* Missing Wedge Completion via Unsupervised Learning with
676 Coordinate Networks. *Int. J. Mol. Sci.* **25**, (2024).
- 677 11. Rice, G. *et al.* TomoTwin: generalized 3D localization of macromolecules in cryo-
678 electron tomograms with structural data mining. *Nat. Methods* **20**, 871–880 (2023).
- 679 12. Moebel, E. *et al.* Deep learning improves macromolecule identification in 3D cellular
680 cryo-electron tomograms. *Nat. Methods* **18**, 1386–1394 (2021).
- 681 13. de Teresa-Trueba, I. *et al.* Convolutional networks for supervised mining of molecular
682 patterns within cellular context. *Nat. Methods* **20**, 284–294 (2023).
- 683 14. Bepler, T. *et al.* Positive-unlabeled convolutional neural networks for particle picking in
684 cryo-electron micrographs. *Nat. Methods* **16**, 1153–1160 (2019).

- 685 15. Lamm, L. *et al.* MemBrain: A deep learning-aided pipeline for detection of membrane
686 proteins in Cryo-electron tomograms. *Comput. Methods Programs Biomed.* **224**, 106990
687 (2022).
- 688 16. Ermel, U. *et al.* A data portal for providing standardized annotations for cryo-electron
689 tomography. *Nat. Methods* (2024) doi:10.1038/s41592-024-02477-2.
- 690 17. Gubins, I. *et al.* SHREC 2020: Classification in cryo-electron tomograms. *Comput.*
691 *Graph.* **91**, 279–289 (2020).
- 692 18. Gubins, I. *et al.* SHREC 2021: Classification in Cryo-electron Tomograms. *Eurographics*
693 *Workshop 3D Object Retr.* (2021) doi:10.2312/3DOR.20211307.
- 694 19. Jeon, M. *et al.* CryoBench: Diverse and challenging datasets for the heterogeneity
695 problem in cryo-EM. *ArXiv Prepr. ArXiv240805526* (2024).
- 696 20. Hanson, S. M., Woppard, G. & Astore, M. The Inaugural Flatiron Institute Cryo-EM
697 Heterogeneity Community Challenge. (2024) doi:10.17605/OSF.IO/8H6FZ.
- 698 21. Keshavarz-Joud, P. *et al.* Exploring the Landscape of the PP7 Virus-like Particle for
699 Peptide Display. *ACS Nano* **17**, 18470–18480 (2023).
- 700 22. Kim, K. *et al.* The structure of natively iodinated bovine thyroglobulin. *Acta Crystallogr.*
701 *Sect. D* **77**, 1451–1459 (2021).
- 702 23. Zivanov, J. *et al.* New tools for automated high-resolution cryo-EM structure
703 determination in RELION-3. *eLife* **7**, e42166 (2018).
- 704 24. Maki-Yonekura, S., Kawakami, K., Takaba, K., Hamaguchi, T. & Yonekura, K.
705 Measurement of charges and chemical bonding in a cryo-EM structure. *Commun. Chem.* **6**,
706 98 (2023).
- 707 25. Natchiar, S. K., Myasnikov, A. G., Kratzat, H., Hazemann, I. & Klaholz, B. P.
708 Visualization of chemical modifications in the human 80S ribosome structure. *Nature* **551**,
709 472–477 (2017).
- 710 26. Zheng, S. *et al.* AreTomo: An integrated software package for automated marker-free,

