

1 **TomoTwin: Generalized 3D Localization of Macromolecules in Cryo-electron**  
2 **Tomograms with Structural Data Mining**

3

4 Gavin Rice<sup>‡</sup>, Thorsten Wagner<sup>‡</sup>, Markus Stabrin, Stefan Raunser\*

5

6 Department of Structural Biochemistry, Max Planck Institute of Molecular Physiology, Otto-Hahn-

7 Str. 11, 44227 Dortmund, Germany

8

9 <sup>‡</sup>Both authors contributed equally

10 \*Corresponding author: [stefan.raunser@mpi-dortmund.mpg.de](mailto:stefan.raunser@mpi-dortmund.mpg.de)

11

12

13 **Abstract:**

14 **Cryoelectron tomography enables the visualization of cellular environments in**  
15 **extreme detail through the lens of a benign observer; what remains lacking**  
16 **however are tools to analyze the full amount of information contained within these**  
17 **densely packed volumes. Detailed analysis of macromolecules through**  
18 **subtomogram averaging requires particles to first be localized within the**  
19 **tomogram volume, a task complicated by several factors including a low signal to**  
20 **noise ratio and crowding of the cellular space. Available methods for this task**  
21 **suffer either from being error prone or requiring manual annotation of training data.**  
22 **To assist in this crucial particle picking step, we present TomoTwin: a robust, first**  
23 **in class general picking model for cryo-electron tomograms based on deep metric**  
24 **learning. By embedding tomograms in an information-rich, high-dimensional**  
25 **space which separates macromolecules according to their 3-dimensional**  
26 **structure, TomoTwin allows users to identify proteins in tomograms *de novo***  
27 **without manually creating training data or retraining the network each time a new**  
28 **protein is to be located. TomoTwin is open source and available at**  
29 **<https://github.com/MPI-Dortmund/tomotwin-cryoet>.**

30

## 31 **Main Text:**

## 32 **Introduction**

33 In recent years, cryo-electron tomography (cryo-ET) has emerged as a landmark  
34 technique for the visualization of macromolecules within their native cellular  
35 environment<sup>1-7</sup>. Advances in high-pressure freezing and the advent of focused ion beam  
36 (FIB) milling at cryogenic temperatures now allow for the routine preparation of thin (<  
37 200 nm) lamellae from cells or even small organisms<sup>8-10</sup>. Performing cryo-ET on these  
38 thin lamellae offers a unique opportunity to capture cellular processes in 3D and in  
39 unprecedented detail. Subsequent analysis of specific macromolecules from tomographic  
40 volumes through subtomogram averaging (STA) allows in depth structural determination  
41 of macromolecular complexes in their native environment<sup>11-14</sup>. Particularly when  
42 complemented by recent advances in structure prediction such as alphafold2, STA forms  
43 a powerful crossbridge between protein biochemistry and cellular proteomics<sup>15-17</sup>. In  
44 order to perform STA however, particles of a macromolecule of interest must first be  
45 located within the tomographic volume, a task complicated by the 3D nature of these data.

46 The accurate localization of macromolecules inside cryo-electron tomograms is a  
47 well-recognized barrier for studying cellular life at the mesoscopic level, sparking  
48 competitions such as the annual Classification in Cryo-Electron Tomograms (SHREC)  
49 competition where contestants submit algorithms to localize proteins in tomograms with  
50 a benchmark set by template matching<sup>18</sup>. This has led to the development of several deep  
51 learning-based tools with high picking accuracies often achieved by leveraging popular  
52 3D-Unet convolutional neural network (CNN) architectures<sup>19-21</sup>. Each of these

53 approaches is unified however in the fact that they share a non-generalizing workflow,  
54 meaning that for each protein of interest, users must first manually pick the protein in at  
55 least one tomogram and retrain the neural network to identify that protein. Not only is this  
56 incompatible with the future directions of automated tomogram reconstruction and STA,  
57 but for many proteins picking sufficient training data by eye is not possible. With a minimal  
58 requirement for user-input, template matching<sup>22,23</sup> is still often utilized in cryo-ET  
59 processing workflows that place an emphasis on throughput<sup>24</sup> although at the cost of  
60 picking accuracy.

61 One method to retain the accuracy of deep learning-based picking while  
62 circumventing the requirement of manually annotating training data for each protein of  
63 interest is to train a model to learn a generalized representation of 3D molecular shape  
64 that then can differentiate between macromolecules based on their structure. Such  
65 approaches have demonstrated profound impact for particle picking in 2D for single  
66 particle cryo-electron microscopy analysis<sup>25–28</sup>.

67 Particularly well suited for this type of generalization is deep metric learning in  
68 which data are encoded as a high-dimensional representation, called an embedding,  
69 where one or more learned characteristics of the data are related to distance in the  
70 embedding space<sup>29,30</sup>. During training, the model is penalized for placing data from  
71 different classes near to one another and rewarded for placing data from the same class  
72 close together in the embedding space<sup>31</sup>. Therefore, over the training process the model  
73 learns to place data from each class within a distinct region of the embedding space  
74 where more similar classes are placed closer together and dissimilar ones further apart.  
75 In some cases, the embeddings of a dataset are sufficiently ordered to allow for *de novo*

76 identification of classes based on their clustering in the embedding space<sup>31</sup>. By  
77 understanding similarity relationships, deep metric learning models have demonstrated  
78 acute adaptability when presented with new classes of data, being able to place them in  
79 the embedding space according to their similarity to known classes without requiring  
80 retraining<sup>31–33</sup>.

81 Here we present TomoTwin, a generalized particle picking model and deep metric  
82 learning toolkit for structural data mining of cryo-electron tomograms. We supply two  
83 workflows for macromolecular localization with TomoTwin, a reference-based workflow in  
84 which a single molecule is picked for each protein of interest and used as a target, and a  
85 *de novo* clustering workflow where macromolecular structures of interest are identified on  
86 a 2D manifold. Trained on a diverse set of simulated tomograms, the picking model of  
87 TomoTwin is able to locate new proteins with high accuracy in not only simulated data,  
88 but experimental and cellular tomograms as well. TomoTwin combines the high accuracy  
89 of deep learning-based particle picking with high throughput processing by removing the  
90 step of manual annotation of training data and model training, and allows simultaneous  
91 picking of several proteins of interest in each tomogram.

92

## 93 **Results:**

### 94 **Overview of functions, build, and philosophy behind TomoTwin**

95 The machine learning backbone of TomoTwin is built on the principle of learning  
96 generalized representations of 3-dimensional shapes in tomograms ([Supplementary Fig.](#)  
97 [1a,b](#)). Trained with deep metric learning, the 3D-CNN is able to locate not only  
98 macromolecules contained in the training set, but novel macromolecules in tomograms

99 as well. This allows TomoTwin to retain the high fidelity of deep learning-based particle  
100 picking while avoiding the burden of requiring retraining for each protein of interest. The  
101 trained model plots tomogram subvolumes as points in a high-dimensional embedding  
102 space organized according to the similarity of their macromolecular contents  
103 ([Supplementary Fig. 1c](#)). Once this high-dimensional space is mapped for a tomogram,  
104 particles of each macromolecule can be picked by identifying their associated region in  
105 of the embedding space. This can be done either by identifying a single example of each  
106 protein of interest in a tomogram and using them to mark the region of the space where  
107 they are embedded to create a target embedding (reference-based workflow), or by  
108 plotting the tomogram embeddings onto a 2D manifold where clusters of subvolumes for  
109 each macromolecule can be identified by eye (clustering workflow) ([Fig. 1a,b](#)). Once the  
110 subvolumes containing a protein of interest are identified in the embedding space, they  
111 must be mapped back to real space in the tomogram where overlapping picks of the same  
112 molecule can be consolidated into one centralized pick per molecule ([Fig. 1c](#)). Finally,  
113 TomoTwin allows users to interactively filter the picked particles for each macromolecule  
114 of interest based on the particle size and the network's confidence level, which is encoded  
115 as the distance between each subvolume and the target embedding for that  
116 macromolecule in the embedding space ([Fig. 1d](#), [Supplementary Fig. 2](#)).

117

## 118 **Two workflows to identify and locate macromolecules in tomograms**

119 TomoTwin represents tomograms in a high-dimensional space where subvolumes of  
120 each macromolecule are embedded in a distinct region of the space. In order to identify  
121 which region of the embedding space a macromolecule is located in, we provide the user

122 with two workflows – a reference-based workflow and a clustering workflow. Each  
123 workflow picks particles with high accuracy, but the reference-based approach begins  
124 with identifying an example of the protein of interest in the tomogram and mapping this to  
125 the embedding space whereas the clustering workflow begins with identifying a region of  
126 the embedding space and mapping this to the tomogram. Which workflow is most suitable  
127 for any given application depends on how easily the protein(s) of interest can be identified  
128 in the tomogram versus the embeddings. Both workflows share the common first step of  
129 using the embedding function of TomoTwin to generate a high-dimensional embedding  
130 of the entire volume of the tomogram called a tomogram’s representation map (Fig. 1a).

131 In the reference-based workflow, users identify a single molecule of each protein  
132 of interest in a tomogram and embed it to generate a target in the embedding space for  
133 that protein. In the clustering workflow, TomoTwin approximates the representation map  
134 of the tomogram onto a 2-dimensional manifold. This 2D manifold can then be directly  
135 visualized by the user who can then outline one or more clusters of interest using the  
136 Lasso function of TomoTwin. The Lasso function then computes the center coordinate of  
137 the drawn cluster in the high-dimensional embedding space to be used as a target  
138 embedding in lieu of a reference (Fig. 1b). The map function of TomoTwin takes as input  
139 the tomogram embeddings and target embeddings and calculates the distance matrix  
140 between the target(s) and each point in the tomogram embeddings. The distances are  
141 mapped to the coordinates of each subvolume, constructing a similarity map of proposed  
142 particle locations within the tomogram for each protein of interest (Fig. 1c). The Locate  
143 function uses this similarity map to localize peaks of high similarity and generate  
144 candidate particle positions. Finally, the Pick function of TomoTwin uses these candidate

145 positions as well as adjustable size and confidence thresholds to pick particles in the  
146 tomogram producing a coordinate file for each protein of interest to be then used for  
147 subtomogram averaging or other analysis ([Fig. 1d](#)).

148

### 149 **Training of the general picking model**

150 To produce a picking model capable of localizing novel macromolecules within  
151 tomograms without requiring retraining, TomoTwin is trained using deep metric learning  
152 on triplets of subvolumes from simulated tomograms. A set of 120 structurally dissimilar  
153 proteins procured from the Protein Data Bank (PDB) ranging in size from 30 kDa to 2.7  
154 mDa were used to simulate 84 tomograms containing a total of 120,000 subtomogram  
155 particles ([Supplementary Fig. 3](#)). During training, batches of subvolumes are embedded  
156 by a custom-built 3D CNN which transforms each 37x37x37 realspace 3D subvolume to  
157 a 1D, 32-length coordinate vector located on a high-dimensional embedding manifold  
158 molded to the surface of a 32D hypersphere ([Supplementary Fig. 1a](#)).

159         These coordinate vectors are used in the metric learning process which rewards  
160 the model for placing the anchor and positive close together in the embedding space and  
161 penalizes it for placing the anchor and negative near one another. Therefore, through  
162 training TomoTwin learns to place subvolumes of each macromolecule within a distinct  
163 region of the embedding space, where more structurally similar macromolecules are  
164 placed closer together and dissimilar ones further apart ([Supplementary Fig. 1c](#)). By  
165 training on a large, diverse set of 3D macromolecular shapes and sizes, TomoTwin learns  
166 a generalized representation of 3D macromolecular shapes which it leverages to place



167 novel macromolecules in their own region of this embedding space relative to their  
168 structural similarity to known proteins without requiring retraining.

