1  **Instance segmentation of mitochondria in electron microscopy images with a generalist**

2  **deep learning model.**

3

4  Ryan Conrad [1,2] and Kedar Narayan [1,2]

5  [1] Center for Molecular Microscopy, Center for Cancer Research, National Cancer Institute,

6  National Institutes of Health, Bethesda 20892, Maryland, USA.

7  [2] Cancer Research Technology Program, Frederick National Laboratory for Cancer Research,

8  Frederick 21702, Maryland, USA.

9

10  **Abstract**

11  Mitochondria are extremely pleomorphic in biology. Automatically identifying each one

12  precisely and accurately from any 2D or volume electron microscopy (volume EM) dataset is an

13  unsolved computational challenge. Current deep learning (DL) models are trained within limited

14  contexts, restricting their widescale utility and potential as a universal or generalist solution for

15  mitochondrial segmentation. To address this, we amass a highly heterogeneous ~1.5 x $10^6$

16  unlabeled cellular EM image dataset and a ~22,000, partially crowdsource-labeled,

17  mitochondrial instance segmentation dataset. We release MitoNet, a DL model trained on these

18  data, which performs well on new and challenging volume EM benchmarks. An accompanying

19  Python package and napari plugin, called empanada, can be used for efficient training, inference,

20  and clean-up of instance segmentations on EM images.

21

22  Keywords

23  Volume EM, segmentation, deep learning, electron microscopy, Panoptic, benchmark, image

24  dataset, crowdsourcing

25

26    **Introduction**

27

28    Electron microscopy (EM) reveals 2D snapshots of cellular and subcellular ultrastructure at

29    unrivaled resolutions, and more recent volume EM approaches extend these into the third

30    dimension[1,2]. Technological advances in volume EM have dramatically increased the size of

31    volumes that can be interrogated, most notably in connectomics[3–5]. Developing meaningful

32    quantitative insights from EM data often requires the segmentation of features of interest. For

33    these large neuronal datasets, machine learning (ML) or deep learning (DL) based tracing

34    algorithms have helped generate detailed wiring diagrams[6–10]. Similar DL approaches have also

35    provided insights into a variety of sub-cellular structures in other systems[11–15].

36

37    As one of the most ubiquitous and morphologically complex organelles in biology[16–18] and a

38    critical player in cellular physiology[19–21] and pathological processes[22–24], mitochondria are

39    frequent targets of analyses. Mitochondria have extraordinarily variable sizes, shapes,

40    distributions and intra-organellar structures, yet they are instantly recognizable by their

41    ultrastructure. This suggests that a universal DL model that accurately and precisely recognizes

42    individual mitochondria in any given EM image should be possible.

43

44    Current approaches largely take the following strategy: a region of interest (ROI) from a dataset

45    is chosen for dense manual segmentation, a model is trained on this ROI and then inference is

46    run on the larger dataset, often followed by some manual "polishing"[6,7,11,13,25,26]. This strategy

47    can lead to visually impressive results[11,27]; however, the constrained contexts presented during

48    model training are one reason why such models have failed to generalize. As a result,

49    segmentation quality drops when the models are presented with unseen cell and tissue types,

50    sample preparation and imaging protocols, or even in some cases, regions of a dataset spatially

51    distant from the training ROI[11,13,28]. Poor generalization limits the usefulness of models, as

52    unfamiliar datasets require revisiting the cycle of manual annotation, model training,

53    hyperparameter tuning, and prediction updating. Achieving good generalization is a challenging

54    task. Mitochondria are present in nearly all eukaryotic cells such that the landscape of cellular

55    contexts that a generalist segmentation model must account for is immense. This challenge

56     fundamentally distinguishes mitochondrial segmentation from other tasks such as tracing

57     neurites or synapses that exist in a single tissue context (e.g., brain).

58

59     There are several publicly available datasets with segmented mitochondria from a narrow range

60     of tissues: MitoEM,[26] Perez et al.,[29] Lucchi++ and Kasthuri++[25] are from brain tissue, Guay et

61     al.[30] from human platelets, UroCell[13] from urinary bladder and Heinrich et al.[11] from three kinds

62     of *in vitro* cells (HeLa, Jurkat and macrophage). Only one of these datasets (MitoEM) includes

63     instance segmentations (i.e., where each mitochondrion is assigned a different label). We

64     hypothesized that sparse 2D instance segmentations from an eclectic set of EM images would

65     more effectively expand the range of contexts as compared to dense volumetric segmentation.

66     Moreover, instance segmentation can unlock methods for the quantification of mitochondrial

67     morphologies[31–34], networks[18,35], and fusion and fission dynamics[36,37], for which semantic

68     segmentation is insufficient.

69

70     Here, we have curated a heterogeneous, non-redundant, information-rich and relevant unlabeled

71     EM image dataset (at ~1.5 x $10^6$ images, the largest of its kind to our knowledge), called

72     *CEM1.5M,* for use as a database from which to sample images for mitochondrial segmentation.

73     Through the unification of existing labeled datasets and crowdsourced annotations of images

74     from CEM1.5M, we have created a similarly diverse labeled image dataset called *CEM-MitoLab*

75     for training ML/DL models. *MitoNet* is the model resulting from self-supervised pre-training on

76     CEM1.5M followed by supervised training on CEM-MitoLab. We show that MitoNet

77     outperforms other candidate datasets when tested on new and challenging benchmarks (we also

78     share these benchmarks). Finally, in line with our overall mission is to democratize these efforts,

79     we have created a Python package and napari plugin called *empanada*, which allows easy model

80     training and inference as well as prediction clean-ups for MitoNet and any other ML/DL models.

81

82    **Results**

83

84    *Dataset Collation*

85

86    The first step in creating this resource was the generation of a database of unlabeled EM images

87    of cells and tissues both for self-supervised model pre-training and as a source for sampling

88    images for mitochondrial annotation. We collected a total of 463 volume EM datasets, of which

89    352 were produced in-house at the Center for Molecular Microscopy (CMM) and 111 were

90    generated and publicly deposited by external sources. Datasets larger than 5 GB were randomly

91    cropped into 320 cubes of 256x256x256 (equaling 5 GB of unsigned 8-bit voxels) to limit their

92    overrepresentation. Smaller datasets were kept at their original dimensions. In total, this yielded

93    15,152 ROIs, corresponding to 338 GB. To these we added 5,390 2D EM images acquired at

94    comparable pixel samplings (traditional TEM and STEM images of stained and resin-embedded

95    biological samples), of which 1,738 images were from CMM, 2,261 were from the NCI electron

96    microscopy lab and 1,391 were from external sources. Lastly, 27 video files of volume EM data

97    from online publications were incorporated. Metadata and appropriate attribution were recorded,

98    where possible, for all datasets (**Supplementary File 1**). We then applied the data preparation

99    and curation pipeline previously developed in our work on the CEM500K dataset[38] to generate

100   approximately $1.5 \times 10^6$ 2D image patches of 224x224 pixels. Patches from isotropic voxel

101   datasets were derived from xy, xz and yz planes. The curation pipeline removed near-duplicate

102   patches within each dataset before a neural network classified patches as "informative" or

103   "uninformative", the latter of which were discarded.

104

105   Of the combined 490 volume and video datasets, we selected 478 for annotation. The excluded

106   volumes were those with pixel sizes greater than 40 nm or severe artifacts that likely would have

107   made accurate annotation impossible. 20 of the volumes had 2D or 3D ROIs with mitochondrial

108   instance segmentations (13 externally sourced; seven from CMM). The 3D segmentations were

109   cropped into 2D patches and combined with 3,289 previously generated in-house segmentations

110   of 2D data to form the "legacy" portion of the supervised dataset (see **Materials and Methods**).

111   For all the remaining volumes, which had no corresponding mitochondrial annotations, we

112   sampled patches for later manual segmentation. The overall curation process is outlined in

113 **Figure 1a**. Even with the 5 GB cap, the number of patches derived from each dataset within

114 CEM1.5M varied significantly (**Figure 1b**) such that a naïve random sampling would have

115 resulted in drastically oversampling the largest ones. Therefore, when picking annotation targets,

116 we chose a maximum of 15 random patches from every dataset. For large and heterogeneous

117 datasets this very sparse sampling likely misses some unique cellular contexts. However, we

118 decided that for testing our hypothesis and the release of these resources, this level of sampling

119 would adequately represent the data landscape while avoiding the time and labor cost of denser

120 annotation. For the 2D images in CEM1.5M, we circumvented unequal sampling by coarsely

121 grouping sets of 2D images together by known publication or imaging project (see **Materials**

122 **and Methods**). This resulted in the creation of 83 groups of 2D EM images. To compensate for

123 the possibility that groups may each represent multiple biological contexts, we selected a

124 maximum of up to 25 instead of 15 random patches.

125

126 At the end of data collation, we had 1,494,389 image patches for self-supervised pre-training,

127 and for supervised training, nearly 15,897 previously annotated images from legacy data and

128 5,841 images set aside for crowdsourced annotation. Unlike CEM1.5M and legacy data, which

129 were significantly imbalanced toward a few large datasets, the images set aside for annotation

130 almost evenly represented every volume EM dataset. During model training, weighted sampling

131 from datasets was used to correct for imbalance in the combined legacy and crowdsourced data,

132 see **Figure 1b, Materials and Methods**). The total 21,738 annotated images (of various sizes,

133 **Supplementary Figure 1a**) that constituted CEM-MitoLab contained 134,812 mitochondrial

134 instances and represented myriad pixel resolutions, imaging techniques, sample preparation

135 protocols, cells, tissues, and organisms (**Figure 1c-f, Supplementary Figure 1b**).

136

137 Our dataset consists exclusively of 2D images and favors a broad but superficial sampling of

138 very many cellular contexts. **Supplementary Figure 1c** shows a random sample of patches from

139 the Heinrich et al.,[11] and MitoEM[26] datasets, revealing a similar appearance across images. The

140 mitochondria also showed little variation in contrast (**Supplementary Figure 1d**), as may be

141 expected with the standardized sample preparation and imaging techniques used by these groups.

142 By eschewing 3D segmentation strategies, we were able to quickly annotate widely, building a

143 significantly larger and more heterogeneous dataset than the current alternatives. As a resource,

144   the CEM1.5M dataset can be used to build similarly diverse labeled datasets for other

145   downstream segmentation tasks.