- 711 motion-corrected cryo-electron tomographic alignment and reconstruction. *J. Struct. Biol.* **6**,
712 100068 (2022).
- 713 27. Ishemgulova, A., Noble, A. J., Bepler, T. & De Marco, A. Preparation Of Labeled Cryo-
714 ET Datasets For Training And Evaluation Of Machine Learning Models.
- 715 28. Goodenough, D. *et al.* Method and phantom to study combined effects of in-plane (x,y)
716 and z-axis resolution for 3D CT imaging. *J. Appl. Clin. Med. Phys.* **17**, 440–452 (2016).
- 717 29. Cho, N. H. *et al.* OpenCell: Endogenous tagging for the cartography of human cellular
718 organization. *Science* **375**, eabi6983 (2022).
- 719 30. Hein, M. Y. *et al.* Global organelle profiling reveals subcellular localization and
720 remodeling at proteome scale. *bioRxiv* (2023) doi:10.1101/2023.12.18.572249.
- 721 31. Wang, F. *et al.* Amino and PEG-amino graphene oxide grids enrich and protect samples
722 for high-resolution single particle cryo-electron microscopy. *J. Struct. Biol.* **209**, 107437
723 (2020).
- 724 32. Wang, F. *et al.* General and robust covalently linked graphene oxide affinity grids for
725 high-resolution cryo-EM. *Proc. Natl. Acad. Sci.* **117**, 24269–24273 (2020).
- 726 33. Chaillet, M. L. *et al.* Extensive Angular Sampling Enables the Sensitive Localization of
727 Macromolecules in Electron Tomograms. *Int. J. Mol. Sci.* **24**, (2023).
- 728 34. Hrabe, T. *et al.* PyTom: A python-based toolbox for localization of macromolecules in
729 cryo-electron tomograms and subtomogram analysis. *J. Struct. Biol.* **178**, 177–188 (2012).
- 730 35. Meng, E. C. *et al.* UCSF ChimeraX: Tools for structure building and analysis. *Protein*
731 *Sci.* **32**, e4792 (2023).
- 732 36. Ermel, U. H., Arghittu, S. M. & Frangakis, A. S. ArtiaX: An electron tomography toolbox
733 for the interactive handling of sub-tomograms in UCSF ChimeraX. *Protein Sci.* **31**, e4472
734 (2022).
- 735 37. Lehtinen, J. *et al.* Noise2Noise: Learning Image Restoration without Clean Data. Preprint
736 at <https://doi.org/10.48550/arXiv.1803.04189> (2018).

- 737 38. Bepler, T., Kelley, K., Noble, A. J. & Berger, B. Topaz-Denoise: general deep denoising
738 models for cryoEM and cryoET. *Nat. Commun.* **11**, 5208 (2020).
- 739 39. Sofroniew, N. *et al.* napari: a multi-dimensional image viewer for Python. Zenodo
740 <https://doi.org/10.5281/zenodo.13619583> (2024).
- 741 40. Moore, J. *et al.* OME-Zarr: a cloud-optimized bioimaging file format with international
742 community support. *Histochem. Cell Biol.* **160**, 223–251 (2023).
- 743 41. Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC: algorithms for
744 rapid unsupervised cryo-EM structure determination. *Nat. Methods* **14**, 290–296 (2017).
- 745 42. Tang, Y. *et al.* Self-Supervised Pre-Training of Swin Transformers for 3D Medical Image
746 Analysis. in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
747 (CVPR) 20698–20708 (2022). doi:10.1109/CVPR52688.2022.02007.
- 748 43. Martinez-Sanchez, A., Lamm, L., Jasnin, M. & Phelipeau, H. Simulating the cellular
749 context in synthetic datasets for cryo-electron tomography. *IEEE Trans. Med. Imaging* 1–1
750 (2024) doi:10.1109/TMI.2024.3398401.
- 751 44. Iudin, A. *et al.* EMPIAR: the Electron Microscopy Public Image Archive. *Nucleic Acids*
752 *Res.* **51**, D1503–D1511 (2023).
- 753 45. The wwPDB Consortium. EMDB—the Electron Microscopy Data Bank. *Nucleic Acids*
754 *Res.* **52**, D456–D465 (2024).
- 755 46. Paszke, A. *et al.* PyTorch: An Imperative Style, High-Performance Deep Learning
756 Library. Preprint at <https://doi.org/10.48550/arXiv.1912.01703> (2019).
- 757 47. Martín Abadi *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous
758 Systems. (2015).
- 759 48. Neuroglancer. <https://github.com/google/neuroglancer> (accessed Feb 9, 2019). [Google](#)
760 [Scholar](#).
- 761 49. Moult, J., Pedersen, J. T., Judson, R. & Fidelis, K. A large-scale experiment to assess
762 protein structure prediction methods. *Proteins Struct. Funct. Bioinforma.* **23**, ii–iv (1995).