169

170 **The general picking model accurately locates particles across a wide range of**  
171 **shapes and sizes**

172 Because *a priori* information on the ground truth locations of all molecules in a tomogram  
173 is not possible to obtain for experimental data, the picking performance of the trained  
174 model was first assessed on the simulated tomograms containing proteins from the  
175 training set where the F1 picking score was calculated from the true positive, false  
176 positive, true negative, and false negative picks as described in Methods.

177 The median F1 picking score across all validation tomograms was 0.88 with a  
178 range from 0.76 to 0.98 ([Supplementary Fig. 4a](#)). Across all proteins in the training set,  
179 the median validation F1 picking score is 0.92 ([Supplementary Fig. 4b](#)). In rare cases,  
180 outlier scores were observed where specific proteins were unable to be picked across a  
181 range of sizes ([Supplementary Fig. 4c](#)). Closer inspection of these outliers revealed that  
182 in the simulated tomograms, each of these proteins display a particularly weak signal  
183 when compared to proteins of similar size ([Supplementary Fig. 4d](#)). In these cases, it  
184 appears that these proteins display a shape that is not recovered well during tomogram  
185 reconstruction by weighted back-projection. Despite this, picking on the validation  
186 tomograms demonstrated high accuracy for proteins across a wide array of shapes and  
187 sizes ranging from 30 kDa to 2.7 mDa.

188

189

## 190 **TomoTwin generalizes to unseen proteins**

191 In order to assess the generalization of the general picking model to particles that were  
192 not in the training data set, we measured the picking performance with other, previously  
193 unseen proteins in a simulated tomogram with the reference-based workflow. We  
194 measured the F1 score of proposed particle locations against ground-truth boxes for  
195 seven proteins not included in the training data for which TomoTwin was therefore naïve  
196 ([Fig. 2a,b](#)). This assessment revealed that when trained on a set of 120 dissimilar proteins  
197 ([Supplementary Fig. 3](#)), the resulting model was able to locate all seven proteins  
198 accurately with a median F1 score of 0.82 despite a lack of previous training on these  
199 proteins ([Fig. 2d](#)). To measure the effect of training set size on generalization accuracy  
200 we performed this analysis on picking models trained on 20, 50, 100, and 120 proteins  
201 where we observed a logarithmic increase in generalization accuracy with the number of  
202 proteins in the training set ([Fig. 2c](#)). This high accuracy in locating novel proteins indicates  
203 a high generalization capability of TomoTwin.

204 As TomoTwin is trained entirely on simulated data, it is paramount to investigate  
205 its ability to pick proteins of interest in experimental tomograms. To evaluate this, cryo-  
206 ET was performed on a sample containing a mixture of three proteins, namely  
207 apoferritin<sup>34</sup>, the Type VI secretion effector RhsA from *Pseudomonas protegens*<sup>34</sup>, and  
208 the Tc toxin A component TcdA1 from *Photobacterium luminescens*<sup>35</sup> as well as liposomes  
209 (DOPC/POPC) ([Fig. 3a](#)). This mixture was chosen to create an environment with several  
210 proteins of different sizes as well as liposomes to mimic non-protein structures that may  
211 confound picking accuracy. Ten reconstructed tomograms were picked for apoferritin,  
212 RhsA, and TcdA1 using the pretrained general model of TomoTwin. In each case, the

213 reference-based workflow was employed in which a target embedding was created by  
214 picking a single example of each protein as they are readily observable in tomograms  
215 with sufficient contrast. The target embedding from one tomogram was then applied  
216 across all tomograms in each dataset. Direct visualization of the picking similarity maps  
217 and final picking reveals high fidelity localization of each protein within the tomograms  
218 despite none of these proteins being included in the training set ([Fig. 3b](#)).

219 As ground-truth particle coordinates are not available for experimental data, the  
220 accuracy of the picking was assessed by extracting subvolumes at the picked coordinates  
221 of each protein, projection of the 3D subvolumes to 2D using SPHIRE<sup>36</sup>, and performing  
222 2D classification<sup>37</sup> to evaluate the picked particles in a reference-free manner  
223 ([Supplementary Fig. 5](#)). For each protein of interest, the number of particles in the 2D  
224 classes displaying a high similarity to 2D classes of the protein previously determined by  
225 single particle analysis were recorded and represented as a percentage of the total  
226 number of particles picked ([Fig. 3c](#)). The high proportion of particles in all positive classes  
227 indicates that the picking is of high accuracy, confirming the visual impression of the  
228 picking result.

229 One of the principal advantages of cryo-ET is the ability to directly visualize  
230 proteins in their native cellular environments. Due to crowding of the cellular space and  
231 the poor contrast caused by thick specimens however, particle localization within a  
232 cellular environment presents a significant challenge. To assess its ability to locate  
233 particles in cellular tomograms, we applied TomoTwin to a dataset of tomograms  
234 containing *Mycoplasma pneumoniae*<sup>38</sup> ([EMPIAR 10499](#)) ([Fig. 4a](#)). Using the TomoTwin  
235 general model, we picked 70S ribosomes in 65 tomograms with the reference-based

236 workflow in which a reference was identified on one tomogram and used to generate a  
237 target embedding that was then applied to the entire dataset (Fig. 4b,c). To visualize the  
238 results, we extracted pseudo-subtomograms<sup>39</sup> and performed 3D classification using a  
239 70S ribosome cryo-EM structure (EMD 11650) lowpass filtered to 30 Å as a reference.  
240 As all 3D classes resemble ribosomes refined to ~15 Å, it clearly indicates that TomoTwin  
241 also picks highly accurately in cellular tomograms (Fig. 4d).

242

### 243 **Structural Data Mining on the Embedding Manifold**

244 The embedding feature of TomoTwin constructs a representation of a tomogram as a  
245 series of high-dimensional embeddings. These high-dimensional embeddings can be  
246 directly visualized by approximation on a 2D manifold (Fig. 5a,c). As a result of our deep  
247 metric learning-based approach, these representations contain a wealth of information  
248 about the contents of a tomogram where the distance between two subvolume  
249 embeddings directly correlates to the similarity of the 3-dimensional macromolecular  
250 shapes contained within. Typically, these representations contain a large mass  
251 corresponding to background noise, or a particularly prominent feature of the tomogram  
252 volume as well as additional well-defined clusters corresponding to different shapes such  
253 as proteins, membranes, or fiducials. By directly visualizing these representations on a  
254 2D manifold, the clustering workflow of TomoTwin allows interactive, structural data  
255 mining of tomograms, where clusters of subvolumes on the embedding manifold are used  
256 to locate different macromolecular populations within the tomogram.

257 To evaluate the accuracy of clustering-based picking quantitatively, we again  
258 utilized our simulated generalization tomogram where we evaluated the results of the

259 clustering-based picking of each protein against the ground-truth coordinates with the F1  
260 picking score as well as directly comparing it against the reference-based workflow (Fig.  
261 5b,e). The clustering-based picking identified each protein with high accuracy across a  
262 range of sizes. Notably, it outperforms the reference-based workflow for glutamine  
263 synthetase<sup>40</sup> (PDB ID: 1FPY) indicating that this workflow provides complementary  
264 advantages to the reference-based workflow. Additionally notable in the manifold  
265 projection of the representation map is the fact that individual protein clusters are globally  
266 organized by size, with the three largest protein clusters located in one area of the map,  
267 clusters for medium sized proteins in another, and the clusters for the two smallest  
268 proteins located furthest away from those of the large proteins, demonstrating that the  
269 model accurately represents complex similarity relationships in terms of protein structures  
270 as distance in the embedding space (Fig. 5a).

271 We additionally compared the clustering-based picking workflow directly against  
272 the reference-based approach for the cellular tomograms containing *M. pneumoniae* (Fig.  
273 5c). Examining the representation maps of these tomograms, several clusters are visible.  
274 One of which, when picked, produces accurate particle locations for 70S ribosomes  
275 nearly identical to those produced by the reference-based approach once again  
276 underlining the robustness of both workflows (Fig. 5d).

277

## 278 **Conclusion:**

279 Despite offering the potential to study proteins in their native, cellular environment in  
280 unprecedented detail, it remains that, presently, only a select few proteins have been  
281 successfully studied in detail by cryo-ET with STA. In part, this is because with increased

282 cellular context, the formation of macromolecular complexes, and poorer contrast caused  
283 by thicker specimens, comes the challenge of picking individual proteins for subsequent  
284 subtomogram averaging. To assist in this crucial particle picking step, we developed  
285 TomoTwin, a robust, first in class general picking model for cryo-electron tomograms  
286 based on deep metric learning. TomoTwin allows users to identify proteins in tomograms  
287 *de novo* without manually creating training data or retraining the network each time a new  
288 protein is to be located.

289 The innovation landscape for algorithm development in both cryo-EM and cryo-ET  
290 bears a heavy emphasis on automated processing for increased data throughput<sup>26,37,41–</sup>  
291 <sup>43</sup>. With its highly generalizable picking model, TomoTwin is the first tool based on deep  
292 learning that can be readily integrated with high throughput tomogram reconstruction and  
293 STA workflows. Additionally, when combined with unsupervised cluster detection  
294 algorithms<sup>44</sup>, the clustering workflow of TomoTwin paves the way for unsupervised STA  
295 analysis on a whole-tomogram level ([Supplementary Fig. 6](#)).

296 TomoTwin is a robust, open-source tool for particle localization in cryo-electron  
297 tomograms. The code used to develop and train TomoTwin as well as the general picking  
298 model and tools to use it for generalized particle picking are available at  
299 <https://github.com/MPI-Dortmund/tomotwin-cryoet> with future updates including  
300 extensive user documentation available soon.

301

## 302 **Methods**

### 303 **Training Data Generation**

304 TomoTwin was trained on 123 data classes comprised of subvolumes of 120 different  
305 proteins, membranes, noise, and fiducials from simulated tomograms. To ensure that  
306 TomoTwin is trained on the most diverse set of proteins possible, 108 proteins were  
307 selected from the PDB with sizes ranging from 30 kDa to 2.7 mDa and the cross  
308 correlation between pairs of 10 Å low-pass filtered maps of each protein was calculated  
309 ([Supplementary Figure 3](#)). Any protein with a high similarity (greater than 0.6) to another  
310 protein in the training set was marked for replacement. Additionally included were the  
311 data from the 2021 SHREC competition including 12 proteins<sup>18</sup> to yield a total of 120  
312 proteins for training. A training/validation split was achieved with 800 subvolumes for each  
313 data class in the training set and 200 in the validation set, yielding a total training set size  
314 of 98,400 subvolumes and a validation set size of 24,600 subvolumes.

315

### 316 **Tomogram simulation**

317 Tomogram simulation was done using TEM Simulator<sup>45</sup> which calculates the scattering  
318 potential of individual proteins and places them in definable positions within the volume.  
319 The output of the simulation is a tilt series which is then reconstructed using IMOD<sup>46</sup>. A  
320 configuration file was generated with properties for the electron beam, optics of the  
321 microscope, the detector, the tilt geometry and the sample volume. The default detector  
322 was adjusted to reflect the MTF curve of a modern Gatan K3 Camera with a quantum  
323 efficiency of 0.9. The detector size was set to 1024x1024 with a pixel size of 5 micrometer.  
324 The magnification was set to 9800, the spherical aberration and chromatic aberration

325 were adjusted to 2.7 mm and 2 mm respectively to mimic popular modern TEMs. A  
326 condenser aperture size of 80 micrometer was chosen. For each tomogram the defocus  
327 value was randomly chosen between -2.5  $\mu\text{m}$  and -5  $\mu\text{m}$ . A tilting scheme of -60° to +60°  
328 with a step size of 2° was used. To simplify and streamline the simulation we wrote a set  
329 of open-source programs called “tem-simulator-scripts” ([https://github.com/MPI-](https://github.com/MPI-Dortmund/tem-simulator-scripts)  
330 [Dortmund/tem-simulator-scripts](https://github.com/MPI-Dortmund/tem-simulator-scripts)). They contain scripts that require as input the PDB files  
331 to be simulated and the number of particles to simulate per PDB. The program then  
332 generates reconstructed tomograms as they were used for this study using the following  
333 pipeline:

- 334 1. Generation of densely packed random particle positions within the volume  
335 where individual particles do not overlap.
- 336 2. Generation of an occupancy map - a volume where each voxel is labeled  
337 according the protein identity.
- 338 3. Generation of fiducial maps.
- 339 4. Generation of vesicle maps.
- 340 5. Generation of the configuration file for TEM-simulator
- 341 6. Simulation of the tiltseries using TEM-simulator.
- 342 7. Alignment and reconstruction using IMOD.