146

147   *Crowdsourced Annotation*

148

149   We were inspired by the citizen science project Etch-a-cell[15] on the Zooniverse platform, in

150   which thousands of volunteers annotated organelles like the nuclear envelope, mitochondria and

151   endoplasmic reticulum (ER) in a single or small number of volume EM reconstructions. To

152   create crowdsourced segmentations for mitochondria in CEM-MitoLab, we created a "private"

153   project on Zooniverse, where 34 students recruited from local high schools were briefly trained

154   to identify and annotate mitochondria in diverse and challenging EM images. **Supplementary**

155   **Figure 1c** demonstrates the difficulty of the task: The appearance of a mitochondrion depends

156   strongly on cell and tissue type, sample preparation and imaging approach. To give some 3D

157   context and help resolve ambiguous images and edge cases, we presented images extracted from

158   volumetric datasets as short "flipbooks" of five consecutive slices. For this task students were

159   instructed to only annotate the middle slice of the flipbook (screenshots of the user interface

160   shown in **Supplementary Figure 2a**). We also implemented three controls to enhance the

161   quality of annotations. First, a retirement limit of ten was set such that each image was annotated

162   independently by ten students, and not shown again. These ten annotations were combined to

163   form a consensus instance segmentation (see **Materials and Methods, Supplementary Figure**

164   **3**). Second, we asked annotators to rate their confidence in their annotation on a scale from 1 to 5

165   (where 1 was not at all confident and 5 was very confident). Third, all consensus annotations

166   were reviewed by at least two experienced CMM researchers to create the final ground truth

167   segmentation.

168

169   For this project, nine sets of 500-1000 images were uploaded to the Zooniverse platform over a

170   period of six months. These sets were organized into three groups with each group containing

171   patches from different 2D and 3D datasets. Group 1 contained a mixture of in-house and external

172   volume EM datasets, Group 2 contained in-house and NCI 2D datasets and Group 3 exclusively

173   contained in-house volume EM datasets acquired since the launch of the project. As a baseline,

174   and to test the annotation interface, the first two sets of images (referred to as Group 1a) were

175 annotated by seven experienced EM researchers, with an image retirement limit of five. The

176 remaining three sets of images in Group 1 (referred to as Group 1b) along with Group 2 and

177 Group 3 were annotated exclusively by students.

178

179 We found that the consensus annotations for mitochondria with clearly defined outer membranes

180 and cristae were often high-quality, especially at the pixel-level **(Figure 2a)**. 39.6% of consensus

181 annotations required no corrections upon expert review, 12.5% required cosmetic pixel-level

182 corrections and the remaining 47.9% required some combination of pixel-level, false positive

183 and false negative corrections. Among this last set, the average correction added 3.7 false

184 negatives and removed 0.4 false positives per image. The total number of corrections for each

185 group of images are tabulated in **Supplementary Figure 2b**. Annotations in need of the most

186 corrections were often the result of systematic gaps in the annotators' knowledge. Such errors

187 were manifested as annotations with strong mutual agreements which nevertheless were far from

188 the expert-reviewed ground truth **(Supplementary Figure 2c)**. Nevertheless, student annotation

189 performance as a whole improved over the course of the project. The students' first annotations

190 had much lower F1 scores measured at an intersection-over-union (IoU) threshold of 0.5

191 (F1@50) than experts **(Figure 2b)**. By Group 3, the median student F1@50 scores had risen

192 from 0.57 to 0.72, and impressively, the top 50% of students achieved F1@50 scores that were

193 within the range of experienced annotator scores for Group 1a. This improvement occurred

194 despite the introduction of previously unseen cellular contexts in each group and may be

195 attributed to the heterogeneous data encountered earlier in the project. The improvements were

196 not uniform for all students and this was reflected by a widening range in annotator performance

197 with each new group of images. Seven of the 29 annotators (five dropped out before the last

198 group of images) actually did worse over time **(Figure 2b).**

199

200 Despite the variability in annotator quality, consensus annotations were always near or above the

201 75$^{th}$ percentile of annotators **(Figure 2b)**. This underscores the effectiveness of acquiring

202 multiple independent annotations and applying a robust consensus algorithm. Although there was

203 a slight upward trend with increasing the limit for students, the gains in annotation quality

204 appeared to start leveling off after five or six regardless of the group of images being annotated

205 **(Figure 2c)**. The lower limit of five for experienced annotators was sufficient to reach nearly

206    perfect annotations, with consensus scores reaching 0.9. With a retirement limit of ten for the

207    students annotators, we observed that the best consensus annotations on a per-image basis were

208    created by setting a vote threshold of five (i.e., at least five annotators must have agreed on the

209    label of a pixel for that label to be used in the consensus) (**Figure 2d**). This vote threshold also

210    achieved a reasonable balance between true positive, false positive and false negative detections

211    (**Supplementary Figure 2d**). While this vote threshold was set after evaluating results of Group

212    1, we subsequently observed that F1@50 over all mitochondrial instances that were annotated by

213    students, instead of over all images, was marginally higher at a threshold of three (**Figure 2d**).

214    However, even with a vote threshold of five for all consensus annotations, our approach created

215    each instance segmentation as accurately as possible, independent of number of mitochondria per

216    image.

217

218    Lastly, we hoped that self-reporting of annotation confidence could be used to filter out low-

219    quality predictions, instead we found this factor to be a poor proxy for annotation quality. While

220    there was an increase in average F1@50 with increasing confidence scores, there was also a wide

221    variance. Thus discarding "low confidence" annotations would have thrown out nearly as many

222    good as bad annotations. (**Figure 2e**). There was only a weak correlation (r=0.322) between an

223    annotator's average reported confidence score and that annotator's average annotation quality

224    (**Figure 2f**), suggesting different levels of baseline confidence within the student cohort. For

225    larger crowdsourced datasets where proofreading every annotation is infeasible, confidence

226    scores may still be useful as a guide for identifying annotations most likely to require correction.

227    Surprisingly, we also observed that there was no correlation (r=-0.096) between the number of

228    annotations a student completed and their average F1@50 score (**Figure 2g**). These findings

229    highlight the difficulty of gathering annotations from the public even with a motivated and

230    trained cohort. Future crowdsource projects may benefit from targeted coaching or

231    exclusion/weighting of lower-performing annotators. Time and resources must be invested to

232    teach the skills necessary for handling such tasks, and still, reliable and high-performing

233    annotators can only be identified after expert review and corrections over several rounds.

234

235    *Benchmark Datasets*

236

237     To assess how well models trained on the collected mitochondrial data might generalize to future

238     unseen images, we withheld six instance segmented volumes from the above pipelines. These

239     volumes were chosen to encapsulate a few different cellular contexts, mitochondrial

240     morphologies, and sample preparation protocols. Crucially, they also represented varying levels

241     of difficulty to help identify deficiencies in the training dataset, model architectures and

242     postprocessing algorithms. The six volumes, five of which were not previously annotated,

243     included conventionally fixed, heavy metal stained and resin-embedded samples (different

244     sample preparation protocols, and imaged independently) of fly[39] and mouse brain[25,40] tissue,

245     glycolytic muscle tissue[35], mouse salivary gland, and a HeLa cell. The sixth volume was a high-

246     pressure frozen (HPF), freeze substituted C. elegans embryo specimen. The first three of these

247     were generated externally and the last three were generated in-house. Representative 2D images

248     and 3D renderings of mitochondria from each dataset are shown in **Figure 3a.** The mouse brain

249     volume was the test set from the commonly used Lucchi++ benchmark (we converted the

250     semantically segmented mitochondria to instance segmentation, see **Materials and Methods**).

251     This benchmark was useful to calibrate performance because there were established results for

252     comparison. Importantly, unlike prior research on this dataset, the entire Lucchi++ dataset was

253     excluded from our training data. Aside from the Lucchi++ benchmark, which was imaged at

254     isotropic 5 nm voxel resolutions, we intentionally chose or resampled the other volumes to have

255     intermediate resolutions (10-25 nm) in line with the most common resolutions in the training

256     dataset. In this range, mitochondria are generally still easy to identify on the whole but may be

257     small.

258

259     Each of the new benchmark volumes presented unique challenges. The mitochondria in the fly

260     brain volume were morphologically simple but appeared in two distinct variants: lightly stained

261     with poorly defined cristae and darkly stained with well-defined cristae (**orange and blue**

262     **arrows respectively**). The HeLa cell volume, in addition to being cluttered with subcellular

263     features such as Golgi, ER and a variety of vesicles, had areas of localized heavy metal

264     precipitation from sample preparation (**green arrow**). The C. elegans embryo volume has a

265     mixture of two mitochondrial morphologies that are difficult to segment: small puncta and

266     skinny tubules, with a small median cross-sectional radius of six voxels (**Figure 3b**). Moreover,

267     this volume also contained membranous organelles[41] that might induce false positive detections

268     as they appear similar to mitochondria with swollen cristae at these resolutions (**yellow arrow**).

269     Finally, overall contrast was relatively lower, as expected for high-pressure freezing (HPF) and

270     freeze-substitution used in the sample preparation protocol. The glycolytic muscle volume had a

271     relatively uncluttered background but more complex and elaborate branched morphologies.

272     About a third of the mitochondria had more than one branch and ten had four or more long

273     branches, with the most branched mitochondrion having 11 (**Figure 3b**). The salivary gland

274     volume was by far the most challenging. Besides the unusual tissue context, the mitochondria

275     were weakly stained relative to the cytoplasm, which was packed with ER. Many mitochondria

276     were also flat or bowl-shaped and tightly pressed against and around acini (**red arrow**).

277     Especially difficult for accurate instance segmentation, mitochondria were closely packed

278     together with the boundaries of over 75% being within just five voxels of the next nearest

279     mitochondrion's boundary (**Figure 3b**).