- 763 50. Senior, A. W. *et al.* Improved protein structure prediction using potentials from deep
764 learning. *Nature* **577**, 706–710 (2020).
- 765 51. Kovalevskiy, O., Mateos-Garcia, J. & Tunyasuvunakool, K. AlphaFold two years on:
766 Validation and impact. *Proc. Natl. Acad. Sci.* **121**, e2315002121 (2024).
- 767 52. Varadi, M. *et al.* AlphaFold Protein Structure Database in 2024: providing structure
768 coverage for over 214 million protein sequences. *Nucleic Acids Res.* **52**, D368–D375 (2024).
- 769 53. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,
770 583–589 (2021).
- 771 54. Rosenthal, P. B. & Henderson, R. Optimal Determination of Particle Orientation,
772 Absolute Hand, and Contrast Loss in Single-particle Electron Cryomicroscopy. *J. Mol. Biol.*
773 **333**, 721–745 (2003).

774 Software Availability

775 Software packages used in this study for reconstruction of tomograms, annotation generation
776 and annotation curation are released under open source licenses and available in public
777 repositories. Refer to **Supplementary Table 1** for a summary of packages.

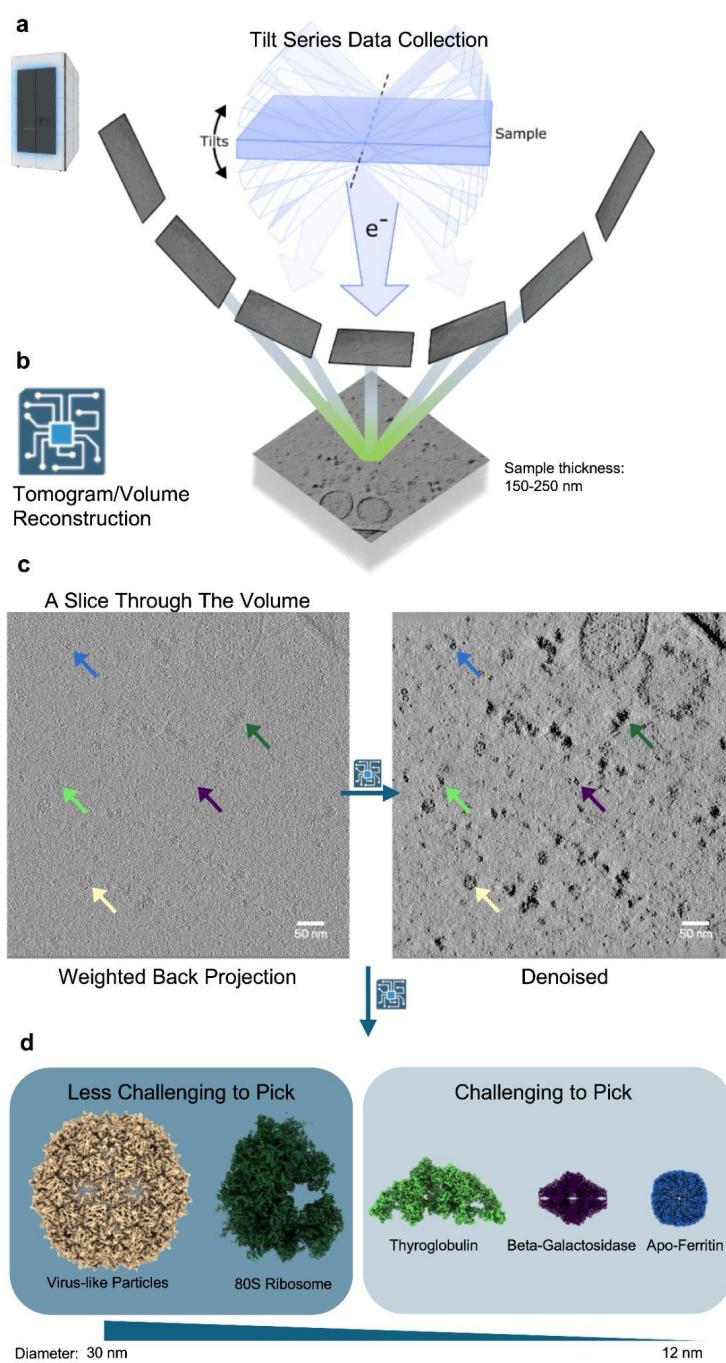
778 Contributions

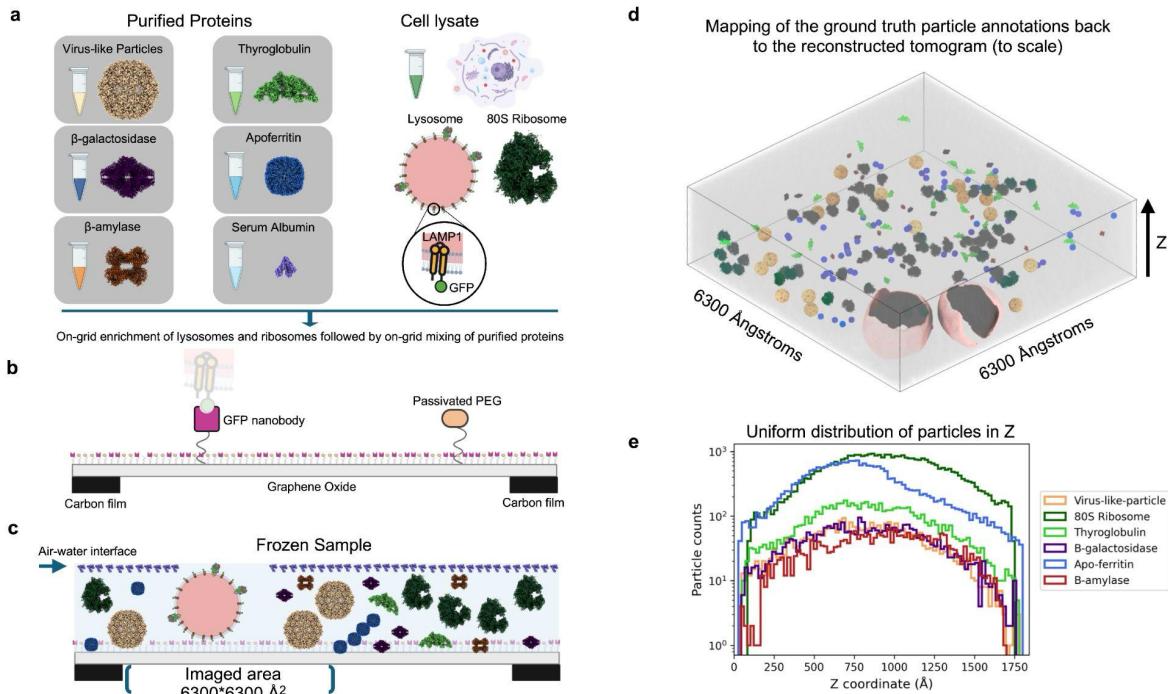
779 **M.P.**, **K.H.**, and **B.C.** developed the idea for a machine learning challenge and managed the