343 However, all steps can also be carried out individually to have full control over all  
344 parameters.



345           Using this procedure, we simulated 11 sets of proteins. The sets contain in total  
346 108 different proteins with each set covering proteins of various sizes. For each set we  
347 simulated 8 tomograms of size 512x512x200 voxels with a pixel size of 1.02 nm and  
348 varying protein density. For tomogram 1, 2 and 8, 150 particles per protein were  
349 generated, for tomogram 3 and 4, 125 particles per protein, for tomogram 5 and 6, 100  
350 particles per protein and for tomogram 7, 75 particles per protein. Tomograms 1-7 were  
351 used for training and tomogram 8 for validation. The generated tomograms used in this  
352 study with all meta-data are publicly available<sup>47</sup>. These simulated data were used to  
353 construct the training and validation sets<sup>48</sup> to evaluate network training, particle  
354 localization, and model generalizability.

355

## 356 **Convolutional Network Architecture**

357 To encode volumetric cryo-ET data as embedding vectors in a high-dimensional space,  
358 TomoTwin employs a 3D CNN consisting of five convolutional blocks followed by a head  
359 network ([Supplementary Fig. 1a](#)). Each convolution block consists of two 3D  
360 convolutional layers with a kernel size of 3x3x3. Each convolutional layer is followed by  
361 a normalization layer and a leaky rectified linear (ReLU) activation function. In the first  
362 convolutional layer of each convolutional block, the number of output channels is twice  
363 the input channels and in the second convolutional layer the number of output channels  
364 matches the output from the previous layer. Max pooling is performed with a kernel size  
365 of 2x2x2 after the first convolutional block and adaptive max pooling to a size of 2x2x2 is  
366 performed after the final convolutional block. As a result, when provided with a 37x37x37  
367 subvolume with 1 channel as a normalized, 37x37x37x1 array, the convolutional blocks

368 transform the input to a 2x2x2x1024 feature vector which is then fed to the head network.

369 In the head network, the feature vector is first flattened channel-wise before being subject  
370 to a dropout layer and then passed through a series of fully-connected layers that  
371 transform the flattened vector to a 1-dimensional, 32-length feature vector. Finally, this  
372 feature vector is L2-normalized to yield an output embedding vector for the subvolume.

373

### 374 **Triplet Generation**

375 TomoTwin is trained on triplets of subvolumes consisting of an anchor volume A, a  
376 positive volume P, and a negative volume N ([Supplementary Fig. 1b](#)). Each subvolume  
377 is assigned to a data class corresponding to the macromolecule contained within and has  
378 a size of 37x37x37 voxels. Triplets are constructed where A and P are sampled from the  
379 same data class and N from a different data class. Given a distance function D and an  
380 embedding function f, the triplet loss is defined as:

$$381 \quad L(A,P,N) = \max(D(f(A), f(P)) - D(f(A), f(N)) + \alpha, 0)$$

382 where the hyperparameter  $\alpha$  is the margin value. As distance function D we use cosine  
383 similarity which is defined as

$$384 \quad D(Q,P) = \frac{Q \cdot P}{\|Q\| \times \|P\|}$$

385 where Q and P are arbitrary embedding vectors,  $\cdot$  is the dot product and  $\|\cdot\|$  the length of  
386 the vector. During training, triplets are generated by online semihard triplet mining  
387 wherein a batch of subvolumes are embedded and triplets generated automatically with  
388 the negative subvolume embedding being selected from those only with a distance to the

389 anchor greater than the positive subvolume embedding but not greater than a margin  
390  $\alpha_{miner}$ :

$$391 \quad D(a, p) < D(a, n) < d(a, p) + \alpha_{miner}$$

392 Where a, p and n are the embedding vectors of the anchor, positive and negative  
393 respectively and  $\alpha_{miner}$  is the margin of the miner.

394

### 395 **Training of the General Picking Model**

396 Training of the 3D CNN was performed for 600 epochs using an adaptive moment  
397 estimation (ADAM) optimizer<sup>49</sup>. The model from the epoch with the best F1 score on the  
398 subvolumes in the validation set was further evaluated in the localization and  
399 generalization tasks and used as the general picking model.

### 400 **Data augmentation**

401 To prevent overfitting during training and to improve generalization of the model, online  
402 data augmentations were applied to each normalized volume before its embedding was  
403 calculated including rotation, dropout, translation, and the addition of noise. For the  
404 rotation augmentation, subvolumes were rotated by a random angle in the X-Y plane but  
405 not X-Z or Y-Z to prevent reorientation of the missing wedge. In the dropout augmentation,  
406 a random portion between 5 and 20% of the voxels were set to the subvolume mean  
407 value. In the translation augmentation, the subvolume was shifted by 1-2 pixels in each  
408 direction. The addition of noise augmentation added Gaussian noise with a randomly  
409 chosen standard deviation between 0 and 0.3 to the subvolume.

410

## 411 **Hyperparameter optimization**

412 The training of modern convolutional neural networks involves the selection of many  
413 hyperparameters, some of these choices affect the architecture while others affect the  
414 learning process itself. While some heuristics exist to guide hyperparameter selection,  
415 finding a combination of settings that maximize the utility of a machine learning tool by  
416 hand quickly becomes intractable. Optuna<sup>50</sup> was applied to explore the hyperparameter  
417 search space and identify an optimized set of parameters for training . Models were  
418 trained on a subset of the training data for 200 epochs and the F1 score calculated on  
419 the validation set after each epoch. Pruning was performed after 50 epochs for training  
420 runs with an F1 score lower than the global median. In total, searches were applied for  
421 the hyperparameters of learning rate, dropout rate, optimizer, batch size, weight decay,  
422 size of the first convolution kernel, number of output layer nodes, online triplet mining  
423 strategy (semihard<sup>51</sup>, easyhard<sup>52</sup> , none), normalization type (group norm<sup>53</sup>, batch  
424 norm<sup>54</sup>), loss function (TripletLoss<sup>31</sup>, SphereFace<sup>55</sup>, ArcFace<sup>56</sup>), and loss margin  
425 ([Supplementary Fig. 7](#)).

426 Most notably and unexpectedly, the type of normalization applied during training  
427 was the largest overall affecter of performance with group normalization<sup>53</sup> outperforming  
428 the more common batch normalization<sup>54</sup> strategy ([Supplementary Fig. 7b](#)). Additionally  
429 noted was the increased performance of a standard triplet loss function over the  
430 theoretically superior SphereFace<sup>55</sup> and ArcFace<sup>56</sup> loss functions ([Supplementary Fig.](#)  
431 [7c](#)). These findings underpin the necessity to explore a wide range of hyperparameters  
432 during training as heuristics alone are not enough to guide optimal hyperparameter  
433 selection for the training of modern convolutional neural networks.

434

### 435 **Particle picking workflow with the general model**

436 For each dataset picked with the general model, first all tomograms were embedded. To  
437 achieve this, the tomograms were subdivided into a series of overlapping 37x37x37  
438 subvolumes with a stride of 2 voxels. For the reference-based workflow, a random particle  
439 for each protein of interest was selected as reference and embedded to generate a target  
440 embedding. The tomogram and target embeddings were provided to TomoTwin Map  
441 which calculated the distance matrix between each target embedding and each  
442 subvolume embedding from the tomogram and returned this along with a similarity map  
443 for each target embedding. This matrix was then provided to TomoTwin Locate which  
444 identified areas of high confidence as target locations using a region-growing based  
445 maximum detection procedure followed by non-maxima suppression. The returned  
446 candidate positions were then subject to confidence and size thresholding with TomoTwin  
447 pick to produce final coordinates for each protein of interest.

448

### 449 **Evaluation of simulated data**

450 The performance of particle localization was calculated from three metrics: recall,  
451 precision, and, the harmonic mean of the two, the F1 score which are defines as:

$$452 \text{ precision} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$$

$$453 \text{ recall} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$

$$F1 = 2 \frac{\textit{precision} \cdot \textit{recall}}{\textit{precision} + \textit{recall}}$$

454

455 Selected particle locations counted as true positives if the intersection over union (IOU)  
456 of the box of the selected particle location and the ground truth box was greater than 0.6.  
457 The IOU is defined as the ratio of the intersecting volume of two bounding boxes and the  
458 volume of their union.

459 The particle localization accuracy of the trained model was assessed for each tomogram  
460 in the validation set ([Supplementary Fig. 4a](#)). To test model generalization, the  
461 localization task was performed on a tomogram containing 7 proteins not included in the  
462 training set for which TomoTwin was therefore naïve ([Fig. 2](#)).

463

## 464 **Clustering**

465 For clustering analysis, a random sample of 400,000 embeddings from the high-  
466 dimensional tomogram embeddings were fit to a uniform 2D manifold with Uniform  
467 Manifold Approximation (UMAP) with GPU-acceleration provided by the RAPIDS  
468 package<sup>57</sup>. The UMAP model was used as the basis to transform the entire tomogram  
469 embeddings and the results plotted ([Figure 5a,c](#)). Clusters were identified by eye and  
470 selected by drawing a closed shape containing the desired points. The enclosed points  
471 were then traced back to their original high-dimensional embeddings and the average  
472 embedding of them was calculated. This average embedding was then used as a target  
473 embedding for classification, localization, and picking in the same manner as for the  
474 reference-based workflow.

475

## 476 **Preparation of Experimental Samples**

477 The components of the mixture were either thawed from long-term storage at -80 °C or  
478 freshly prepared. *Photorhabdus luminescens* holotoxin was expressed, purified and the  
479 holotoxin formed as described previously<sup>58</sup> and used at a stock concentration of 0.49  
480 mg/mL. RhsA from *Pseudomonas protegens* was expressed and purified as described  
481 previously<sup>34</sup> and used at 4 mg/mL concentration. Liposomes were prepared by extrusion.  
482 4 mg/mL of each POPC (1-palmitoyl-2-oleoyl-glycero-3-phosphocholine, Avanti Polar  
483 Lipids) and DOPS (1,2-dioleoyl-sn-glycero-3-phospho-L-serine, Avanti Polar Lipids) were  
484 mixed in buffer (50 mM Tris, pH 8, 150 NaCl, 0.05% Tween20) and after brief sonication  
485 (1 min in water bath) and three cycles of freeze-thawing (-196 °C and 50 °C), the liposome  
486 solution was passed 11 times through a polycarbonate membrane with a 400 nm pore  
487 size in a mini extruder (Avanti Polar Lipids). Total lipid concentration was diluted with  
488 buffer to 0.16 mg/mL. The freeze-dried content of one vial Tobacco mosaic virus (TMV)  
489 (DSMZ GmbH Braunschweig, Germany, PC-0107) was solved in 1 mL buffer and diluted  
490 500 times as working solution. The Apoferritin (ApoF) plasmid was a kind gift by Dr.  
491 Christos Savva (Midlands Regional Cryo-Electron Microscopy Facility). Expression and  
492 purification of ApoF was optimized based on the protocol described earlier<sup>59</sup> and final  
493 concentration of frozen stock was 3 mg/mL.

494 Different ratios of the mixture were prepared and then examined after vitrification using  
495 cryo-EM. For cryo-ET, a ratio of 1:2:2:20:10 (TMV:ApoF:Liposomes:TcToxin:RhsA)  
496 was chosen.

497

## 498 **Grid Preparation**

499 Grids were prepared using a Vitrobot Mark IV (Thermo Fisher Scientific) at 4 °C and 100%  
500 humidity. 4 µL of the freshly prepared mixture were applied to glow-discharged (Quorum  
501 GloQube) R1.2/1.4 Cu 200 (Quantifoil) grids. After blotting (3.5 s at blot force -1, no drain  
502 time) the specimen was vitrified in liquid ethane.