280

281     *Instance Segmentation Algorithm*

282

283     For 2D instance segmentation of mitochondria, we adopted the Panoptic-DeepLab[42] (PDL)

284     architecture. PDL is an encoder-decoder architecture, similar to the standard U-Net[43]. As is

285     common for biological images, PDL employs a bottom-up approach to instance segmentation. A

286     key advantage of bottom-up algorithms is that they can segment an arbitrarily large number of

287     objects within a field of view, in contrast to top-down algorithms like Mask R-CNN[44] or

288     Mask2Former[45] that set limits on the number of objects. We trained PDL models to infer

289     mitochondrial semantic segmentations, centers and per-pixel x and y offsets from each center

290     (**Figure 4a**). To create semantic segmentations at resolutions higher than the input image

291     resolution, we also employed the PointRend[46] module. PointRend iteratively interpolates the

292     segmentation, identifies the most uncertain pixels, and then reevaluates the label for those pixels

293     to refine the segmentation (**Supplementary Figure 4**). Thus PointRend allows the model to

294     accept images that have been downsampled to lower resolutions – preferably in the 10-20 nm

295     range to match our labeled dataset – and output segmentations with crisp boundaries at the

296     original resolution without the need for a resolution-specific model architecture. Finally, the

297     semantic segmentation, offsets and object centers were merged into an instance segmentation

298     (**Figure 4a**) via a simple postprocessing algorithm (see **Materials and Methods).**

299

300    Next, we generalized the 2D instance segmentation algorithm to handle volume EM data. As

301    implemented, PDL did not track objects through an image stack. Instead, we matched objects

302    across consecutive slices by computing IoU scores between all pairs of instances and applying

303    the Hungarian algorithm[47]. We observed that long mitochondrial instances in 2D were

304    sometimes erroneously split into multiple fragments. Since the Hungarian algorithm assumes a

305    1-to-1 matching, these fragments were left unmatched with an object in the preceding slice. By

306    computing intersection-over-area (IoA) scores we were able to identify fragments that were

307    enclosed by a larger object on the preceding slice. These were then merged to the label of that

308    larger object (**Figure 4b**). Since the matching algorithm relied only on consecutive slices, a false

309    negative detection on even a single slice could break the connectivity of that object in 3D. To

310    correct this, we used median within a short stack of the last 3, 5 or 7 recorded semantic

311    segmentation probabilities (**Figure 4c**) to infer the correct segmentation for the current slice. As

312    a final step, matching was performed in the reverse direction through the stack. This was

313    essential to correctly merge branched mitochondria **(Figure 4d)**.

314

315    Optionally, for isotropic-voxel volumes, inference may be performed equivalently on images

316    from the xy, xz and yz planes. Following our previous work[48], we refer to this as ortho-plane

317    inference, although the same idea has been applied elsewhere[15,49]. In this case, the outputs of

318    ortho-plane inference were three independent instance segmentations. We merged these into a

319    single consensus instance segmentation by applying the same algorithm used during

320    crowdsourcing (**Figure 4e**). In all experiments, unless otherwise noted, we applied ortho-plane

321    inference to isotropic-voxel volumes. Other model architectures and postprocessing algorithms

322    like watershed[50] could likely work comparably with the architecture presented here. Our design

323    decisions were heavily influenced by computational efficiency. We used simple run-length

324    encoding compression to represent all 2D and 3D instance segmentations. As a result, forward

325    and backward matching of instances and the consensus algorithm, which must process three

326    complete volumetric instance segmentations, was run on large datasets using minimal compute

327    resources. Such considerations are critical for ensuring that this tool can be adopted widely

328    within the EM community.

329

330    *Model training and evaluation*

331

332    Before training any models on the mitochondrial segmentations, we pretrained a ResNet50[51]

333    model on CEM1.5M using the SwAV[52] self-supervised learning algorithm. This pretrained

334    model was employed as the encoder network in our experiments. Our best performing PDL

335    model was trained for 120 epochs on CEM-MitoLab (see **Materials and Methods**). Our results

336    showed that our model, MitoNet, was able to segment a variety of mitochondrial instances

337    accurately in both 2D and 3D (**Figure 5a, b**). Model performance was best on the brain tissue

338    datasets (fly brain and Lucchi++) with both semantic IoU and F1@75 scores of about 0.9 and

339    above (**Figure 5c**). Detailed performance metrics of MitoNet over all benchmarks are included in

340    Table 1. The high contrast, simplicity of mitochondrial morphologies, and relative homogeneity

341    of tissue ultrastructure in these datasets presented an easier challenge. The model also correctly

342    detected both the lightly and darkly stained variants of mitochondria in the fly brain volume as

343    described above. Strikingly, the F1@75 score of 0.88 achieved on the Lucchi++ test set by our

344    generalist model matched the result from models that were trained exclusively on the Lucchi++

345    training set[53]. Moreover, we observed that all false negative detections at an IoU threshold of 0.5

346    were from small instances that were truncated at the edge of the volume (data not shown).

347

348    MitoNet performance on the HeLa and glycolytic muscle benchmarks was also strong even

349    though these datasets were more difficult. While semantic IoU scores still held within the 0.8-0.9

350    range, F1@75 scores were 0.50 and 0.60 respectively. For the HeLa benchmark, the model

351    successfully ignored most of the intracellular clutter that was expected to be challenging (**Figure

352    5a**). Mitochondria in the glycolytic muscle dataset had clearly defined membranes and cristae

353    that should make them easy to distinguish. That said, model performance was surprisingly good

354    given that no images of glycolytic muscle were present in the training dataset at all and only a

355    handful of images were from muscle tissue (mitochondrial ultrastructure is significantly different

356    depending on the oxidative state of muscle)[54]. The lower F1 scores observed for instance

357    segmentation, despite good semantic segmentation, were evidence that individual mitochondria

358    were harder for our model and postprocessing to distinguish in these datasets. At an IoU

359    threshold of 0.5, we observed that over 75% of the false negatives in the HeLa volume and about

360    50% in the glycolytic muscle volume occurred for mitochondria within five or fewer voxels of

361    their nearest neighbor **(Supplementary Figure 5)**. **Figure 5b** shows examples of the over

362    merging errors that our method was prone to make. It should be noted, however, that the model

363    still correctly detected 78% of the closely apposed mitochondria in the HeLa volume and 52% in

364    the glycolytic muscle, and performed well with branched mitochondria. As an example, aside

365    from one minor split, the model and postprocessing correctly handled a heavily branched

366    mitochondrion in the glycolytic muscle volume **(Figure 5b,** blue arrow**).**

367

368    The C. elegans dataset presented a number of challenges for our model. The mitochondria were

369    low contrast, and small and tightly packed together; this was made worse by the low pixel

370    sampling (24 nm). To get the best performance on this dataset we found that it was necessary to

371    use a lower vote threshold of only one plane for ortho-plane inference. With this lower vote

372    threshold, the IoU score increased to 0.60 from 0.42, though the F1@75 score was lower at 0.35

373    (still a substantial increase from 0.19)**.** An IoU threshold of 0.5, the model failed to accurately

374    detect any mitochondria at that were smaller than 1,000 voxels, or ~ 0.25 $\mu m^3$, in size

375    **(Supplementary Figure 5)**; such small instances accounted for about 25% of all false negatives.

376    The model also performed slightly worse on mitochondria with lower contrast (mean grayscale

377    of correctly identified mitochondria (True Positive), 75.9; wrongly missed mitochondria (False

378    Negative), 70.5, $p=3x10^{-4}$). Supplementary Table 1 lists the mean mitochondrial parameters

379    corresponding to TP and FN detections across all benchmarks. Again, closely apposed instances

380    accounted for over 75% of all false negatives. A cluster of a few tightly packed mitochondria

381    could induce multiple false negatives and false positives – and rapidly decrease F1 scores – even

382    when the voxel-level segmentation appeared quite good (**Figure 5b, black arrow**). The model

383    did manage, however, to mostly avoid erroneously labeling the so-called membranous

384    organelles, but occasionally labeled some vesicles as mitochondria (**Figure 5a**).

385

386    On the most challenging of the benchmarks, the salivary gland volume, MitoNet achieved an IoU

387    score of just 0.1 and consequently negligible F1 scores. The unusual appearance and

388    morphologies of the mitochondria and abundance of ER made this benchmark extremely

389    difficult. Therefore, here we tested the ability of the model after minimal fine-tuning, i.e., after a

390    second round of training on a small subset of the volume's data. We extracted a small fraction of

391    ground truth data, 64 random patches of 224x224 pixels from this salivary gland dataset (~0.2%

392    of the total volume). Even with a very brief finetuning step of 500 iterations (~5 minutes on our

393    system), we achieved dramatically improved semantic segmentation quality with IoU increasing

394    to 0.65 **(Figure 5d).** F1@50 scores also improved, from 0.04 to 0.22. The lower F1 score overall

395    was because, uniquely among these benchmarks, a large fraction of the mitochondria in this

396    dataset were part of a network of closely apposed instances; our model and postprocessing

397    merged nearly all of these mitochondria into a single large instance (**Figure 5a, Supplementary**

398    **Figure 5**). However, we underscore an important point: Even when the generalist model fails on

399    a given dataset, it is still a strong starting point for training a specialized model, and can succeed

400    with a modest number of examples and compute power. The main challenges that remain for

401    these benchmarks appear to be the handling of small objects and tightly packed instances

402    (Supplementary Table 1). We found that instance branch length, measured in voxels, was not a

403    significant factor that led to false negative detection **(Supplementary Figure 5)**, meaning that

404    MitoNet and our 3D postprocessing algorithms were able to track mitochondrial instances over

405    many slices of volume EM datasets without losing connectivity.

406

407    Finally, we directly probed the efficacy of CEM-MitoLab against other training datasets by

408    measuring performance of models across all benchmarks except the salivary gland volume. Since

409    each dataset included a different number of patches, for this experiment we trained all models for

410    approximately 10,000 iterations to control bias toward larger training datasets (the 120 epochs

411    used for MitoNet was the equivalent of about 40,000 iterations). The labeled datasets for training

412    included MitoEM[26] and Heinrich et al.[11], which were the largest and most heterogeneous

413    publicly available datasets, respectively, plus our legacy dataset and the crowdsourced dataset

414    both with and without proofreading. We found that the Heinrich et al. dataset was not useful for

415    our benchmark volumes, likely because of its small size and limited breadth. The model trained

416    on it achieved F1 scores of zero at all IoU thresholds and the highest IoU score was $1 \times 10^{-3}$ on

417    the Lucchi++ benchmark. The model trained on MitoEM achieved respectable – but expected

418    since they were similar neuronal data – F1@75 scores of 0.86 and 0.75 on the fly brain and

419    Lucchi++ volumes, but scores of 0, 0.07, and 0.01 on the C. elegans, HeLa cell and glycolytic

420    muscle volumes. The average performance metrics over all benchmarks for MitoNet versions

421    trained on various annotated datasets is shown in Table 2.

422

423    Our legacy dataset, described above, was 5x larger than the crowdsourced dataset and reasonably
424    heterogeneous with data from 17 volumes, over 500 TEM images and over 2,000 small 2D
425    patches from CEM500K. Still, we observed that the model trained on this dataset had slightly
426    lower F1 and IoU scores than the model trained on the expert proofread crowdsourced data
427    (which excluded images in Group 1a). This surprising finding underscores how critical data
428    diversity is to training generalist models. A dataset of nothing but small 2D patches sparsely
429    sampled from numerous datasets was more effective on our benchmarks than any of the much
430    larger but somewhat less heterogeneous datasets. As a corollary to this finding, we observed that
431    the student consensus annotations, despite over 50% requiring some expert correction, yielded a
432    more general model than training on the MitoEM dataset. This suggests that even noisy and
433    inaccurate, but heterogeneous, labeled data contains enough training signal for a model to learn a
434    better representation of mitochondrial ultrastructure than could be achieved from large, high-
435    quality but homogeneous data.