780 project. **M.P.**, **H.S.**, **D.S.**, and **F.W.** optimized grid functionalization and freezing of organelles
781 and carried out the sample preparation. **M.P.** and **E.M.** carried out the data collection. **M.P.**
782 curated the tomograms based on sample quality and tilt series alignment quality. **S.Z.**, **A.P.**,
783 **Y.Y.**, **J.S.**, and **M.P.** developed the data processing pipeline. **S.Z.** and **A.P.** developed
784 AreTomo3 CTF deconvolution and DenoisET. **S.Z.**, **A.P.**, and **Y.Y.** developed the slab method.
785 **U.H.E.** developed copick and ChimeraX-copick plugin. **K.H.** and **Z.Z.** developed CellCanvas and
786 copicklive. **J.S.** developed DeepFindET. **D.K.** and **J.S.** developed the 3D refinement pipelines.
787 **K.H.**, **Z.Z.**, and **J.S.** developed the challenge metrics and the dataset splitting. **M.P.** and **Y.Y.**
788 manually picked initial training sets. **M.P.**, **A.P.**, **Y.Y.**, and **J.S.** prepared the final curations for all
789 the picks for all the species. **S.K.**, **J.S.**, **K.H.**, and **Z.Z.** prepared the example notebooks. **U.H.E.**
790 and **A.C.** developed the cryoET data portal. **M.P.**, **A.P.**, **U.H.E.**, **K.H.**, **Z.Z.**, and **B.C.** wrote the
791 manuscript. All authors reviewed the manuscript. **M.H.**, **D.A.**, **C.P.** and **B.C.** provide overall
792 leadership for projects at the Chan Zuckerberg Imaging Institute.