503

## 504 **Cryo-ET**

505 Grids of different mixing ratios were screened using a Talos Arctica electron microscope  
506 (Thermo Fisher Scientific) equipped with a X-FEG and Falcon 3 camera. Small datasets  
507 of 100-200 images were collected using the software EPU (Thermo Fisher Scientific). The  
508 best specimen was transferred to a Titan Krios G3 electron microscope equipped with X-  
509 FEG. Images were recorded on a K3 camera (Gatan) operated in counting mode at a  
510 nominal magnification of 63,000, resulting in a pixel size of 1.484 Å/pix. A Bioquantum  
511 post-column energy (Gatan) was used for zero loss imaging with a slit width of 20 eV.  
512 Tilt series were acquired using SerialEM<sup>60</sup> with the Plugin PACEtomo<sup>61</sup> and with a dose  
513 symmetric tilt scheme<sup>62</sup> from 60° to 60° with a step size of 3°. Each movie was collected  
514 as an exposure of 0.2 seconds subdivided into 10 frames. Frames were then exported to  
515 Warp 1.0.9<sup>26</sup> for motion correction, CTF estimation and generation of tilt series. Tilt series  
516 were aligned with patch tracking and tomograms reconstructed by weighted back-  
517 projection in IMOD<sup>47</sup> with a pixel size of 5.936. Tomograms were scaled by Fourier  
518 shrinking to 10 Å/pix for embedding with TomoTwin.



519 Raw frames of *M. pneumoniae* cells were downloaded from EMPIAR (EMPIAR-10499).  
520 Motion correction and CTF estimation were performed in Warp 1.0.9 which was then used  
521 to generate tilt series. These tilt series were aligned with patch tracking and tomograms  
522 reconstructed by weighted back-projection in IMOD with a pixel size of 6.802 Å/pix.  
523 Tomograms were then scaled by Fourier shrinking to 13.6 Å/pix for embedding with  
524 TomoTwin.

525

### 526 **Evaluation of experimental data**

527 For tomograms from samples prepared in-house, coordinates of particles identified with  
528 TomoTwin were scaled to a pixel size of 5.936 to match the originally reconstructed  
529 tomograms. The tomograms were imported and these coordinates were used to extract  
530 subtomograms in Relion 3.0<sup>37</sup>. For reference-free analysis, 3D subtomograms were  
531 projected to 2D with SPHIRE<sup>36</sup> and then used for 2D classification. For tomograms  
532 attained from EMPIAR, coordinates of particles identified with TomoTwin were scaled to  
533 a pixel size of 6.802 Å/pix to match the originally reconstructed tomograms. The  
534 tomograms were imported and coordinates were imported and used to reconstruct  
535 pseudo-subtomograms in Relion 4.0<sup>40</sup>. A reference was created from a 70S ribosome  
536 ([EMD-11650](#)) by lowpass filtering to 30 Å and then scaling the pixel size to 6.802 Å/pix.  
537 This reference was used for 3D classification with the pseudo-subtomograms in Relion  
538 4.0.

539

### 540 **Hardware**

541 Two computational setups were utilized for calculations, a distributed computing system  
542 and a local workstation. The distributed computing system consisted of the Max Planck  
543 Gesellschaft Supercomputer ‘Raven’ using up to 30 Nvidia A100 GPUs, where each GPU  
544 has 40 GB memory. Each process had 18 cores of Intel Xeon IceLake-SP 8360Y  
545 processors and 128GB system memory available. The local workstation consisted of a  
546 local unit equipped with a Nvidia Titan V (12 GB memory) GPU and a Intel i9-7920X CPU  
547 with 64 GB system memory.

548

549 Hyperparameter optimization was done in parallel for 7 days on the distributed  
550 computing setup and embeddings were calculated on this set up as well using 2 GPUs.

551 In all cases a box size of 37 and stride of 2 were used for embedding.

552 The inhouse workstation was used for miscellaneous tasks and for calculating timings  
553 using 2 GPUs.

554

## 555 **Timings**

556 The calculation of the embeddings is the only function of TomoTwin requiring significant  
557 processing time. To measure this, we embedded our largest experimental tomogram  
558 (608x855x148 after Fourier shrinking) on a local workstation and a distributed computing  
559 system. Using 2 GPUs, tomogram embedding took 80 minutes for the local setup and 30  
560 minutes for the distributed setup, corresponding to the total time to pick all proteins of  
561 interest per tomogram on each setup.

562

563

## 564 **Data Availability**

565 All simulated tomograms used in this study are available here:  
566 <https://doi.org/10.5281/zenodo.6637357>.

567 The extracted subvolumes used to train and evaluate the performance of TomoTwin are  
568 available at: <https://doi.org/10.5281/zenodo.6637456>.

569 The TEM-Simulator-Scripts package used for automated tilt-series simulation and  
570 reconstruction is available at: <https://github.com/MPI-Dortmund/tem-simulator-scripts>.

571 TomoTwin is available under an open-source license at: [https://github.com/MPI-](https://github.com/MPI-Dortmund/tomotwin-cryoet)  
572 [Dortmund/tomotwin-cryoet](https://github.com/MPI-Dortmund/tomotwin-cryoet).

573

## 574 **Author contributions**

575 Conceptualization: T.W. and S.R.;

576 Software implementation: T.W., G.R., M.S.;

577 Software – Testing: G.R., T.W., M.S.;

578 Formal analysis: G.R., T.W.;

579 Supervision: S.R.;

580 Writing – original draft: G.R., T.W.;

581 Writing – review and editing: G.R., T.W., and S.R.;

582 Funding acquisition: S.R.;

583

## 584 **Acknowledgements:**

585 We thank D. Prumbaum for collecting cryo-ET data collection, P. Günther for providing  
586 RhsA and liposomes, P. Njenga Ng` Ang` A for providing TcdA1, and K. Vogel-Bachmayr

587 for purifying apoferritin, C. Savva for providing the apoferritin plasmid and A. Prajica for  
588 support in figure and logo creation. The work was supported by the Max Planck Society  
589 (to S.R.).  
590

591 **References:**

- 592 1. Wan, W. & Briggs, J. A. G. Cryo-Electron Tomography and Subtomogram Averaging. in  
593 *Methods in Enzymology* vol. 579 329–367 (Elsevier, 2016).
- 594 2. Lučić, V., Rigort, A. & Baumeister, W. Cryo-electron tomography: The challenge of doing  
595 structural biology in situ. *Journal of Cell Biology* **202**, 407–419 (2013).
- 596 3. Bharat, T. A. M. & Scheres, S. H. W. Resolving macromolecular structures from electron  
597 cryo-tomography data using subtomogram averaging in RELION. *Nat Protoc* **11**, 2054–2065  
598 (2016).
- 599 4. Koning, R. I., Koster, A. J. & Sharp, T. H. Advances in cryo-electron tomography for biology  
600 and medicine. *Ann Anat* **217**, 82–96 (2018).
- 601 5. Asano, S., Engel, B. D. & Baumeister, W. In Situ Cryo-Electron Tomography: A Post-  
602 Reductionist Approach to Structural Biology. *Journal of Molecular Biology* **428**, 332–343  
603 (2016).
- 604 6. Wang, Z. *et al.* Structures from intact myofibrils reveal mechanism of thin filament  
605 regulation through nebulin. *Science* **375**, eabn1934 (2022).
- 606 7. Wang, Z. *et al.* The molecular basis for sarcomere organization in vertebrate skeletal  
607 muscle. *Cell* **184**, 2135-2150.e13 (2021).
- 608 8. Schaffer, M. *et al.* Cryo-focused Ion Beam Sample Preparation for Imaging Vitreous Cells by  
609 Cryo-electron Tomography. *Bio Protoc* **5**, e1575 (2015).
- 610 9. Wagner, F. R. *et al.* Preparing samples from whole cells using focused-ion-beam milling for  
611 cryo-electron tomography. *Nat Protoc* **15**, 2041–2070 (2020).

- 612 10. Tacke, S. *et al.* A streamlined workflow for automated cryo focused ion beam milling.  
613 *Journal of Structural Biology* **213**, 107743 (2021).
- 614 11. Sutton, G. *et al.* Assembly intermediates of orthoreovirus captured in the cell. *Nat Commun*  
615 **11**, 4445 (2020).
- 616 12. Schaffer, M. *et al.* A cryo-FIB lift-out technique enables molecular-resolution cryo-ET within  
617 native *Caenorhabditis elegans* tissue. *Nat Methods* **16**, 757–762 (2019).
- 618 13. Li, M., Ma, J., Li, X. & Sui, S.-F. In situ cryo-ET structure of phycobilisome–photosystem II  
619 supercomplex from red alga. *eLife* **10**, e69635 (2021).
- 620 14. Burbaum, L. *et al.* Molecular-scale visualization of sarcomere contraction within native  
621 cardiomyocytes. *Nat Commun* **12**, 4086 (2021).
- 622 15. Mosalaganti, S. *et al.* AI-based structure prediction empowers integrative structural analysis  
623 of human nuclear pores. *Science* **376**, eabm9506 (2022).
- 624 16. Schwartz, T. U. Solving the nuclear pore puzzle. *Science* **376**, 1158–1159 (2022).
- 625 17. Zhu, X. *et al.* Structure of the cytoplasmic ring of the *Xenopus laevis* nuclear pore complex.  
626 *Science* **376**, eabl8280 (2022).
- 627 18. Gubins, I. *et al.* SHREC 2021: Classification in Cryo-electron Tomograms. *Eurographics*  
628 *Workshop on 3D Object Retrieval* 13 pages (2021) doi:10.2312/3DOR.20211307.
- 629 19. Moebel, E. *et al.* Deep learning improves macromolecule identification in 3D cellular cryo-  
630 electron tomograms. *Nat Methods* **18**, 1386–1394 (2021).
- 631 20. Wu, S., Liu, G. & Yang, G. Fast Particle Picking For Cryo-Electron Tomography Using One-  
632 Stage Detection. in *2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI)*  
633 1–5 (IEEE, 2022). doi:10.1109/ISBI52829.2022.9761580.

- 634 21. Hao, Y. *et al.* VP-Detector: A 3D multi-scale dense convolutional neural network for  
635 macromolecule localization and classification in cryo-electron tomograms. *Comput Methods*  
636 *Programs Biomed* **221**, 106871 (2022).
- 637 22. Frangakis, A. S. *et al.* Identification of macromolecular complexes in cryoelectron  
638 tomograms of phantom cells. *Proc. Natl. Acad. Sci. U.S.A.* **99**, 14153–14158 (2002).
- 639 23. Lucas, B. A. *et al.* Locating macromolecular assemblies in cells by 2D template matching  
640 with cisTEM. *eLife* **10**, e68946 (2021).
- 641 24. Balyschew, N. *TomoBEAR: an automated, configurable and customizable full processing*  
642 *pipeline for tomographic cryo electron microscopy data and subtomogram averaging.*  
643 (2021).
- 644 25. Wagner, T. *et al.* SPHIRE-crYOLO is a fast and accurate fully automated particle picker for  
645 cryo-EM. *Commun Biol* **2**, 218 (2019).
- 646 26. Tegunov, D. & Cramer, P. Real-time cryo-electron microscopy data preprocessing with  
647 Warp. *Nat Methods* **16**, 1146–1152 (2019).
- 648 27. Bepler, T. *et al.* Positive-unlabeled convolutional neural networks for particle picking in  
649 cryo-electron micrographs. *Nat Methods* **16**, 1153–1160 (2019).
- 650 28. Wagner, T. & Raunser, S. The evolution of SPHIRE-crYOLO particle picking and its application  
651 in automated cryo-EM processing workflows. *Commun Biol* **3**, 61 (2020).
- 652 29. Kaya & Bilge. Deep Metric Learning: A Survey. *Symmetry* **11**, 1066 (2019).
- 653 30. Ghojogh, B., Ghodsi, A., Karray, F. & Crowley, M. Spectral, Probabilistic, and Deep Metric  
654 Learning: Tutorial and Survey. (2022) doi:10.48550/ARXIV.2201.09267.