436 **Discussion**

437

438 In this paper we have reported several resources for the growing volume EM community: a

439 highly heterogeneous ~1.5 x $10^6$ unlabeled cellular EM image dataset (CEM1.5M), a well-

440 described, large and heterogeneous dataset with instance labeled mitochondria (CEM-MitoLab),

441 a model trained on this dataset (MitoNet), a set of diverse benchmarks to test this and future

442 models, and a Python package and napari plugin (empanada) for immediate use of MitoNet.

443 Increases in throughput enabled by technological advances and the application of volume EM

444 techniques to a wider range of biological systems have exacerbated the segmentation and data

445 analysis bottleneck in the field. Recent DL-based results have been spectacular in narrow

446 contexts, but researchers risk getting stuck in the cycle of "manually label and train a model on a

447 sub-volume, run inference and polish on the full volume". The models that result from such

448 workflows generalize poorly to unrelated datasets and are usually obsolete as soon as the

449 relevant segmentation is completed. With the ultimate aim of breaking out of this cycle and

450 creating reusable tools that benefit the community, we show that very broad but shallow

451 sampling of cellular contexts is a robust strategy to create general organelle segmentation

452 models. We present these datasets for use as-is or for future expansion or versioning for

453 specialized tasks. To this end, the organization and appropriate description of datasets is crucial.

454 We have created a simple and practical implementation of REMBI protocols[55] in the form of a

455 spreadsheet (**Supplementary File 1**) and have filled in metadata fields at minimum

456 corresponding to image descriptors and basic biological information, so that users can search this

457 release or future iterations of the dataset using metadata terms.

458

459 We made the important decision to forego 3D approaches, knowingly trading accuracy for

460 generalization and efficiency. Native 3D models are expected to perform better for instance

461 segmentation because of their expanded spatial context. However, it is exceedingly time-

462 consuming to generate sufficient annotation datasets to adequately train a model. For example,

463 our crowdsourced dataset of ~6,000 2D labeled images is equivalent to roughly 30 3D labeled

464 images of the same size (224 pixels to a side). Choosing to create such a 3D dataset would have

465 left a great variety of contexts unsampled and would have excluded thousands of TEM and

466 STEM images. Here we attempted to thread the needle. We collected annotations in 2D but gave

467     annotators the extra context of a flipbook. We ran model inference in 2D but developed instance
468     matching, median filtering, and ortho-plane techniques to propagate and incorporate 3D
469     information. Moving beyond 2D will likely be necessary to achieve human-level annotation
470     quality on some volumetric datasets. Indeed, one use-case of MitoNet is that it can be
471     incorporated into human-in-the-loop workflows to rapidly produce 3D segmentations needed to
472     train this next-generation of models.

473

474     Regardless of dimensionality, our results demonstrate that heterogeneity is a central
475     characteristic for a training dataset. The version of MitoNet trained on noisy and often inaccurate
476     crowdsourced annotations from students was measurably better on our benchmarks than the
477     equivalent model trained on the much larger and expertly labeled MitoEM dataset (**Figure 5d**).
478     Still, proofreading the annotations for accuracy was necessary to achieve our best results. Getting
479     high-quality crowdsourced segmentations without expert intervention is challenging, even with
480     powerful tools like Zooniverse. Accumulating the knowledge necessary to recognize
481     mitochondria across many tissues, organisms and sample preparation protocols takes time and
482     training. Relatively few experienced annotators (3-5) are enough to create strong consensus
483     segmentations, meaning that broader expert participation in shared segmentation efforts would
484     be beneficial to the field.

485

486     The value of training data heterogeneity is only apparent when the test data is also
487     heterogeneous. There are few volume EM benchmarks, and those derived from connectomics
488     data are all quite similar to each other; possibly the best known of these, Lucchi++[25,40] has now
489     been mined to the point that it is difficult to track improvements in DL performance[25,53,56], and
490     until this report, only one benchmark, MitoEM, provided instance segmentations. Our
491     benchmark with six instance segmented volumes aims to correct these deficiencies. It is still just
492     a small sampling of the entire cellular landscape, but appears to be the most stringent test
493     available, with a range of mitochondria, contexts, quality and overall difficulty. We suggest that
494     future benchmarks should comprise relatively small ROIs from many volume EM datasets to
495     best test generalization. Training and test datasets must also evolve to accurately measure
496     progress of automated solutions and prevent misleading reports of success.

497

498 Finally, at least within the volume EM field, we are acutely aware of the gap between research

499 groups that have ample resources and computational expertise, and smaller labs that do not. Most

500 groups deal with small, discrete, and one-off 3D reconstructions, so it is important to have tools

501 that are relatively easy to use and not compute heavy. We designed our Python package,

502 empanada, and the corresponding napari plugin, with this in mind. Empanada and the MitoNet

503 model are optimized for compute and memory efficiency and our plugin is easy to install and

504 use. Additionally, it can run on consumer-grade laptops without GPUs or on HPC clusters with

505 many. MitoNet can be incorporated into other established image analysis platforms and will

506 ideally give users a better model that can be quickly deployed or finetuned to meet their

507 segmentation needs. We did not focus on trying to get the absolute best architecture possible for

508 MitoNet and we hope and expect others will supersede our results. Rather, we used an efficient

509 architecture and postprocessing scheme that scales well to large images and images with very

510 many objects, as would be expected for mitochondrial instance segmentation in EM.

511

512 Volume EM has been suggested to be in the midst of a "quiet revolution" in cell biology[57]. A raft

513 of papers in connectomics and in cell biology have revealed insights into a variety of systems by

514 enabling the high-resolution 3D imaging of cellular and subcellular features. These imaging

515 techniques record ultrastructural information in a largely unselective manner, meaning that the

516 data, once collected, need to be parsed out to extract features of interest for visualization and

517 analysis. This step, segmentation, must be accurate and precise to faithfully represent the

518 underlying biology and it must also be efficient to match the speed of data generation. Recent

519 developments in DL approaches show great promise, yet these tools can be held back by the lack

520 of large-scale and relevant data resources. In this report we directly address this problem as it

521 pertains to universal segmentation of mitochondria in EM images. With an eye towards creating

522 a universal mitochondrial model, we employ a strategy of sparse sampling of widely

523 heterogeneous datasets and show that the resulting model, MitoNet, when trained on these data,

524 yields promising results on challenging tasks. We release the resources to the community with

525 the hope that continued work on this approach will expand the EM segmentation toolkit and

526 further accelerate discoveries in this exciting field.

527

**Materials and Methods**

*Unlabeled dataset creation*

The expansion of CEM500K to create the CEM1.5M dataset followed the data standardization and curation protocols presented in our previous work[58]. External datasets were either downloaded in their entirety or, for large datasets stored online in next generation file formats (n5 or zarr), accessed either with the CloudVolume or fibsem-tools APIs. Any datasets larger than 5 GB were randomly cropped into 320 cubes of 256x256x256 (equivalent to 5 GB). When available, to ensure that the randomly chosen crops included cellular content and not empty image padding or resin, a low-resolution overview image volume of the dataset was downloaded and used to identify the extents of informative ROIs. Lastly, cellular EM images from videos were converted to mrc files after removing all frames that showed non-EM content. Metadata, and proper attribution, for each dataset following REMBI[59] standards is available in **Supplementary file 1.**

In-house 2D EM images were collected from the National Cancer Institute's Electron Microscopy Laboratory database. The database included over $2x10^5$ TEM images. A sample of 5,000 images with magnifications between 1000x (~16 nm) and 3000x (~6 nm) that excluded negative stain and immunogold label images was randomly selected. This sample represented the output of collaborations with 76 NIH investigators from over the last 12 years. Metadata for each image was limited. Therefore, images were grouped by investigator before further sampling. Additionally, seven more groups of 2D STEM images were added from previous in-house imaging projects; because metadata was available for these images they were grouped by biological context and not investigator.

*Zooniverse Workflow*

Our Zooniverse workflow closely approximated the Etch-a-cell project[15].Patches in the CEM1.5M dataset extracted from volume EM images were reconstructed into short "flipbooks" (tif stacks of five consecutive images of 224x224), where each patch from the CEM dataset was

559   the center image (the adjacent patches were typically filtered out by the CEM deduplication
560   pipeline). Before uploading to Zooniverse each flipbook had its 8-bit intensity range rescaled
561   from 25 to 235 and was interpolated to a size of 480x480 for easier viewing. For 2D images,
562   crops of 512x512 were used. Intensity was rescaled to the same range as the flipbooks but no
563   resizing was performed.

564

565   The first two sets of images uploaded to Zooniverse (Group 1a) had an annotation retirement
566   limit of five (i.e., every image was annotated independently by five people) and were annotated
567   by a mixture of lab members familiar with cellular EM and inexperienced student annotators.
568   Consensus annotations from these sets only used annotations from experts; the students'
569   feedback on the project helped refine the Zooniverse tutorials and interface. For all remaining
570   sets, which were annotated entirely by an expanded cohort of students, the retirement limit was
571   set to ten. For the first two of these sets, 5-10% of images were annotated by a team of at least
572   three EM experts. These "gold-standard" annotations were used to measure the performance of
573   the student annotators and to provide them with training examples via a Google CoLab
574   notebook. Later, as the students' annotations became more accurate, the non-proofread
575   consensus annotations were shared. In addition to this, feedback via the Zooniverse message
576   channel and a formal instruction session was used to review challenging examples and explain
577   common mistakes.

578

579   *Consensus annotation algorithm*

580

581   The input to the algorithm was a set of **N** annotations (the retirement limit) containing **K**
582   mitochondrial detections. An undirected graph was initialized where each node corresponded to
583   one of the **K** detections. Pairwise bounding box IoU scores were calculated for all detections to
584   form a **KxK** matrix. For all pairs of detections with non-zero bounding box intersections, pixel-
585   level mask IoU scores were calculated. Edges were added to the graph to connect detection
586   nodes with mask IoU scores that exceeded a small threshold value of 0.1 (the same algorithm,
587   when applied during ortho-plane inference, used a threshold of 0.01). The resultant connected
588   components in the detection graph were processed independently to determine the consensus
589   object instances.