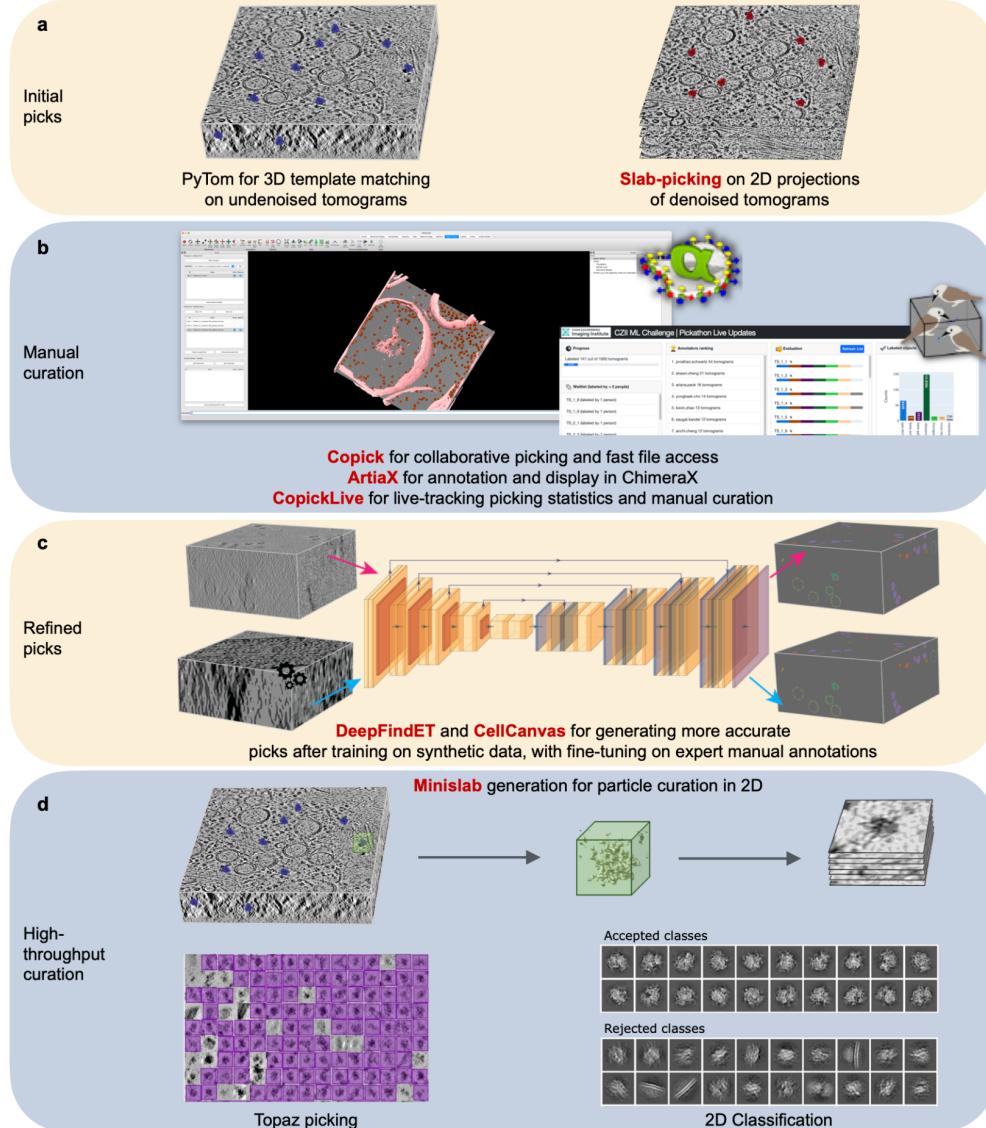
793 Acknowledgements

- 794 ● These authors contributed equally to this work and were listed alphabetically: Anch
795 Cheng, Utz Heinrich Ermel, Saugat Kandel, Dari Kimanius, Elizabeth Montabana, Daniel
796 Serwas, Hannah Siems, Feng Wang, Zhuowen Zhao, Shawn Zheng.
- 797 ● Emma Lundberg (Stanford), Ellen Zhong (Princeton), Thorsten Wagner (MPI of
798 Molecular Physiology), Tristan Bepler (NYSBC), Robert Kiewisz (NYSBC), Alister Burt
799 (Genentech), and Lorenzo Gaifas (Grenoble) who provided valuable insights for the
800 design of the challenge.
- 801 ● Our collaborators at CZ BioHub SF, Manuel Leonetti, Shivanshi Vaid, Madhuri
802 Vangipuram, and Rodrigo Baltazar generated the HEK293T LAMP1-GFP cell lines and
803 shared their cell lysis protocol.
- 804 ● David Agard's lab at UCSF, Feng Wang and Simon Sander, provided their grid
805 functionalization protocol and grids and helped us optimize the protocol further for
806 organelles.
- 807 ● The GFP-nanobody construct was modified and provided by Peng Jin (UCSF, Jan Lab).
- 808 ● Our collaborators at NYSBC, Mykhailo Kopylov and Charlie Dubbledam, provided VLPs.
- 809 ● Chan Zuckerberg Initiative SciTech members who participated in the pickathon: Ashley
810 Anderson, Ben Nelson, Jun Ni, Ellaine Chou, Jessica Gadling, Kandarp Khandwala,