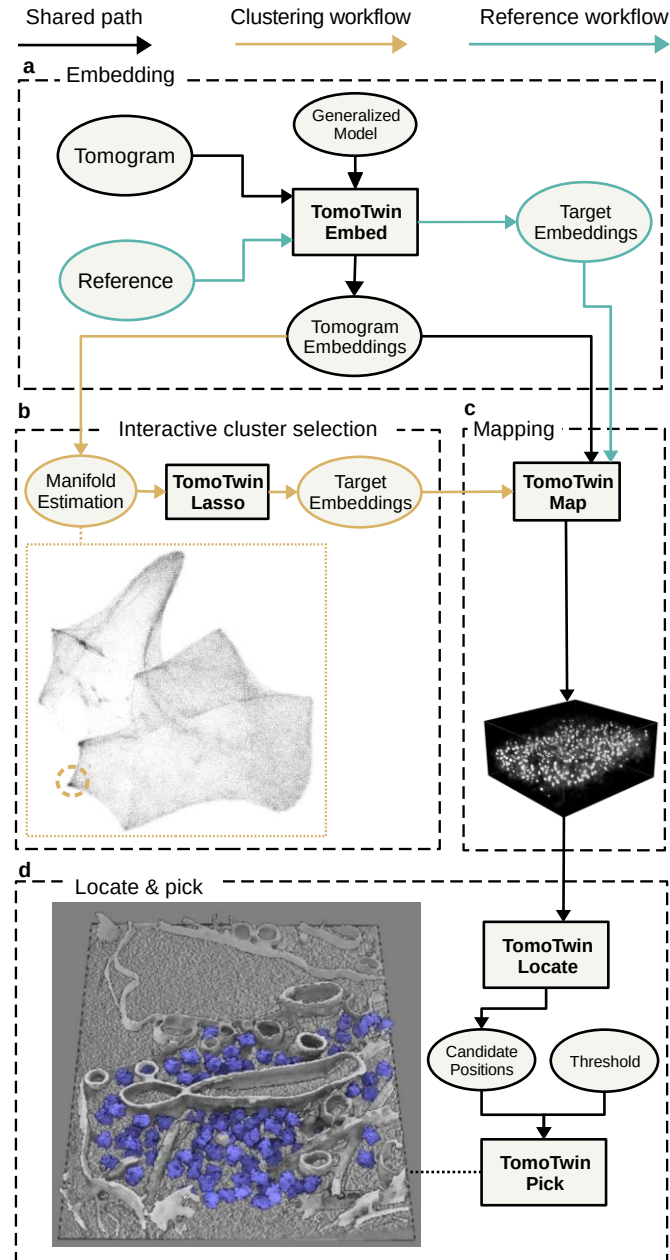
- 655 31. Schroff, F., Kalenichenko, D. & Philbin, J. FaceNet: A unified embedding for face recognition  
656 and clustering. in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*  
657 815–823 (IEEE, 2015). doi:10.1109/CVPR.2015.7298682.
- 658 32. Liu, J., Deng, Y., Bai, T., Wei, Z. & Huang, C. Targeting Ultimate Accuracy: Face Recognition  
659 via Deep Embedding. (2015) doi:10.48550/ARXIV.1506.07310.
- 660 33. Dai, G., Xie, J. & Fang, Y. Deep Correlated Holistic Metric Learning for Sketch-Based 3D  
661 Shape Retrieval. *IEEE Trans. on Image Process.* **27**, 3374–3386 (2018).
- 662 34. Zou, W. *et al.* Expression, purification, and characterization of recombinant human H-chain  
663 ferritin. *Preparative Biochemistry & Biotechnology* **46**, 833–837 (2016).
- 664 35. Günther, P. *et al.* Structure of a bacterial Rhs effector exported by the type VI secretion  
665 system. *PLoS Pathog* **18**, e1010182 (2022).
- 666 36. Gatsogiannis, C. *et al.* Tc toxin activation requires unfolding and refolding of a  $\beta$ -propeller.  
667 *Nature* **563**, 209–213 (2018).
- 668 37. Moriya, T. *et al.* High-resolution Single Particle Analysis from Electron Cryo-microscopy  
669 Images Using SPHIRE. *JoVE* 55448 (2017) doi:10.3791/55448.
- 670 38. Zivanov, J. *et al.* New tools for automated high-resolution cryo-EM structure determination  
671 in RELION-3. *eLife* **7**, e42166 (2018).
- 672 39. 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 40. Kimanius, D., Dong, L., Sharov, G., Nakane, T. & Scheres, S. H. W. New tools for automated  
676 cryo-EM single-particle analysis in RELION-4.0. *Biochemical Journal* **478**, 4169–4185 (2021).



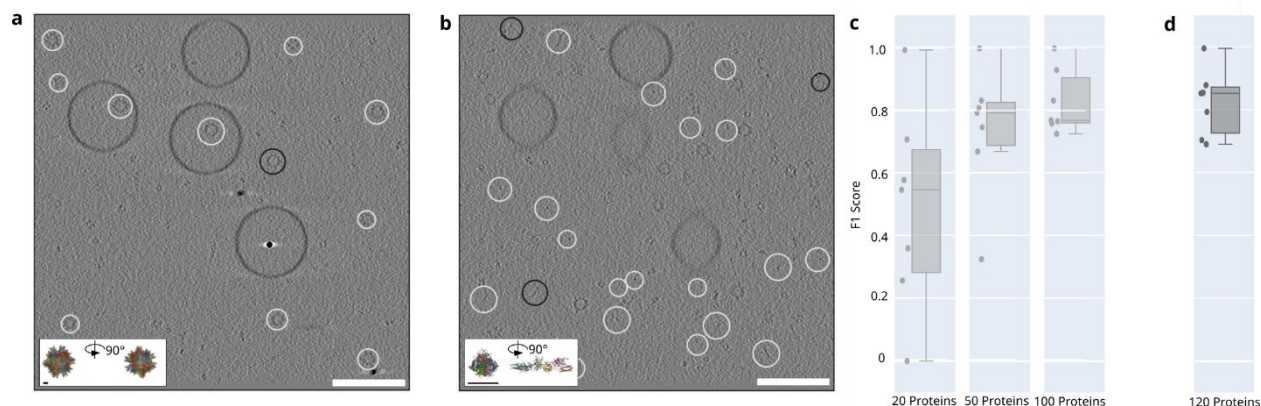
- 677 41. Gill, H. S. & Eisenberg, D. The Crystal Structure of Phosphinothricin in the Active Site of  
678 Glutamine Synthetase Illuminates the Mechanism of Enzymatic Inhibition. *Biochemistry* **40**,  
679 1903–1912 (2001).
- 680 42. Stabrin, M. *et al.* TransPHIRE: automated and feedback-optimized on-the-fly processing for  
681 cryo-EM. *Nat Commun* **11**, 5716 (2020).
- 682 43. Schenk, A. D., Cavadini, S., Thomä, N. H. & Genoud, C. Live Analysis and Reconstruction of  
683 Single-Particle Cryo-Electron Microscopy Data with CryoFLARE. *J. Chem. Inf. Model.* **60**,  
684 2561–2569 (2020).
- 685 44. Maluenda, D. *et al.* Flexible workflows for on-the-fly electron-microscopy single-particle  
686 image processing using *Scipion*. *Acta Crystallogr D Struct Biol* **75**, 882–894 (2019).
- 687 45. Campello, R. J. G. B., Moulavi, D. & Sander, J. Density-Based Clustering Based on  
688 Hierarchical Density Estimates. in *Advances in Knowledge Discovery and Data Mining* (eds.  
689 Pei, J., Tseng, V. S., Cao, L., Motoda, H. & Xu, G.) vol. 7819 160–172 (Springer Berlin  
690 Heidelberg, 2013).
- 691 46. Rullgård, H., Öfverstedt, L.-G., Masich, S., Daneholt, B. & Öktem, O. Simulation of  
692 transmission electron microscope images of biological specimens: SIMULATION OF TEM  
693 IMAGES OF BIOLOGICAL SPECIMENS. *Journal of Microscopy* **243**, 234–256 (2011).
- 694 47. Kremer, J. R., Mastronarde, D. N. & McIntosh, J. R. Computer visualization of three-  
695 dimensional image data using IMOD. *J Struct Biol* **116**, 71–76 (1996).
- 696 48. Wagner, Thorsten, Rice, Gavin, Stabrin, Markus & Raunser, Stefan. Simulated Tomograms  
697 used for TomoTwin. (2022) doi:10.5281/ZENODO.6637357.

- 698 49. Wagner, Thorsten, Rice, Gavin, Stabrin, Markus & Raunser, Stefan. Training and validation  
699 data used to produce the pre-trained model for the TomoTwin paper. (2022)  
700 doi:10.5281/ZENODO.6637456.
- 701 50. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. (2017).
- 702 51. Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A Next-generation  
703 Hyperparameter Optimization Framework. in *Proceedings of the 25th ACM SIGKDD  
704 International Conference on Knowledge Discovery & Data Mining* 2623–2631 (ACM, 2019).  
705 doi:10.1145/3292500.3330701.
- 706 52. Musgrave, K., Belongie, S. & Lim, S.-N. PyTorch Metric Learning. (2020)  
707 doi:10.48550/ARXIV.2008.09164.
- 708 53. Xuan, H., Stylianou, A. & Pless, R. Improved Embeddings with Easy Positive Triplet Mining.  
709 (2019) doi:10.48550/ARXIV.1904.04370.
- 710 54. Wu, Y. & He, K. Group Normalization. (2018) doi:10.48550/ARXIV.1803.08494.
- 711 55. Ioffe, S. & Szegedy, C. Batch Normalization: Accelerating Deep Network Training by  
712 Reducing Internal Covariate Shift. (2015) doi:10.48550/ARXIV.1502.03167.
- 713 56. Liu, W. *et al.* SphereFace: Deep Hypersphere Embedding for Face Recognition. (2017)  
714 doi:10.48550/ARXIV.1704.08063.
- 715 57. Deng, J., Guo, J., Xue, N. & Zafeiriou, S. ArcFace: Additive Angular Margin Loss for Deep Face  
716 Recognition. (2018) doi:10.48550/ARXIV.1801.07698.
- 717 58. RAPIDS Development Team. *RAPIDS: Collection of Libraries for End to End GPU Data  
718 Science*. (2018).