590

591    Nodes in a connected component were organized into cliques where all detections within a clique

592    had IoU scores greater than 0.75. An iterative algorithm was then applied to determine whether

593    cliques should be merged or remain split. First, the clique whose detections shared the most

594    edges with other cliques, i.e. the most connected clique, was selected (node C in Supplementary

595    Figure 3 ii, A in 3 iii). Second, if any of the most connected clique's neighbors contained more

596    detections, then the most connected clique was dissolved and its detections were pushed out to

597    each of its neighbors (Supplementary Figure 3 iv). Otherwise, all the most connected clique's

598    neighbors were dissolved and their detections were pulled in by the most connected clique

599    (between B and D in Supplementary Figure 3 iv, and B and C into A in 3 v). These two steps

600    were repeated until no cliques had outgoing edges. Each clique represented a segmented object

601    instance. The detection masks within a clique were added together to form an image where each

602    pixel had a value from 1 to **N** denoting the number of annotators who labeled it. A vote threshold

603    of 3 when the retirement limit was 5 or 5 when the retirement limit was 10 was applied to create

604    the final binary instance mask. In cases where there was weak instance-level consensus between

605    annotators, the binary instance masks from separate cliques within the same connected

606    component could have non-trivial overlaps with each other. If these overlaps resulted in IoU

607    scores greater than 0.1 between binary instance masks, then those masks were combined into a

608    single mask. In the final consensus instance segmentation, all non-zero binary masks were

609    assigned new labels and merged. A schematic of this algorithm is shown in **Supplementary**

610    **Figure 3.** In summary, this algorithm assumed that the most connected clique at each step

611    roughly corresponded to a maximally merged instance. When one of this clique's neighbors

612    contained more detections, this implied that more votes were in favor of splitting than maximal

613    merging. Pushing the most connected clique to its neighbors ensured that no detections were

614    deleted before being counted towards the final pixel-level majority vote. Conversely, when the

615    most connected clique had the most detections, this implied that more votes were in favor of

616    maximal merging and therefore all detections in neighboring cliques were pulled into the most

617    connected clique.

618

619    *Proofreading*

620

621    Once all images in a set were retired, a consensus strength and instance segmentation were

622    calculated for each. The consensus strength was the average F1@50 score of each individual

623    annotation relative to the consensus instance segmentation. Lower consensus strengths indicated

624    greater disagreement between annotators. Next, all images and instance segmentations were

625    ordered by consensus strength and stacked into a single tif file for later proofreading. If the set

626    used flipbooks, then the images and their consensus instance segmentations were resized from

627    480x480 back to their original sizes. Padding was added as needed to ensure that all images in

628    the stack had the same dimensions. Typically, the single proofreading stack was split into chunks

629    of 40-50 images/flipbooks such that multiple experts could perform proofreading in parallel.

630    Proofreaders were instructed to emphasize the correction of false positive and false negative

631    detections over fine-grained instance boundary corrections. Annotations in each chunk of images

632    were verified by a second proofreader and disagreements were discussed and resolved by at least

633    three proofreaders. 3D Slicer[60] was used to visualize and correct annotations. The chunks of

634    proofread images and annotations were again combined into a single large stack. Every image, or

635    every third image for flipbooks, and corresponding instance segmentation was stored as a tiff

636    image with any padding cropped out. Importantly, all annotations were passed through a

637    connected components filter to guarantee that there were no disconnected instances.

638

639    *Training dataset creation*

640

641    The training dataset included all consensus annotations gathered from Zooniverse as well as the

642    legacy data. The latter included annotations from publicly available mitochondrial segmentation

643    benchmarks (Kasthuri++[25,61], Guay[30], UroCell[62], MitoEM[26], Heinrich et al.[11] and Perez et al.[29])

644    and previous in-house projects. Of the benchmark datasets, only the MitoEM benchmark had

645    individual instances labeled; instance segmentations for all others were generated manually after

646    applying a connected components filter. In-house annotations for various unrelated projects

647    included 11 volume EM image reconstructions (see **Supplementary File 1**) as well as 529 TEM

648    images and 2,231 image patches from the CEM500K dataset. None of these images were derived

649    from any of the above, or from benchmark datasets. The 11 volume EM images had previously

650    been annotated by a different deep learning model and then manually corrected, while all the 2D

651    images were manually annotated from scratch. Legacy volume EM images and instance

652     segmentations were cropped into patches of 512x512 and passed through the deduplication and

653     filtering pipeline from CEM500K[58] to remove redundant and uninformative patches. Patches

654     from isotropic voxel volumes were taken from xy, xz and yz planes while patches from

655     anisotropic volumes were taken from xy only. For external data comprising the legacy datasets,

656     all images were used as-is, except: prior to cropping patches, the MitoEM-H and MitoEM-R

657     volumes were binned by a factor of two in x and y (from 8 nm to 16 nm pixels), and all Heinrich

658     et al. images were downloaded at 8 nm resolution.

659

660     *Benchmark dataset creation*

661

662     Four of the six benchmark volume EM datasets were independently annotated by ariadne.ai

663     (HeLa cell, C. elegans, fly brain, salivary gland). The glycolytic muscle[35] dataset was binned by

664     a factor of 2 in all dimensions (from 9 nm to 18 nm voxels) and then automatically annotated by

665     a deep learning model that was trained solely on patches from the volume, and then manually

666     corrected. The remaining benchmark is derived from the Lucchi++[25,40] benchmark. First, the 2D

667     images and labelmaps were stacked to form volumes. Then, the binary mitochondrial labelmaps

668     were converted to rough instance segmentations by applying a connected components filter and

669     then manually corrected to derive the final instance segmentation.

670

671     Not only were all benchmark volumes excluded from both CEM1.5M and crowdsourced

672     annotation, **we also excluded** any images related to them to rigorously test generalization,

673     **Supplementary Figure 6.** OpenOrganelle[39], which was the source of the fly brain benchmark,

674     included multiple volumes of fly brain tissue from different brain areas. We chose the

675     "Drosophila brain: Fan-shaped body" volume and excluded the "Drosophila brain: Entire alpha

676     and alpha' lobes of a mushroom body" and "Drosophila brain: Accessory calyx" volumes. For

677     the C. elegans volume, a related FIB-SEM and STEM dataset were excluded. And lastly, for the

678     Lucchi++[25,40] benchmark, the Lucchi++ training set was excluded.

679

680     *Benchmark datasets measurements*

681

682     Mitochondrial volumes were calculated by trivially counting the number of voxels per instance.

683     Minimum distances between neighboring mitochondria were calculated by computing the

684     distance transform of each instance's complement (i.e., set of all voxels not inside the instance).

685     These results were stored in a new volume (called "distance volume") where each voxel was

686     updated to the minimum distance calculated so far from all the measured complements at that

687     location. The minimum distance to a mitochondrion's nearest neighbor was then the minimum

688     value enclosed by that instance in the distance volume. Branch lengths and cross-sectional

689     diameters were calculated by first skeletonizing[63] each mitochondrion and computing its distance

690     transform. Branches shorter than 60 nm were pruned. Lengths were simply the number of voxels

691     in each branch, and the mean cross-sectional diameter was the average of all values in the

692     distance transform that overlapped with the skeleton.

693

694     *CEM1.5M Pretraining*

695

696     Unsupervised pretraining followed the SwAV algorithm[64]. A ResNet50[65] model was trained for

697     200 epochs with a batch size of 256. All other hyperparameters used the default values defined in

698     **https://github.com/facebookresearch/swav.** Image augmentations included 360-degree

699     rotations, randomly resized crops (following SwAV defaults), brightness and contrast jitter,

700     random Gaussian blur and noise, and horizontal and vertical flips. To correct for the imbalance

701     in the number of patches per dataset, weighted random sampling was applied. Weights per

702     dataset were calculated by:

$$w_i = \frac{n_i^{-\gamma}}{\sum_{j=1}^{N} w_j}$$

703     $\gamma$ was a float from 0 to 1, $n_i$ was the number of patches from the $i^{th}$ dataset, and N was the total

704     number of datasets. Each patch was assigned a sampling weight based on its source dataset.

705     When $\gamma = 0$, patches from all datasets were sampled with the same probability (true random

706     sampling) and when $\gamma = 1$ the sampling was perfectly balanced between datasets. During

707     pretraining we set $\gamma = 0.5$ (i.e., the square root of the number of patches per dataset).

708

709     *Deep learning model architecture*

710

711   The model was based on Panoptic-DeepLab[42] (PDL). In brief, PDL is an encoder-decoder style
712   network that uses atrous spatial pyramid pooling (ASPP) to integrate features that occur over
713   multiple scales. A ResNet50 encoder, pretrained as detailed above, was used in all experiments
714   unless otherwise noted. To preserve more spatial information, strides in the last downsampling
715   layer were replaced with dilation such that the output from the encoder was 16x smaller than the
716   input (compared to 32x smaller with strides). The baseline model configuration included two
717   ASPP modules and decoders; one for semantic segmentation prediction and the other for
718   instance center and offsets prediction. Since input patches used during training were just
719   256x256, the dilation rates in the ASPP modules were set to 2, 4 and 6 with 256 channels per
720   convolution and dropout with probability 0.5 applied to the output. The semantic and instance
721   decoders each had a single level at which the output from the first layer in the encoder (4x
722   smaller than the input image) was fused via depthwise separable convolution with the
723   interpolated output from the ASPP module. Following the standard for PDL, the convolution
724   channels in the instance decoder were half of those in the semantic decoder (16 compared to 32).
725   The semantic segmentation, instance center and offset heads used a single depthwise separable
726   convolution with kernel size of 5.

727

728   *PointRend:* During training, the semantic segmentation logits predicted by PDL were upsampled
729   by a factor of 4 to the original image resolution and then refined for a single step by the
730   PointRend[66] module. For a batch of images, a set of 3,072 points were randomly sampled from
731   anywhere within the dimensions of the images and the segmentation logits were evaluated at
732   these points. From these sampled logits, 1,024 were selected for each image where 25% of the
733   logits were randomly chosen and the other 75% were the logits with the smallest absolute values.
734   The same 1,024 points were also used to sample features from the semantic decoder output.
735   Sampled features and logits were concatenated and fed through a three-layer fully connected
736   network. When evaluating the loss, the 1,024 points were also used to sample the ground truth
737   semantic segmentation. Loss was calculated as the (binary) cross-entropy between the sampled
738   and refined logits and the sampled ground truth. During evaluation, two refinement steps were
739   applied with the semantic segmentation logits being upsampled by a factor of 2 at each step.
740   Points corresponding to the center of all logit pixels were used to sample the logits. From these,
741   points corresponding to the 8,192 logits with the smallest absolute values were selected. The

742    logits and semantic decoder output were sampled at these points and passed through the fully

743    connected network. The refined logits then replaced the unrefined logits in the upsampled

744    segmentation. Since the number of sampled points was fixed, the cost of running PointRend

745    refinement was approximately constant with respect to the size of the input.