811 Chili Chiu, Ann Jones, Timmy Huang, Janeece Pourroy, Dannielle McCarthy, Andy
812 Sweet, Eric Wang, Kirsty Ewing, Mikala Caton, Manasa Venkatakrishnan, Kira Evans.
- 813 ● CZ Imaging Institute members who participated in the pickathon: Yongbaek Cho, Nina
814 Borja, Norbert Hill, Carmela Villegas, Shu-Hsien Sheu, Gorica Margulis, Noeli
815 Pazzoldan.
- 816 ● Chan Zuckerberg Initiative SciTech team who contributed to the development of the
817 cryoET data portal: Jun Xi Ni, Jessica Gadling, Manasa Venkatakrishnan, Kira Evans,
818 Jeremy Asuncion, Andrew Sweet, Janeece Pourroy, Zun Shi Wang, Kandarp
819 Khandwala, Benjamin Nelson, Dannielle McCarthy, Eric M Wang, Richa Agarwal, Trent
820 Smith, Bryan Chu, Dana Sadgat, Erin Hoops, Justine Larsen.
- 821 ● Kristen Maitland and Stephani Otte for their support in the planning and execution of the
822 competition.
- 823 ● Samantha Yammie who provided valuable feedback on the manuscript text.
- 824 ● Some schematic elements in Figure 1 (panels a,b), Figure 2 (panels a,b,c) and
825 Extended Data Figure 1 were made using BioRender: Created in BioRender. Paraan, R.,
826 Serwas, D. (2024) BioRender.com/r59u258
- 827 ● Grant Reference: CZ Imaging Institute is made possible with support from Chan
828 Zuckerberg Initiative (CZII-2023-327779).