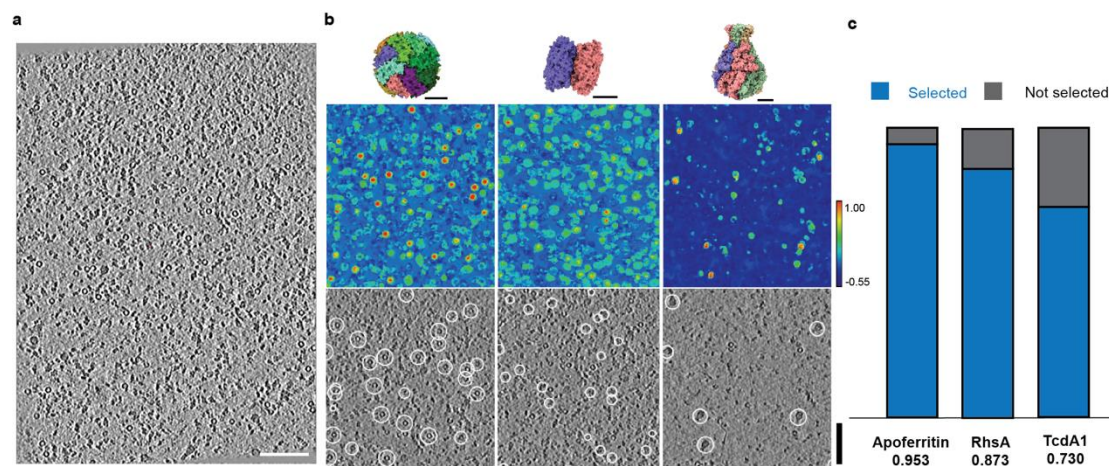
- 719 59. Roderer, D., Hofnagel, O., Benz, R. & Raunser, S. Structure of a Tc holotoxin pore provides  
720 insights into the translocation mechanism. *Proc Natl Acad Sci U S A* **116**, 23083–23090  
721 (2019).
- 722 60. Mastronarde, D. N. Automated electron microscope tomography using robust prediction of  
723 specimen movements. *Journal of Structural Biology* **152**, 36–51 (2005).
- 724 61. Eisenstein, F. *et al.* *Parallel cryo electron tomography on in situ lamellae*.  
725 <http://biorxiv.org/lookup/doi/10.1101/2022.04.07.487557> (2022)  
726 doi:10.1101/2022.04.07.487557.
- 727 62. Hagen, W. J. H., Wan, W. & Briggs, J. A. G. Implementation of a cryo-electron tomography  
728 tilt-scheme optimized for high resolution subtomogram averaging. *Journal of Structural*  
729 *Biology* **197**, 191–198 (2017).
- 730



**Fig. 1: TomoTwin identifies and localizes particles by a clustering or a reference-based workflow.** **a**, The first step in using TomoTwin is to embed the tomogram with the pre-trained model. Optionally, references can be selected and embedded as well to create target embeddings. **b**, For the clustering workflow the tomogram embeddings are projected on a 2D manifold and an interactive lasso tool is used to select clusters of interest to generate target embeddings. **c**, The distance matrix between each target embedding and the embeddings of the tomogram is calculated. **d**, All local maxima are located with TomoTwin Locate and are used to pick final coordinates for each protein of interest using TomoTwin Pick with confidence and size thresholding.

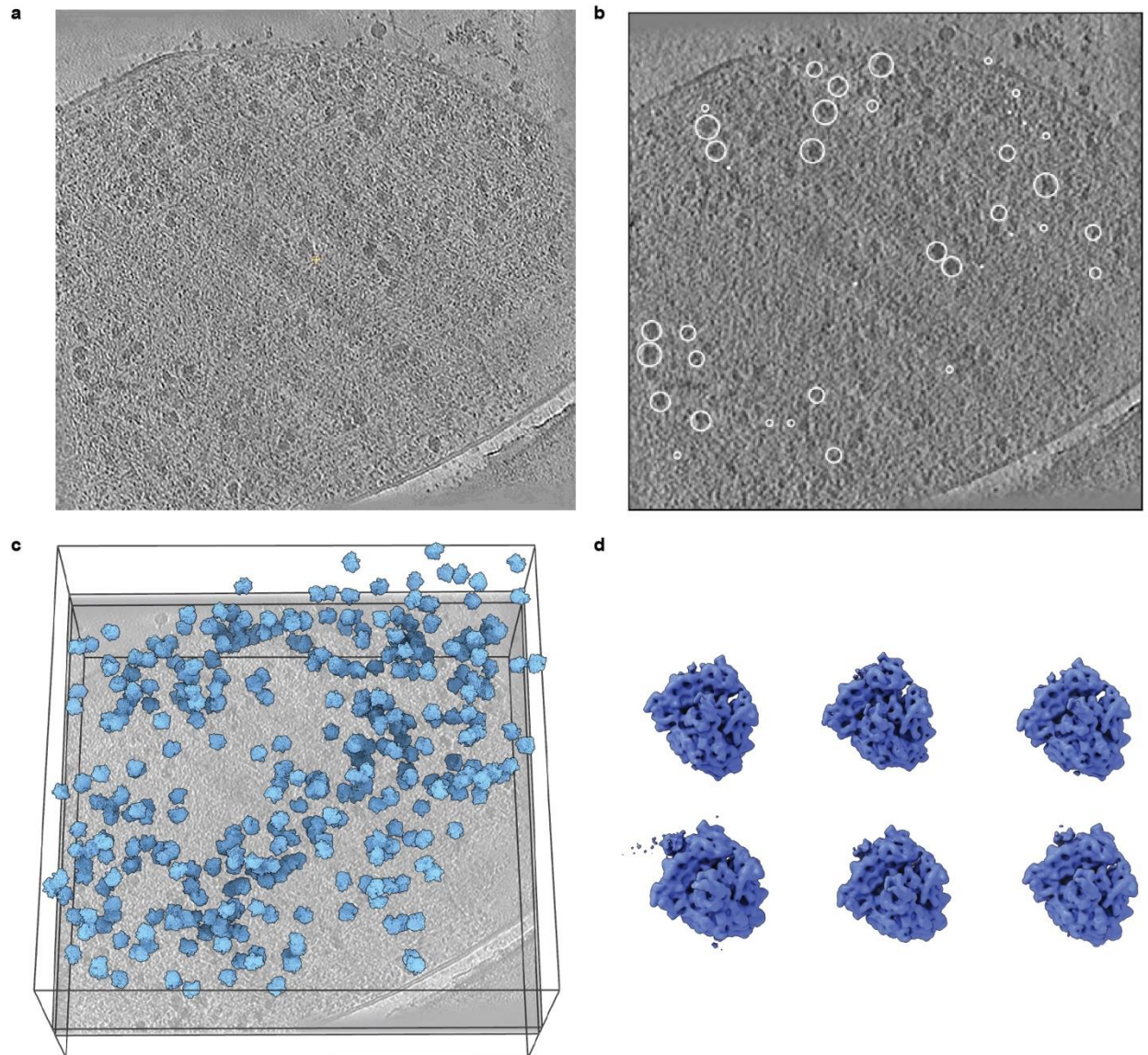


**Fig. 2: TomoTwin generalizes to novel proteins and locates them accurately.** **a**, True positive selected particles (white) and false negative (black) of the largest protein 2DF7 (896 kDa) and **b**, the smallest protein 1FZG (142 kDa) in the generalization tomogram. The F1 scores are 0.99 and 0.88 for 2DF7 and 1FZG respectively. **c**, With increasing number of proteins used during training the mean F1 score on the generalization tomogram increased as well. The mean F1 scores are 0.49, 0.73, 0.82 for a model trained on 20, 50 and 100 proteins respectively. **d**, The model trained on the full training set of 120 proteins reached a mean F1 score of 0.82 but has the highest median F1 score of 0.85. White scale bar 100 nm, black scale bar 5 nm

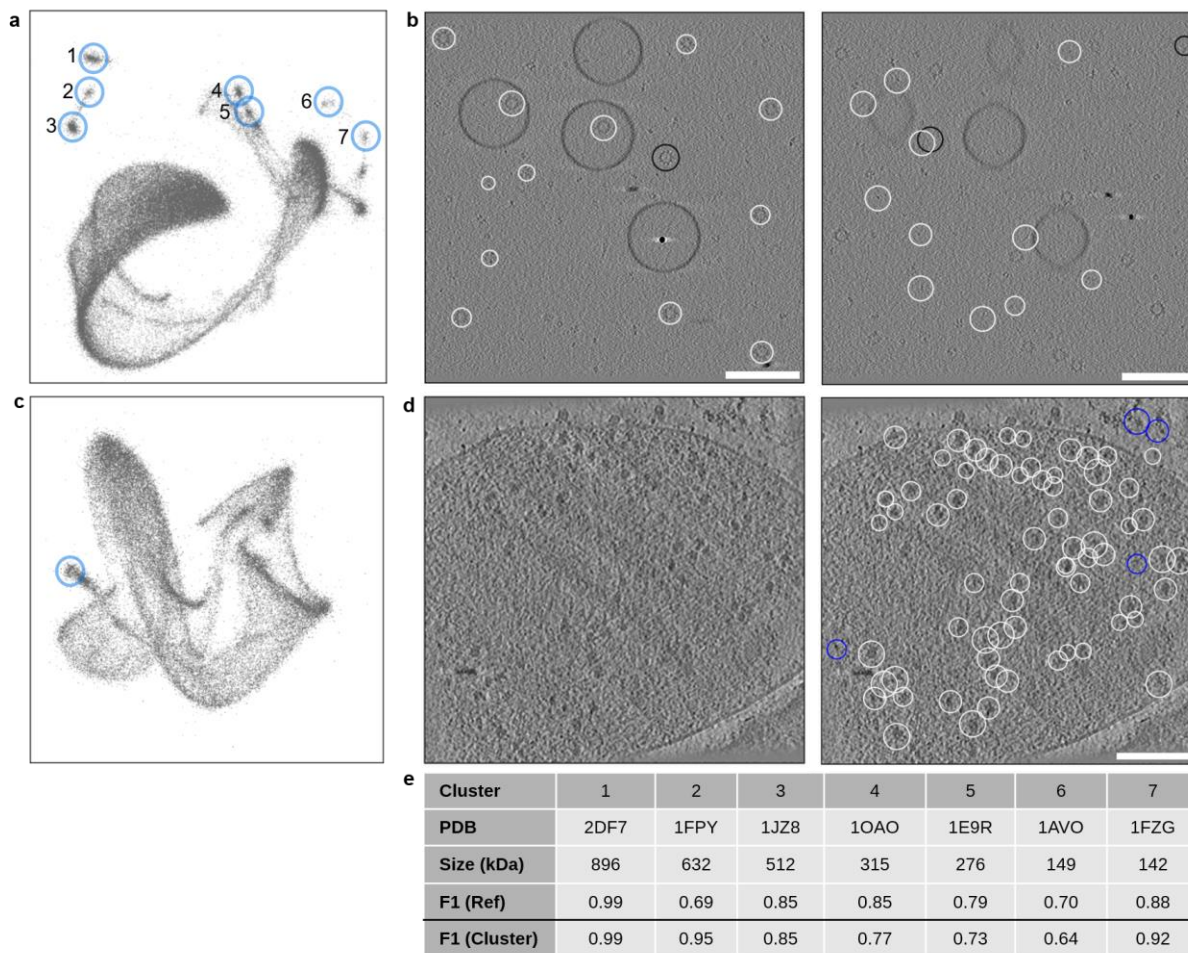


**Fig. 3: TomoTwin accurately localizes multiple proteins simultaneously in crowded tomograms.** **a**, Representative slice of a tomogram containing a mixture of apoferritin, RhsA, and TcdA1; scale bar: 100 nm. **b**, Protein structure, cosine similarity map between tomogram and each target, and representative picking for apoferritin (PDB ID: 1DAT), RhsA (PDB ID: 7Q97), and TcdA1 (PDB ID: 6L7E) respectively. Scale bar for protein structures: 5 nm, scale bar for tomograms: 100 nm, color bar: -0.55 - 1.00 **c**, Proportion of picked subvolumes contained within positive 2D classes. Total subvolumes picked: apoferritin: 848, RhsA: 577, TcdA1: 122.