746

747    *Model training parameters*

748

749    The overall loss function during training was:

$$\mathcal{L} = \mathcal{L}_{sem} + \alpha * \mathcal{L}_{center} + \beta * \mathcal{L}_{offset} + \gamma * \mathcal{L}_{pr}$$

750    $\alpha, \beta, \gamma$ were constants used to weight the relative contributions of each loss and were set to 200,

751    0.01 and 1. $\mathcal{L}_{sem}$ was the semantic segmentation loss computed as the bootstrapped (binary)

752    cross-entropy where only the top 20% of largest cross-entropy values were averaged across a

753    batch. $\mathcal{L}_{center}$ was the instance center regression loss and $\mathcal{L}_{offset}$ was the center offset loss

754    calculated as mean squared error and absolute error (L1), respectively. $\mathcal{L}_{pr}$ was the PointRend

755    loss calculated as the (binary) cross-entropy.

756

757    All models were trained using the One Cycle learning rate policy[67] with AdamW[68]. The max

758    learning rate was set of 0.003 with weight decay of 0.1 and momentum was cycled from 0.85 to

759    0.95. Learning rate warmup lasted for the first 30% of training epochs. Weighted sampling of

760    patches from datasets was implemented as above ($\gamma = 0.3$ here) to correct for overrepresented

761    datasets in CEM-MitoLab. Image augmentations included large scale jitter[69], where images were

762    zoomed in by a maximum of 2x and zoomed out by a maximum of 10x, random cropping to a

763    256x256 patch, 360-degree rotations, brightness and contrast adjustments and vertical and

764    horizontal flips. Like CEM1.5M pre-training, images were selected from datasets based on

765    weighted sampling. Unless otherwise noted, $\gamma = 0.3$ was used for all training runs.

766

767    *Panoptic-DeepLab Postprocessing*

768

769    First, values in the instance center heatmap that were less than the center confidence threshold

770    (0.1 in all experiments) were zeroed. Non-maximum suppression with a kernel size of 7 was

771    applied to filter out peaks in the heatmap that occurred within the same 7x7 grid. The non-zero

772   coordinates that remained in the heatmap were extracted. Next, the instance center offsets were

773   added to their absolute (x, y) positions within the image such that all pixels belonging to the

774   same object ought to be equal to the coordinate that defines the center of that object. The

775   Euclidean distance between this map of each pixel's corresponding object center and the centers

776   extracted from the heatmap were used to group pixels into instances. The index of the center

777   with the minimum distance at each pixel was adopted as the instance label for that pixel. After

778   applying the sigmoid activation to get probabilities from the semantic segmentation logits and

779   hardening the result over a threshold, the instance labels and semantic segmentation were merged

780   to form the final panoptic segmentation.

781

782   *Stack instance matching*

783

784   To match object instances between consecutive slices in a 2D stack, intersection-over-union

785   (IoU) scores were first calculated between instance bounding boxes. Second, for each pair of

786   instances with non-zero bounding box IoU, a mask IoU was calculated and stored in a matrix of

787   size KxN, where K is the number of instances in slice j and N is the number of instances in slice

788   j+1. At this stage, mask intersection-over-area (IoA) scores were also calculated and stored in a

789   separate KxN matrix. The Hungarian algorithm[47] was then applied to the mask IoU matrix to

790   determine the assignment of instances in slice j to instances in slice j+1 that maximized the total

791   IoU. Any unassigned instances in slice j+1, or assigned instances with IoU less than a threshold

792   value (IoU threshold), were considered unmatched. Unmatched instances with a IoA score

793   greater than a threshold value (IoA threshold) were matched to the instance in slice j with which

794   they shared the highest IoA. The remaining unmatched instances were assigned new labels in the

795   stack. After forward matching (i.e., from slice j to slice j+1) was completed for all images in the

796   stack, a backward pass of matching was conducted (i.e., from slice j to slice j-1). No new labels

797   were assigned in the stack during the reverse pass.

798

799   *Benchmark Set Inference*

800

801   Ortho-plane inference[48] on xy, xz and yz planes was applied to each volume in the benchmark

802   test set. During inference over each plane, semantic segmentations were generated from the

803    median probabilities within a short queue of previous and subsequent segmentation results. A

804    queue length of 3 was used for the C. elegans and fly brain volumes, 5 for the HeLa and

805    glycolytic muscle volumes and 7 for the Lucchi++ and salivary gland volumes. These queue

806    lengths roughly track with voxel size (in nm) such that smaller voxels corresponded to longer

807    queues. Median probabilities were hardened at a confidence threshold of 0.3. Since mitochondria

808    were never occluded in any of our benchmark volumes, mitochondrial instances on each 2D slice

809    were required to be connected components. Instance matching across slices used IoU and IoA

810    thresholds of 0.25. After a forward and backward matching pass, a simple size and bounding box

811    extent filter were applied to eliminate likely false positives. The minimum object size was set at

812    the $5^{th}$ percentile of mitochondrial volume for each benchmark while the minimum bounding box

813    extent was fixed at eight voxels for all datasets. Finally, the same algorithm that was used to

814    generate consensus segmentations during crowdsourced annotation was applied to ensemble the

815    three segmentation stacks created by ortho-plane inference **(Supplementary Figure 3)**. A clique

816    IoU threshold of 0.75 was used for all volumes and a vote threshold of 2 out of 3 was used for all

817    benchmarks except the C. elegans and salivary gland volumes. For those a vote threshold of 1

818    out of 3 was used. This was in response to the relatively low consensus strength between each

819    stack and the final ortho-plane result. Lowering the threshold gave a significant increase in

820    performance on the C. elegans volume and a modest increase for the salivary gland volume.

821

822    *Salivary gland finetuning*

823

824    The salivary gland volume was first cropped into patches of 224x224 and then deduplicated and

825    filtered (see *Training dataset creation* section above). 64 patches (~0.2% of the volume) were

826    randomly selected from the filtered subset and set aside as the finetuning training set. The best

827    performing segmentation model was retrained on the finetuning training set for 500 iterations

828    using the same parameters outline in the *Model training parameters* section above.

829

830    *Benchmark Evaluation*

831

832    Ground truth instance segmentations were stored in json files where each entry contained an

833    instance id, bounding box and run-length encoded segmentation. Equivalent json files were

834    created for all predictions. Instances between the ground truth and prediction were assigned as

835    matches using the Hungarian algorithm to maximize IoU over all possible pairs of instances. In

836    addition to the F1, IoU, and PQ scores reported in results, average precision (AP) scores for

837    future comparison.

838

839    *Test Set Evaluation by Training Dataset*

840

841    To prepare datasets for fair comparison, the MitoEM and Heinrich et al. ground truth ROIs were

842    sliced to 2D images. Since the MitoEM volumes were large and anisotropic, patches of size

843    512x512 were cropped from xy slices only. The Heinrich et al. volumes being small and

844    isotropic were cropped into patches of 224x224 or smaller from xy, xz and yz planes. As above,

845    mitochondrial instances were relabeled on each 2D slice to guarantee that they were connected

846    components. All models were trained for approximately 10,000 iterations.

847

**Acknowledgements**

868    **Figure Legends**

869

870    **Figure 1: Creation of a diverse and representative dataset for mitochondrial instance**

871    **segmentation. a.** Schematic of the data curation pipeline. Volume EM reconstructions and 2D

872    EM images were curated[58] to create CEM1.5M, an unlabeled dataset of $\sim 1.5 \times 10^6$ image

873    patches used for self-supervised pretraining. Approximately 6,000 randomly sampled patches

874    from CEM1.5M (green) were uploaded to the Zooniverse platform for crowdsourced annotation.

875    These, and 16,000 randomly selected patches from previously labeled data (legacy annotations,

876    red) were combined to form the supervised mitochondrial training dataset CEM-MitoLab. **b.**

877    Lorenz plots for CEM1.5M (blue, Gini coefficient = 0.802), crowdsourced data (green, 0.209),

878    legacy data (red, 0.840), CEM-MitoLab (dashed gray, 0.686). Black line, perfect equality in

879    distribution (Gini = 0). **c.** Distribution of imaging plane pixel sizes for volume EM images in

880    CEM-MitoLab. The dashed lines denote pixel sizes for 2D EM images in the dataset. **d.** Imaging

881    technique, **e.** Source organism, **f.** Source tissue (vertebrates only, *in vitro* cells grouped under

882    Not Defined) of datasets in CEM-MitoLab (n=478).

883

884    **Supplementary Figure 1: a.** Distribution of longest image side in pixels for patches in the

885    supervised mitochondrial dataset, CEM-MitoLab. **b.** Breakdown of isotropic versus anisotropic

886    data (left) and sample fixation method (right) in source datasets that make up CEM-MitoLab **c.**

887    Random selection of image patches from CEM-MitoLab (left), Mito-EM (middle) and Heinrich

888    et al **d.** Comparison between 2D mitochondrial instance label maps in CEM-MitoLab (blue),

889    MitoEM (purple) and Heinrich et al. (pink) labeled datasets, with respect to (left to right) area in

890    pixels, contrast, eccentricity, perimeter to area ratio, and clustering.

891

892    **Figure 2: Crowdsourced annotation of CEM-MitoLab. a.** Ground truth (top left), consensus

893    annotation (bottom left) and ten independent student annotations of an image showing high

894    degree of consensus. **b.** Instance segmentation quality measured by F1@50 over time. Left,

895    Group 1a (expert annotation), 1b, 2, 3 (students). Blue dots, individual scores, pink, consensus

896    score, box denotes median, $25^{th}$ and $75^{th}$ percentile. Right, change in individual annotator

897    performance over time. **c-f** Instance segmentation quality measured by F1@50 plotted against **c.**

898    Increase in retirement limit. Blue, purple, student scores (retirement set at 10); green, expert

899  score (retirement set at 5), **d.** Vote threshold. Blue, F1 per instance; pink, F1 per image; dashed
900  line, vote threshold of 5 chosen for consensus calculation, **e.** Self-reported confidence score over
901  all instances (connected line, average F1), **f.** Student mean confidence score, **f.** Number of
902  images annotated by each students.

903

904  **Supplementary Figure 2: a.** Screenshots from Zooniverse project, showing the crowdsource
905  annotation user interface **b.** Total count of instance False Positive and False Negative errors
906  made by student annotators. **c.** Ground truth (top left), consensus annotation (bottom left) and ten
907  independent student annotations of an image showing a low degree of consensus; the all-or-
908  nothing pattern suggest differences in individual knowledge or experience. **d.** Trends of true
909  positive (blue), false positive (red), false negative (green) for total instances of mitochondria,
910  plotted against vote threshold. A threshold of 5 was chosen for consensus calculation.