830 **Figure 1 | Cryo-electron
831 tomography and its challenges. a.**
832 Tilt series data collection inside a
833 Krios G4 transmission electron
834 microscope. The frozen sample (blue
835 slab) is rotated to generate 2D
836 projection images on a fixed direct
837 electron detector. The 2D images
838 have an SNR of ~0.1. The range of
839 the tilted images is only
840 representative, typically there are 31
841 to 61 tilted images. **b.** A 3D volume
842 is computed using real space back-
843 projection²⁶. The X and Y dimensions
844 of the volume (630 nm) depend on
845 the magnification, while Z (150-250
846 nm) approximates the physical
847 thickness of the sample. The thicker
848 the sample, the lower the SNR. **c.** A
849 slice through the reconstructed
850 volume before and after denoising.
851 The thickness of the slice is equal to
852 5 Å (Ångstroms = 10^{-10} meters). One
853 example for each of the 5 species
854 that are the targets of the challenge is
855 indicated by an arrow. The colors
856 correspond to the colored structures
857 in panel d. **d.** 3D reconstructed
858 volumes of the protein complexes
859 that are the targets of the challenge.
860 These volumes/maps are from
861 published data (refer to Table 1, also
862 refer to Extended Data Fig. 5 for the
863 maps reconstructed from the
864 phantom data). In a common cryoET
865 workflow, for each species thousands
866 of copies of that protein are annotated,
867 extracted, and further processed to produce a high-
868 resolution map. These high-resolution maps reveal how the proteins are organized in space and
869 how they carry out their functions.
870

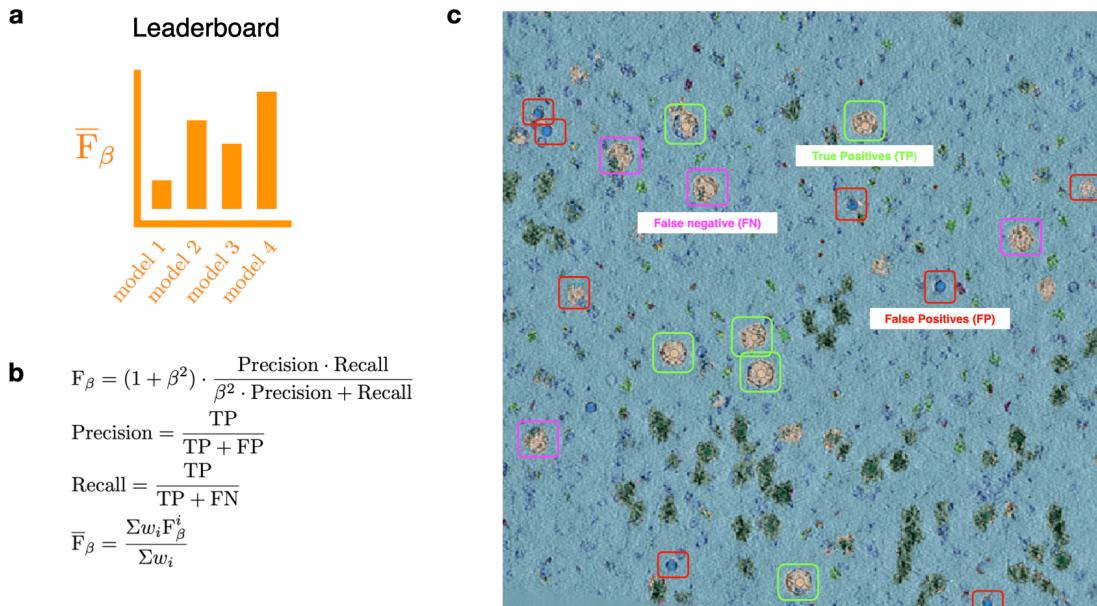






890

891 **Figure 3 | Multi-step workflow developed to generate ground truth labels.** **a.** Particle picks
892 were initially generated either by PyTom^{33,34} 3D template matching algorithm or 2D slab-picking
893 from 300 Å thick projections through denoised tomograms. **b.** These picks were manually
894 curated in ChimeraX³⁵, using ArtiaX³⁶ for annotation, copick for rapid file access, and copicklive
895 to track the progress of curation. **c.** Two ML pipelines, DeepFindET and CellCanvas, were
896 trained on synthetic tomograms and fine-tuned on a selection of the curated picks from expert
897 annotators to generate a more refined and complete set of picks. **d.** Final rounds of curation
898 were performed by applying Topaz, 2D classification, and multi-class *ab initio* reconstruction to
899 per-particle 2D projections (“minislabs”) generated from denoised tomograms to filter out
900 contamination and non-target particles. New tools developed as part of this curation effort are
901 indicated in red.
902
903



904
905
906
907
908
909
910
911
912

Figure 4 | Evaluation of models' picking performance. **a.** The leaderboard ranking is determined by **b.** the weighted average F_β score, where the F_β score is first calculated for individual protein species and then combined across all five species using weighted aggregation. **c.** A tomogram slice annotated in CellCanvas demonstrates model detection outcomes. True positives, where the model correctly identifies VLPs, are marked with green squares. Missed VLPs, or false negatives, are indicated by magenta boxes, while structures mistakenly identified as VLPs (false positives) are highlighted in red.