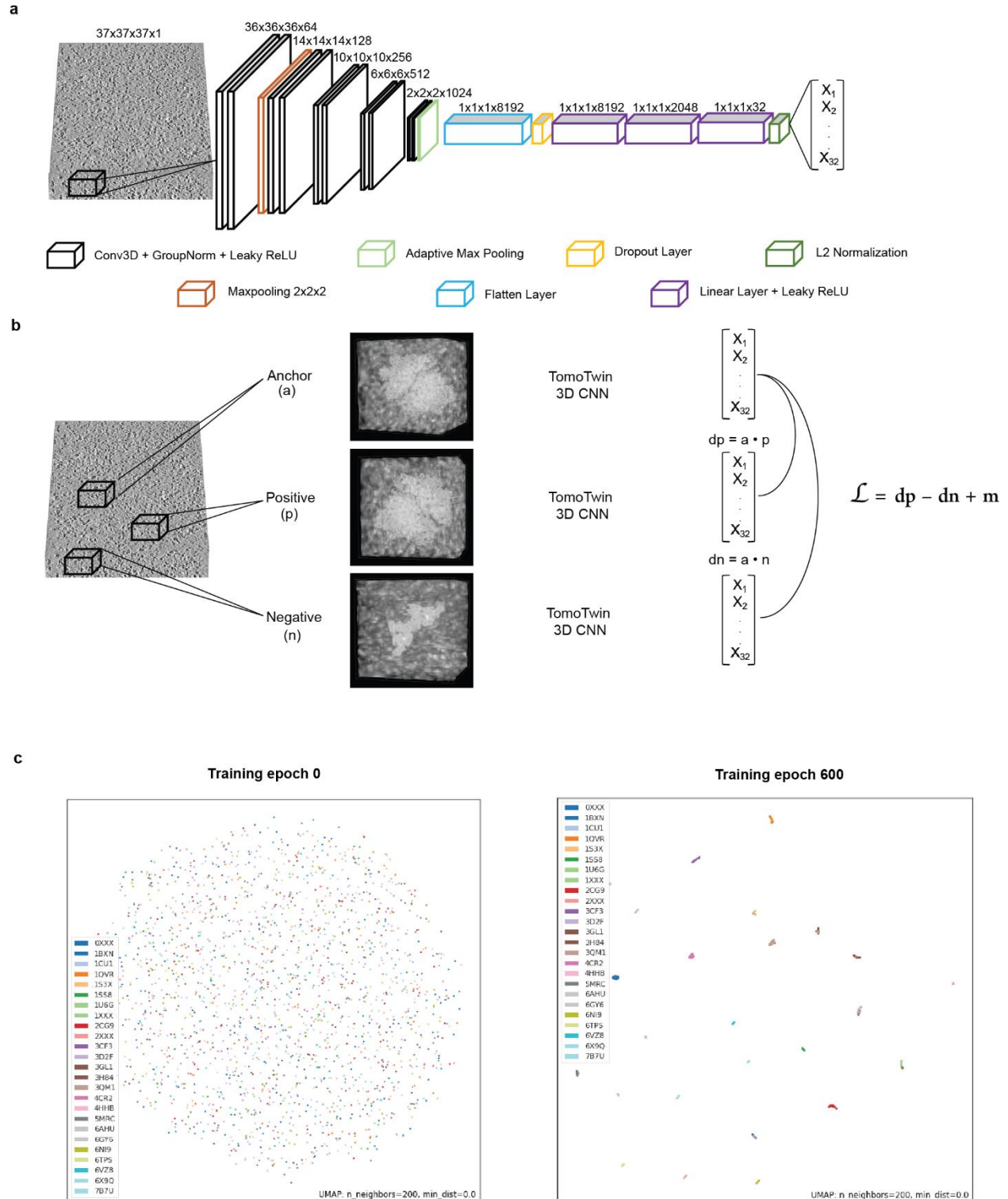


**Fig. 4: TomoTwin locates proteins in a cellular environment.** **a**, Representative slice view of a tomogram containing *Mycoplasma pneumoniae*. **b**, Slice view highlighting positions of picked 70S ribosomes localized in 3D with TomoTwin. Scale bar 100 nm **c**, 3D representation of ribosome positioning within the tomogram, a represented slice is superimposed with 3D classes of ribosomes arranged according to their corresponding coordinates and orientation. **d**, 3D classes from 18,246 particles. Scale bar 10 nm.



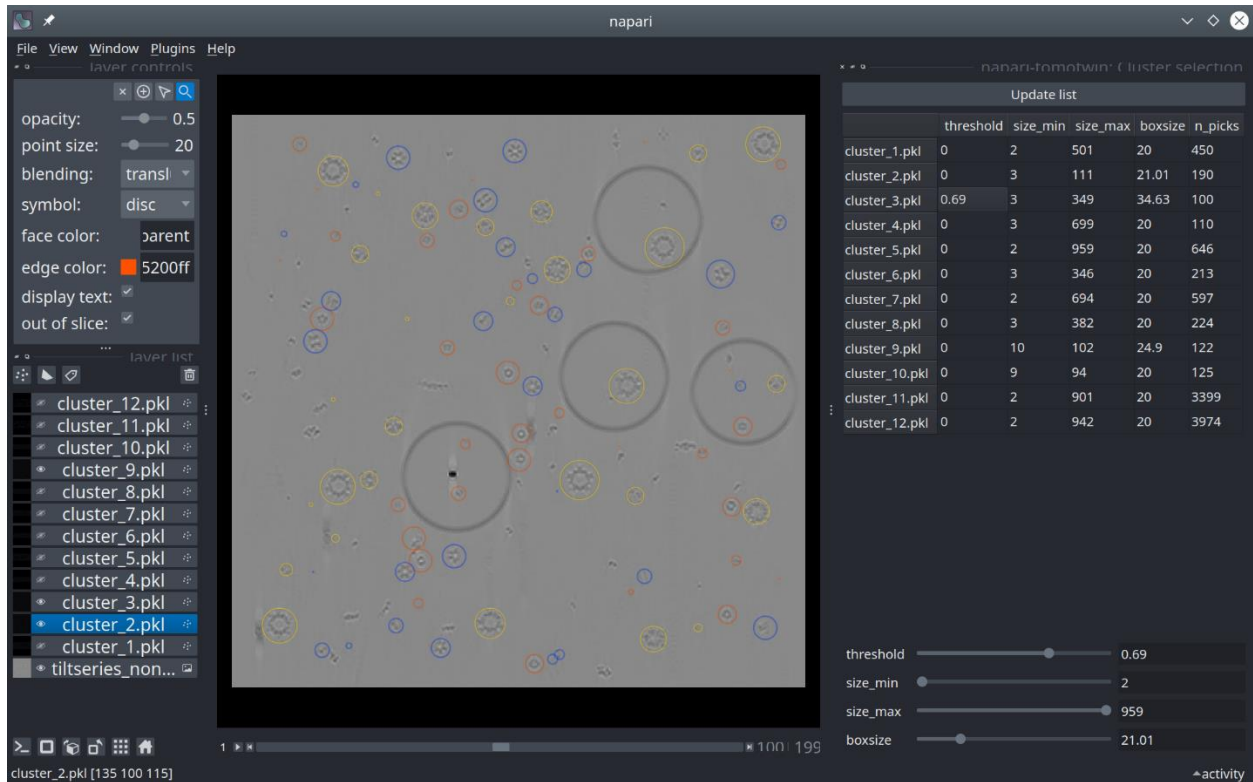
**Fig. 5: TomoTwin enables structural data mining on the embedding manifold.** **a**, Highlighted clusters of all 7 proteins on the generalization tomogram 2D manifold approximation. **b**, Respective particle locations from cluster 3 and 5 which corresponds to the proteins with PDB ID 2DF7 (left) and 1FZG (right). White are true-positive picks and black false-negative. In both cases there were no false-positive selections. **c**, 2D manifold approximation of the embedding space of a tomogram containing *Mycoplasma pneumoniae* (EMPIAR 10499). Highlighted is the manual selected cluster which corresponds to the 70S ribosome. **d**, Using the cluster center for picking identified all ribosomes previously selected by the reference-based picking (white) with a few reference-only selections (blue). **e**, F1 scores for the individual clusters in comparison with the F1 scores for reference-based picking. On average the clustering performed slightly better (0.84 vs 0.82 mean F1 score). However, for some individual proteins the difference was larger (e.g. cluster 7). Scale bar 100 nm



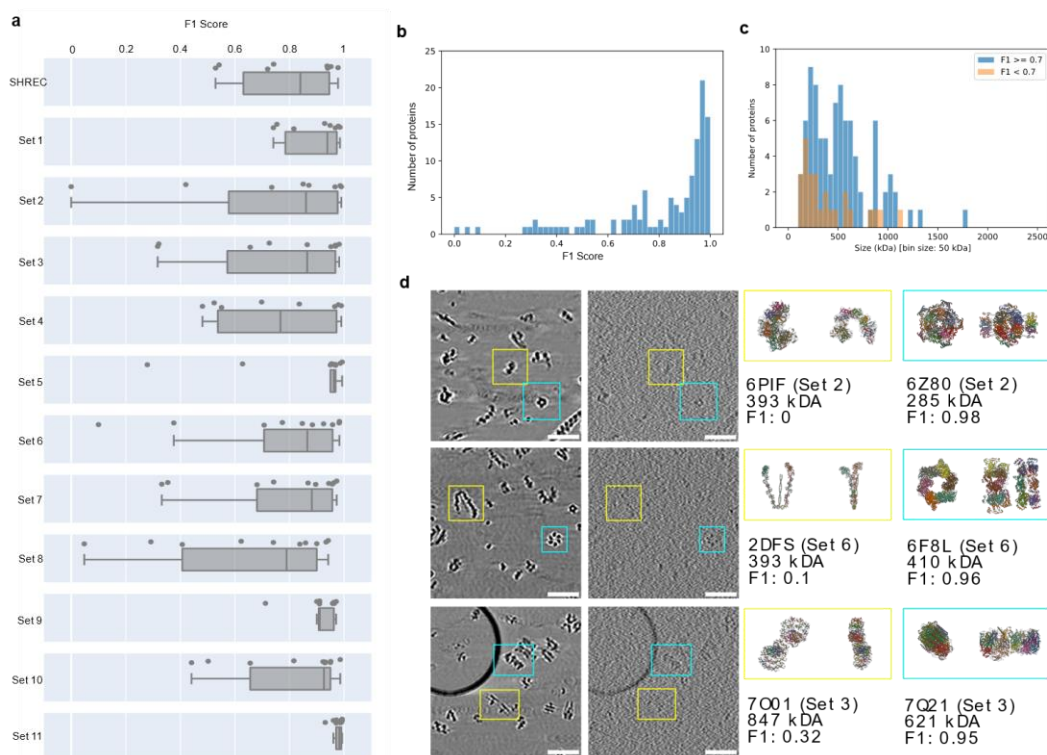


**Supplementary Fig. 1: TomoTwin convolutional architecture and metric learning strategy.** **a**, Architecture of 3D convolutional network utilized by TomoTwin to translate 3D real space tomogram subvolumes into embedding vectors for deep metric learning. **b**, Overview of the deep metric learning training scheme employed by TomoTwin wherein data triplets are constructed of anchor, positive, and negative subvolumes. The triplets of subvolumes are each convolved by the 3D CNN of TomoTwin and the resulting embedding vectors are used to calculate the distance metrics implicit in the triplet loss function. **c**, Uniform manifold approximation of protein subvolume embeddings colored according to protein PDB code from TomoTwin 3D CNN in first training epoch and best model after 600 training epochs.