911

912  **Supplementary Figure 3. Example of consensus algorithm splitting and merging instance**
913  **votes to create accurate consensus label maps.** See Materials and Methods for details.

914

915  **Figure 3: Challenging and diverse benchmarks for evaluating automatic instance**
916  **segmentation performance. a.** 2D representative images (left) and 3D reconstructions (right)
917  for the benchmark test sets. From top to bottom: C. elegans, Fly brain, HeLa cell, Glycolytic
918  muscle, Salivary gland, Lucchi++. Yellow arrow, membranous organelle, orange and blue
919  arrows, lightly and darkly stained mitochondria, green arrow, heavy metal precipitate, red arrow,
920  mitochondrion and tightly apposed salivary granule in the acinus **b.** Comparison of individual
921  mitochondria and box plots across benchmarks by (top to bottom): volume, branch length, cross-
922  section radius, minimum distance to neighbor (all in voxels) and contrast. Pink, C. elegans
923  n=241; yellow, Fly brain n=91; green, HeLa cell n=68; teal, Glycolytic muscle n=104; blue,
924  Salivary gland n=131; purple, Lucchi++ n=33.

925

926  **Figure 4: Deep learning model and postprocessing pipeline to create 2D or 3D instance**
927  **segmentations. a.** Schematic of Panoptic-DeepLab. A grayscale image (left) is passed through
928  the model architecture: blue boxes denote outputs from the encoder layers, black boxes output
929  from atrous spatial pyramid pooling (ASPP) layers, gray boxes output from decoder layers.

930    Outputs of the network are (left to right) semantic segmentation, up-down offsets, left-right
931    offsets and instance centers. Far right, instance segmentation created from the outputs. **b.**
932    Instance matching across adjacent slices with partially overlapping segmentation uses
933    intersection-over-union (IoU) and intersection-over-area (IoA) scores. Top row: predicted
934    segmentation of slice j (left), j+1 (right), bottom row: IoU only merging (left), IoU and IoA
935    merging (right). **c.** Merging of labels with completely missing segmentation uses a median
936    prediction over 3-7 stacks, in direction of black arrow, depending on size of gap  **d.** From top to
937    bottom: Stacked 2D segmentations before matching, after forward matching only and after
938    forward and backward matching, to merge falsely split mitochondria. Black arrows denote
939    direction of matching. **e.** An example of 3D instance segmentation of mitochondria after running
940    inference in (left to right) xy, xz, and yz directions, and far right, merged instance segmentation
941    after combining inferences from orthogonal planes.

942

943    **Supplementary Figure 4.** Example image showing difference in semantic segmentation borders
944    of closely apposed mitochondria, after linear interpolation (top) or PointRend (bottom). Blue
945    dots, PointRend sampling locations.

946

947    **Figure 5. MitoNet results on benchmarks. a.** Representative 2D images showing MitoNet
948    segmentation performance. Top to bottom: C. elegans, Fly brain, HeLa cell, Glycolytic muscle,
949    Lucchi++, Salivary gland before model finetuning, Salivary gland after model finetuning. **b.**
950    Representative ground truth and predicted segmentations from MitoNet. Red and green,
951    predicted mitochondrial instances, blue and orange, ground truth instances. Blue arrow, highly
952    branched mitochondrion, black arrow, example of segmentation expected to return a high IoU
953    but low F1 score. **c.** Left, MitoNet F1score on each of the benchmarks as a function of IoU
954    threshold; right, IoU scores **d.** Left, comparison of mean F1 score for models trained on different
955    datasets plotted against IoU threshold, right, comparison of mean IoUs achieved by models
956    trained on various datasets. All benchmarks except the salivary gland are included.

957

958    **Supplementary Figure 5. Analysis of True Positive and False Negative detections by**
959    **MitoNet on benchmarks, grouped by different mitochondrial attributes.** Top to bottom:
960    mitochondrial volume, minimum distance to nearest neighbor, branch length, branch mean cross-

961    sectional diameter (measurements by voxel) and mitochondrial contrast. Branch length and cross

962    section plots exclude the Lucchi++ benchmark, as the branches for these mitochondria were

963    shorter than the pruning threshold. Fine tuned model results are plotted for salivary gland.

964

965    **Supplementary Figure 6. Representative images of datasets that were excluded from**

966    **CEM1.5M and CEM-MitoLab datasets, in order to maintain integrity of the test set.**

967    Images to the left of the black line are volumes in the test set (fly brain, C. elegans and

968    Lucchi++). Images to the right show datasets that were excluded because they were considered

969    too similar to the benchmarks.

970

971    **Supplementary Figure 7. Screenshot of the napari plugin empanada.**

972

973    **Table 1: Performance metrics of MitoNet across benchmarks. AP, Average Precision, PQ,**

974    **Panoptic Quality**

975

976    **Table 2: Average performance of MitoNet versions pre-trained on CEM1.5M but trained**

977    **on various annotated datasets**

978

979    **Supplementary Table 1: Average mitochondrial measurements for true positive, false**

980    **negative MitoNet predictions across various benchmarks.**

981

982

983 **References**

984 1.    Peddie, C. J. & Collinson, L. M. Exploring the third dimension: Volume electron
985       microscopy comes of age. *Micron* **61**, 9–19 (2014).

986 2.    Titze, B. & Genoud, C. Volume scanning electron microscopy for imaging biological
987       ultrastructure. *Biol. Cell* **108**, 307–323 (2016).

988 3.    Scheffer, L. K. *et al.* A connectome and analysis of the adult drosophila central brain.
989       *Elife* **9**, 1–74 (2020).

990 4.    Turner, N. L. *et al.* Reconstruction of neocortex: Organelles, compartments, cells, circuits,
991       and activity. *Cell* **0**, 17 (2022).

992 5.    Yin, W. *et al.* A petascale automated imaging pipeline for mapping neuronal circuits with
993       high-throughput transmission electron microscopy. *Nat. Commun. 2020 111* **11**, 1–12
994       (2020).

995 6.    Januszewski, M. *et al.* High-precision automated reconstruction of neurons with flood-
996       filling networks. *Nat. Methods* **15**, 605–610 (2018).

997 7.    Berning, M., Boergens, K. M. & Helmstaedter, M. SegEM: Efficient Image Analysis for
998       High-Resolution Connectomics. *Neuron* **87**, 1193–1206 (2015).

999 8.    Abdollahzadeh, A., Belevich, I., Jokitalo, E., Sierra, A. & Tohka, J. DeepACSON
1000      automated segmentation of white matter in 3D electron microscopy. *Commun. Biol. 2021*
1001      *41* **4**, 1–14 (2021).

1002 9.   Dorkenwald, S. *et al.* Automated synaptic connectivity inference for volume electron
1003      microscopy. *Nat. Methods 2017 144* **14**, 435–442 (2017).

1004 10.  Funke, J. *et al.* Large Scale Image Segmentation with Structured Loss Based Deep
1005      Learning for Connectome Reconstruction. *IEEE Trans. Pattern Anal. Mach. Intell.* **41**,
1006      1669–1680 (2019).

1007 11.  Heinrich, L. *et al.* Whole-cell organelle segmentation in volume electron microscopy. *Nat.*
1008      *2021 5997883* **599**, 141–146 (2021).

1009 12.  Guay, M., Emam, Z., Anderson, A., Aronova, M. & Leapman, R. Dense cellular
1010      segmentation using 2D-3D neural network ensembles for electron microscopy. *bioRxiv*
1011      *2020.01.05.895003* (2020) doi:10.1101/2020.01.05.895003.

13. Žerovnik Mekuč, M. *et al.* Automatic segmentation of mitochondria and endolysosomes in volumetric electron microscopy data. *Comput. Biol. Med.* **119**, 103693 (2020).

14. Liu, J. *et al.* Automatic Reconstruction of Mitochondria and Endoplasmic Reticulum in Electron Microscopy Volumes by Deep Learning. *Front. Neurosci.* **14**, 599 (2020).

15. Spiers, H. *et al.* Deep learning for automatic segmentation of the nuclear envelope in electron microscopy data, trained with volunteer segmentations. *Traffic* **22**, 240–253 (2021).

16. Glancy, B., Kim, Y., Katti, P. & Willingham, T. B. The Functional Impact of Mitochondrial Structure Across Subcellular Scales. *Frontiers in Physiology* vol. 11 1462 (2020).

17. Vincent, A. E. *et al.* The Spectrum of Mitochondrial Ultrastructural Defects in Mitochondrial Myopathy. *Sci. Rep.* **6**, 1–12 (2016).

18. Vincent, A. E. *et al.* Quantitative 3D Mapping of the Human Skeletal Muscle Mitochondrial Network. *CellReports* **26**, 996-1009.e4 (2019).

19. Pernas, L. & Scorrano, L. Mito-Morphosis: Mitochondrial Fusion, Fission, and Cristae Remodeling as Key Mediators of Cellular Function. *Annu. Rev. Physiol.* **78**, 505–531 (2016).

20. Delgado, T. *et al.* Comparing 3D ultrastructure of presynaptic and postsynaptic mitochondria. *Biol. Open* **8**, (2019).

21. Stoldt, S. *et al.* Spatial orchestration of mitochondrial translation and OXPHOS complex assembly. *Nat. Cell Biol.* **20**, 528–534 (2018).

22. Meyer, J. N., Leuthner, T. C. & Luz, A. L. Mitochondrial fusion, fission, and mitochondrial toxicity. *Toxicology* **391**, 42–53 (2017).

23. Zhang, L. *et al.* Altered brain energetics induces mitochondrial fission arrest in Alzheimer's Disease. *Sci. Rep.* **6**, 1–12 (2016).

24. Siegmund, S. E. *et al.* Three-Dimensional Analysis of Mitochondrial Crista Ultrastructure in a Patient with Leigh Syndrome by In Situ Cryoelectron Tomography. *ISCIENCE* **6**, 83–91 (2018).

25. Casser, V., Kang, K., Pfister, H. & Haehn, D. Fast Mitochondria Segmentation for Connectomics. *arXiv1812.06024 [cs]* (2018).

26. Wei, D. *et al.* MitoEM Dataset: Large-scale 3D Mitochondria Instance Segmentation from

EM Images. *Med. Image Comput. Comput. Assist. Interv.* **12265**, 66 (2020).

27. Müller, A. *et al.* 3D FIB-SEM reconstruction of microtubule–organelle interaction in whole primary mouse β cells. *J. Cell Biol.* **220**, (2021).