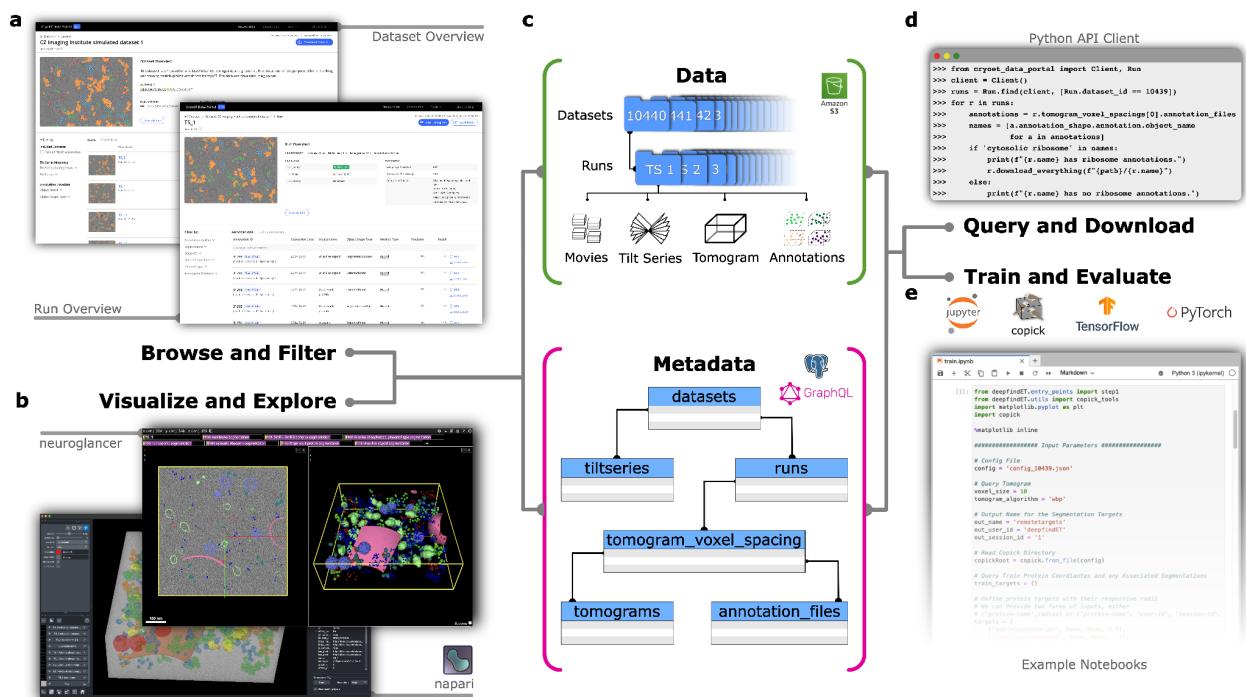


Figure 5 | The cryoET data portal and related infrastructure. **a.** The cryoET data portal web interface allows challenge participants to explore the experimental and simulated Challenge datasets (dataset IDs 10440 and 10441), included runs, annotation and their metadata. **b.** Tomograms and annotations can be visualized directly inside the web browser using Neuroglancer⁴⁸, or in the desktop viewer napari³⁹. **c.** All related image data and metadata are provided from a public AWS S3 bucket and associated metadata can be queried using a publicly exposed GraphQL API. **d.** Both resources can be accessed programmatically using the cryoET data portal python API client, which also provides convenient methods for data download. **e.** All cryoET data portal resources (including other datasets) can be accessed using copick and the provided example notebooks, which feature reference implementations in TensorFlow⁴⁷ and PyTorch⁴⁶.