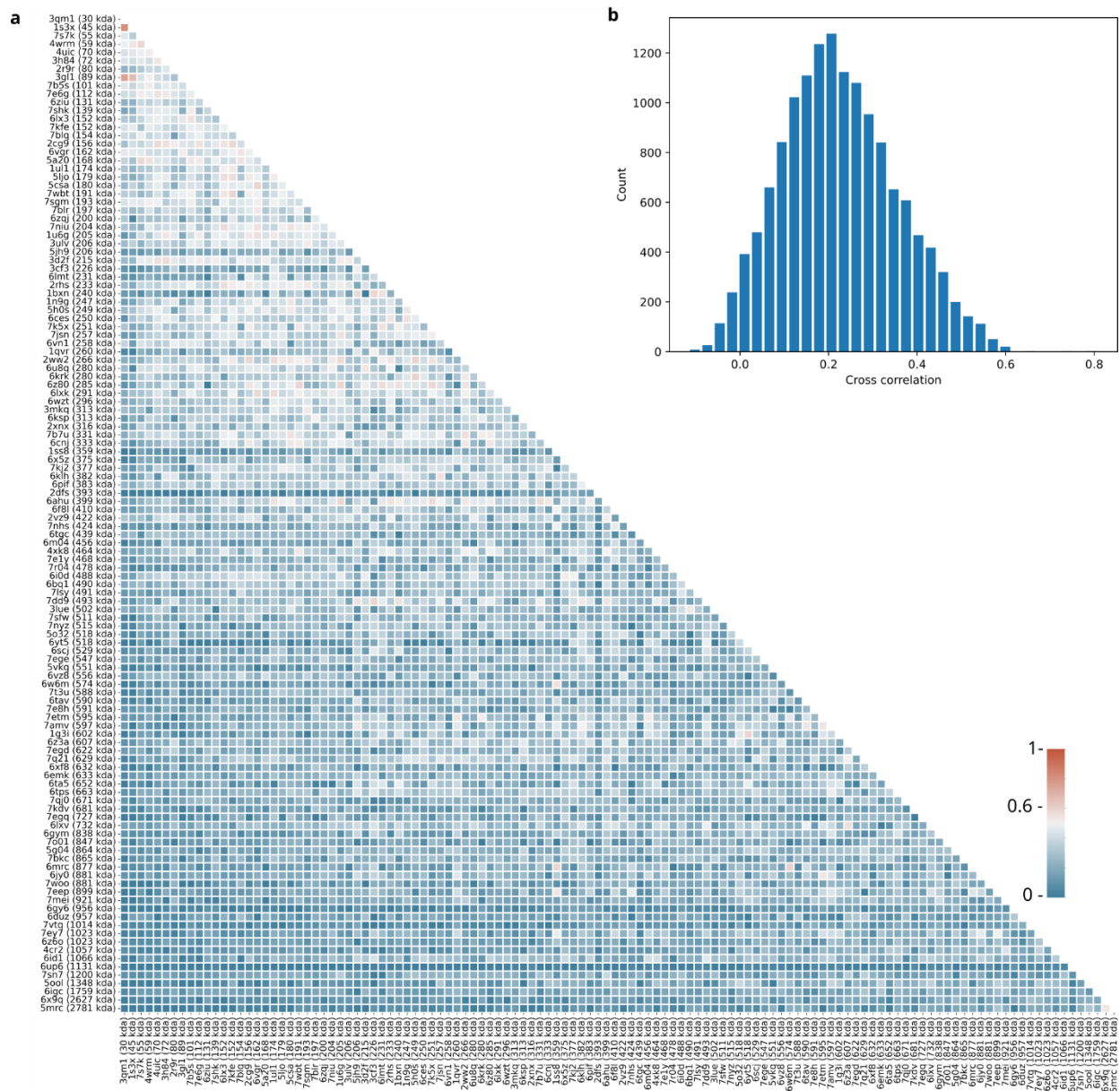


**Supplementary Fig. 2: Graphical user interface of TomoTwin implemented as a Napari plugin.** Visualization of protein picks in simulated generalization tomogram identified by the clustering workflow. Picks for 3 out of 12 clusters are shown as spheres. The lefthand panel allows users to adjust various visualization settings for the tomogram including 3D viewing as an isosurface. The righthand panel allows users to filter picks for each cluster according to similarity threshold, minimum and maximum size, and adjust the box size for viewing.

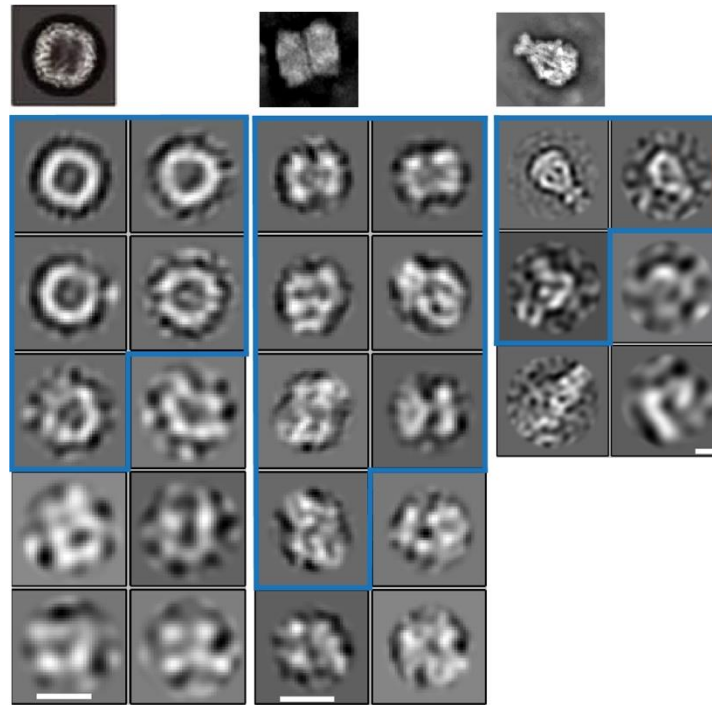


**Supplementary Fig. 3: TomoTwin identifies proteins with high accuracy by using single particle subvolumes as reference.**

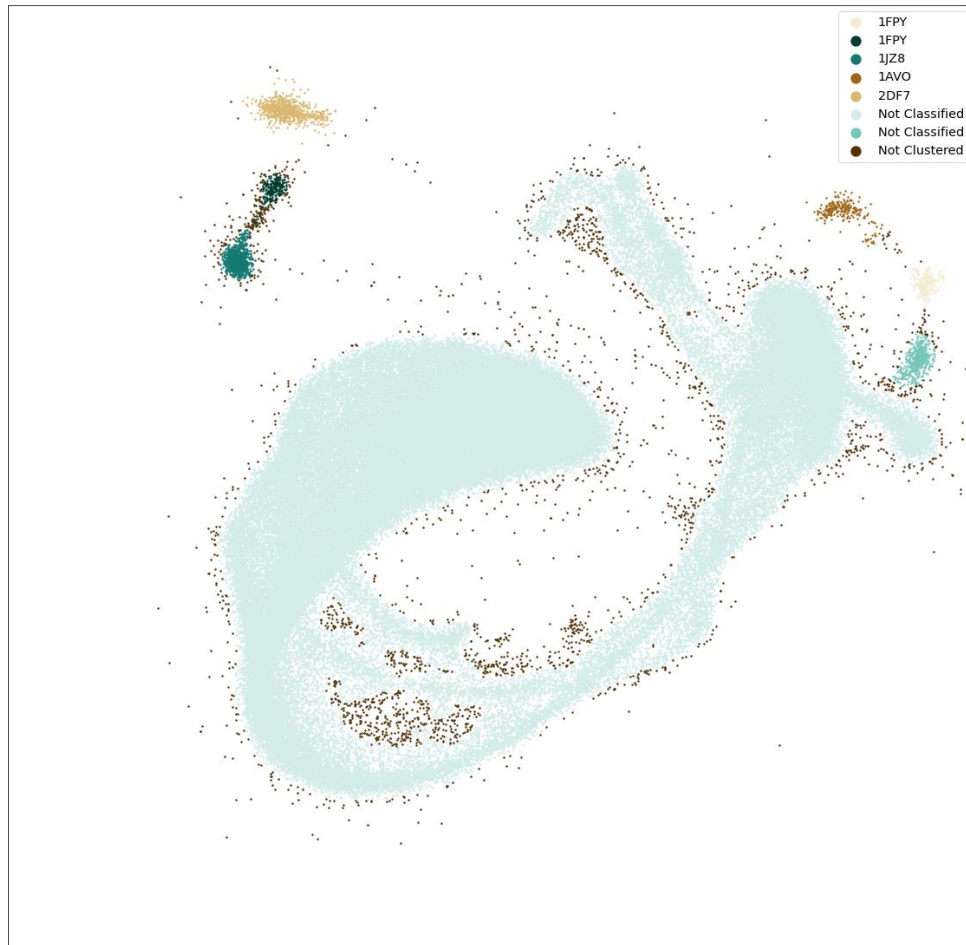
**a**, F1 scores of TomoTwin on the validation tomograms. The median F1 score of the individual sets is most often above 0.8 and not lower than 0.76. **b**, The overall distribution of F1 scores with a median of 0.92. However, a tail of proteins with low F1 scores can be seen. **c**, Size distribution of particles that show good F1 scores ( $F1 \geq 0.7$ ) and those with rather low F1 scores ( $F1 < 0.7$ ). **d**, Examples of proteins of similar size with low (yellow) and high (cyan) F1 score. On the left side the individual particles are depicted in a noisy and noise free reconstruction, respectively. On the right side, the respective structures, F1 scores and sizes are shown. It can be seen that the proteins which were not identified properly by TomoTwin have a lower contrast than the other proteins. Scale bars 100 nm.



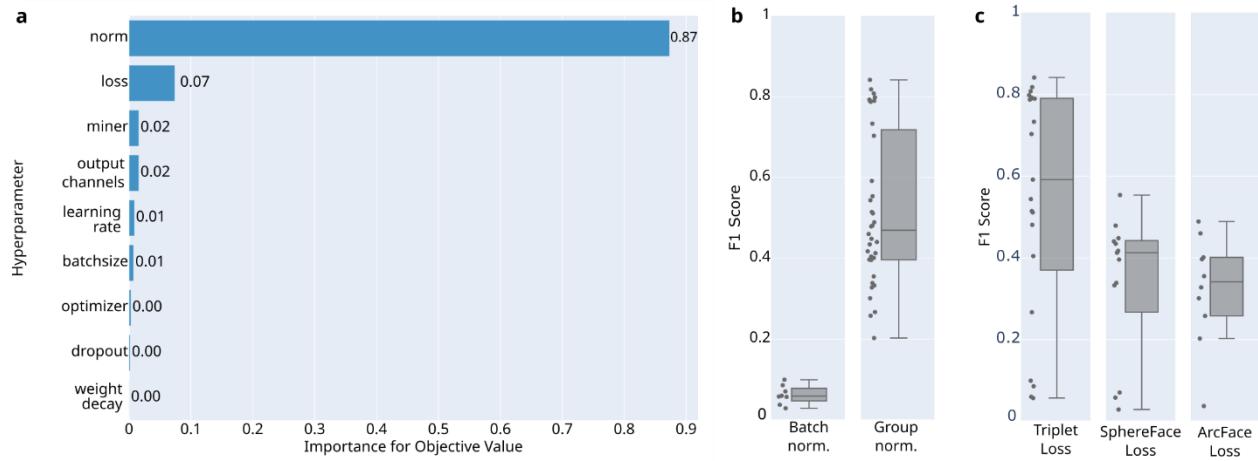
**Supplementary Fig. 4: Characterization of the training data set.** **a**, Pairwise cross-correlation matrix for all 120 proteins sorted by size. Cross-correlations were calculated by converting the individual PDBs to density maps with a pixel size of 1 nm, aligning them pairwise with EMAN2 and calculating the cross-correlation of the aligned pairs. To maximize the value for training, we selected proteins so that all pairs except 3 have a cross correlation value below 0.6. The three pairs with higher correlation are from the SHREC dataset and were not simulated by us. Higher correlation values are more likely for smaller proteins. **b**, Histogram of the pairwise cross correlation values. The mean cross correlation value is 0.22 with a standard deviation of 0.13.



**Supplementary Fig. 5: 2D classes of proteins identified in a mixed tomogram.** Example 2D classes from previous studies by single particle analysis of apoferritin<sup>59</sup>, RhsA<sup>34</sup>, and TcdA1<sup>58</sup> respectively; 2D class averages of TomoTwin picked subvolumes after projection to 2D. Classes outlined in blue were judged to be positive classes by expert inspection, indicating that they contain particles of the appropriate protein. Scale bar: 5 nm



**Supplementary Fig. 6: Automated identification of clusters of interest using HDBSCAN.** A subset of the approximated manifold of Figure 5a was used to run density-based clustering which located 5 out of 7 clusters of interest in an unsupervised. R implementation of HDBSCAN was run with a `min_samples` of 50 and a minimum cluster size of 50.



**Supplementary Fig. 7: Hyperparameter optimization of TomoTwin. a**, Hyperparameter importance estimated by Optuna<sup>50</sup> after 180 trials with different configurations. **b**, F1 scores for trials using either the batch normalization or group normalization layers in convolutional neural network. Points represent the individual trials. Group normalization performed in general better than batch normalization in all cases. **c**, F1 score for trials using either Triplet-, SphereFace-, or ArcFace-Loss.