28. Buhmann, J. *et al.* Automatic Detection of Synaptic Partners in a Whole-Brain Drosophila EM Dataset. *bioRxiv* 2019.12.12.874172 (2019) doi:10.1101/2019.12.12.874172.

29. Perez, A. J. *et al.* A workflow for the automatic segmentation of organelles in electron microscopy image stacks. *Front. Neuroanat.* **8**, 126 (2014).

30. Guay, M. D. *et al.* Dense cellular segmentation for EM using 2D–3D neural network ensembles. *Sci. Reports 2021 111* **11**, 1–11 (2021).

31. Leonard, A. P. *et al.* Quantitative analysis of mitochondrial morphology and membrane potential in living cells using high-content imaging, machine learning, and morphological binning. *Biochim. Biophys. Acta - Mol. Cell Res.* **1853**, 348–360 (2015).

32. Nikolaisen, J. *et al.* Automated Quantification and Integrative Analysis of 2D and 3D Mitochondrial Shape and Network Properties. *PLoS One* **9**, e101365 (2014).

33. Talwar, A. *et al.* A Topological Nomenclature for 3D Shape Analysis in Connectomics.

34. Miyazono, Y. *et al.* Uncoupled mitochondria quickly shorten along their long axis to form indented spheroids, instead of rings, in a fission-independent manner OPEN. *Sci. REPORtS* / **8**, 350 (2018).

35. Bleck, C. K. E., Kim, Y., Willingham, T. B. & Glancy, B. Subcellular connectomic analyses of energy networks in striated muscle. *Nat. Commun.* **9**, (2018).

36. Abrisch, R. G., Gumbin, S. C., Wisniewski, B. T., Lackner, L. L. & Voeltz, G. K. Fission and fusion machineries converge at ER contact sites to regulate mitochondrial morphology. *J. Cell Biol.* **219**, (2020).

37. Tamada, H. *et al.* Three-dimensional analysis of somatic mitochondrial dynamics in fission-deficient injured motor neurons using FIB/SEM. *J. Comp. Neurol.* **525**, 2535–2548 (2017).

38. Conrad, R. & Narayan, K. CEM500K, a large-scale heterogeneous unlabeled cellular electron microscopy image dataset for deep learning. *Elife* **10**, (2021).

39. Xu, C. S. *et al.* An open-access volume electron microscopy atlas of whole cells and tissues. *Nat. 2021 5997883* **599**, 147–151 (2021).

40. Lucchi, A., Li, Y. & Fua, P. Learning for structured prediction using approximate

1074    subgradient descent with working sets. *Proc. IEEE Comput. Soc. Conf. Comput. Vis.*
1075    *Pattern Recognit.* 1987–1994 (2013) doi:10.1109/CVPR.2013.259.

1076 41.  Riddle, D. L. *et al. C. elegans II. Cold Spring Harbor Monograph Series, Vol. 33* (Cold
1077    Spring Harbor Laboratory Press, 1997).

1078 42.  Cheng, B. *et al.* Panoptic-DeepLab: A Simple, Strong, and Fast Baseline for Bottom-Up
1079    Panoptic Segmentation. *arxiv1911.10194v2 [cs]* (2019).

1080 43.  Ronneberger, O., Fischer, P. & Brox, T. U-net: Convolutional networks for biomedical
1081    image segmentation. in *Lecture Notes in Computer Science (including subseries Lecture*
1082    *Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* vol. 9351 234–241
1083    (Springer Verlag, 2015).

1084 44.  He, K., Gkioxari, G., Dollár, P. & Girshick, R. Mask R-CNN. *IEEE Trans. Pattern Anal.*
1085    *Mach. Intell.* **42**, 386–397 (2020).

1086 45.  Cheng, B., Misra, I., Schwing, A. G., Kirillov, A. & Girdhar, R. Masked-attention Mask
1087    Transformer for Universal Image Segmentation. (2021) doi:10.48550/arxiv.2112.01527.

1088 46.  Kirillov, A., Wu, Y., He, K. & Girshick, R. PointRend: Image Segmentation As
1089    Rendering. in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
1090    *Recognition (CVPR)* 9799–9808 (2020).

1091 47.  Kuhn, H. W. The Hungarian method for the assignment problem. *Nav. Res. Logist. Q.* **2**,
1092    83–97 (1955).

1093 48.  Conrad, R., Lee, H. & Narayan, K. Enforcing Prediction Consistency Across Orthogonal
1094    Planes Significantly Improves Segmentation of FIB-SEM Image Volumes by 2D Neural
1095    Networks. *Microsc. Microanal.* 1–4 (2020) doi:10.1017/s143192762002053x.

1096 49.  Stringer, C., Wang, T., Michaelos, M. & Pachitariu, M. Cellpose: a generalist algorithm
1097    for cellular segmentation. *Nat. Methods* **18**, 100–106 (2021).

1098 50.  Vincent, L., Vincent, L. & Soille, P. Watersheds in Digital Spaces: An Efficient
1099    Algorithm Based on Immersion Simulations. *IEEE Trans. Pattern Anal. Mach. Intell.* **13**,
1100    583–598 (1991).

1101 51.  He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. in
1102    *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern*
1103    *Recognition* vols 2016-Decem 770–778 (IEEE Computer Society, 2016).

1104 52.  Caron, M. *et al.* Unsupervised Learning of Visual Features by Contrasting Cluster

1105          Assignments. *Adv. Neural Inf. Process. Syst.* **2020**-**December**, (2020).

1106 53.   Xiao, C. *et al.* Automatic mitochondria segmentation for EM data using a 3D supervised
1107       convolutional network. *Front. Neuroanat.* **12**, 92 (2018).

1108 54.   Bleck, C. K. E., Kim, Y., Willingham, T. B. & Glancy, B. Subcellular connectomic
1109       analyses of energy networks in striated muscle. *Nat. Commun.* **9**, 1–11 (2018).

1110 55.   Sarkans, U. *et al.* REMBI: Recommended Metadata for Biological Images—enabling
1111       reuse of microscopy data in biology. *Nat. Methods* 1–5 (2021) doi:10.1038/s41592-021-
1112       01166-8.

1113 56.   Franco-Barranco, D., Muñoz-Barrutia, A. & Arganda-Carreras, I. Stable deep neural
1114       network architectures for mitochondria segmentation on electron microscopy volumes.
1115       *arXiv:2104.03577 [eess.IV]* (2021).

1116 57.   Narayan, K. & Subramaniam, S. Focused ion beams in biology. *Nat. Methods* **12**, 1021
1117       (2015).

1118 58.   Conrad, R. & Narayan, K. Cem500k, a large-scale heterogeneous unlabeled cellular
1119       electron microscopy image dataset for deep learning. *Elife* **10**, (2021).

1120 59.   Sarkans, U. *et al.* REMBI: Recommended Metadata for Biological Images—enabling
1121       reuse of microscopy data in biology. *Nat. Methods 2021 1812* **18**, 1418–1422 (2021).

1122 60.   Kikinis, R., Pieper, S. D. & Vosburgh, K. G. 3D Slicer: A Platform for Subject-Specific
1123       Image Analysis, Visualization, and Clinical Support. in *Intraoperative Imaging and*
1124       *Image-Guided Therapy* 277–289 (Springer New York, 2014). doi:10.1007/978-1-4614-
1125       7657-3_19.

1126 61.   Kasthuri, N. *et al.* Saturated Reconstruction of a Volume of Neocortex. *Cell* **162**, 648–661
1127       (2015).

1128 62.   Žerovnik Mekuč, M. *et al.* Automatic segmentation of mitochondria and endolysosomes
1129       in volumetric electron microscopy data. *Comput. Biol. Med.* **119**, 103693 (2020).

1130 63.   Lee, T. C., Kashyap, R. L. & Chu, C. N. Building Skeleton Models via 3-D Medial
1131       Surface Axis Thinning Algorithms. *CVGIP Graph. Model. Image Process.* **56**, 462–478
1132       (1994).

1133 64.   Caron, M. *et al.* Unsupervised Learning of Visual Features by Contrasting Cluster
1134       Assignments. *arXiv* (2020).

1135 65.   He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition.

*Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* **2016**-**December**, 770–778 (2015).

66. Kirillov, A., Wu, Y., He, K. & Girshick, R. PointRend: Image Segmentation as Rendering. (2019).

67. Smith, L. N. A disciplined approach to neural network hyper-parameters: Part 1 -- learning rate, batch size, momentum, and weight decay. *arxiv1803.09820 [cs]* (2018).

68. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. *7th Int. Conf. Learn. Represent. ICLR 2019* (2017).

69. Ghiasi, G. *et al. Simple Copy-Paste is a Strong Data Augmentation Method for Instance Segmentation*. https://cocodataset.org/.

# Figure 1
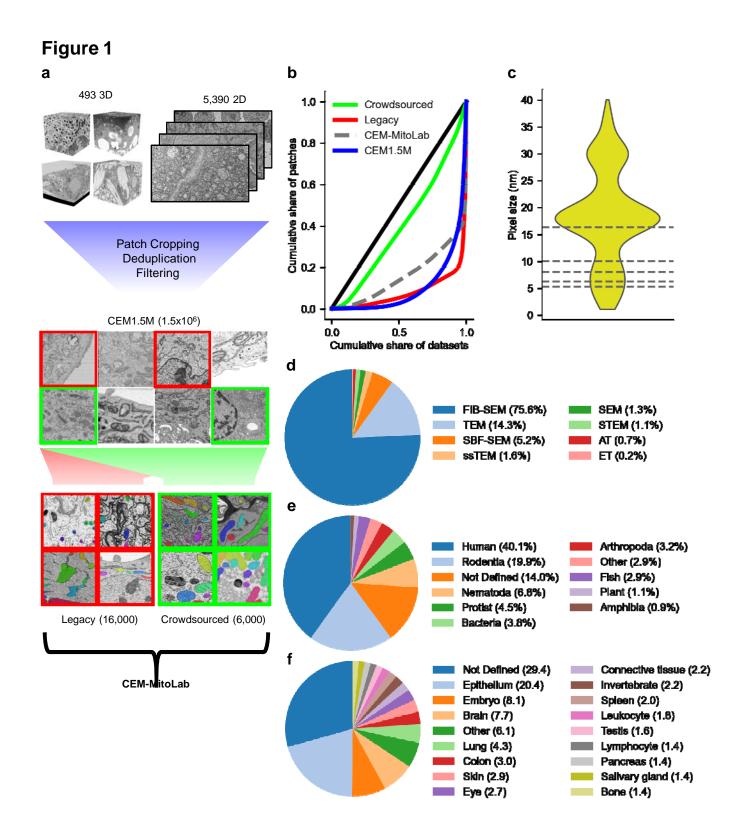
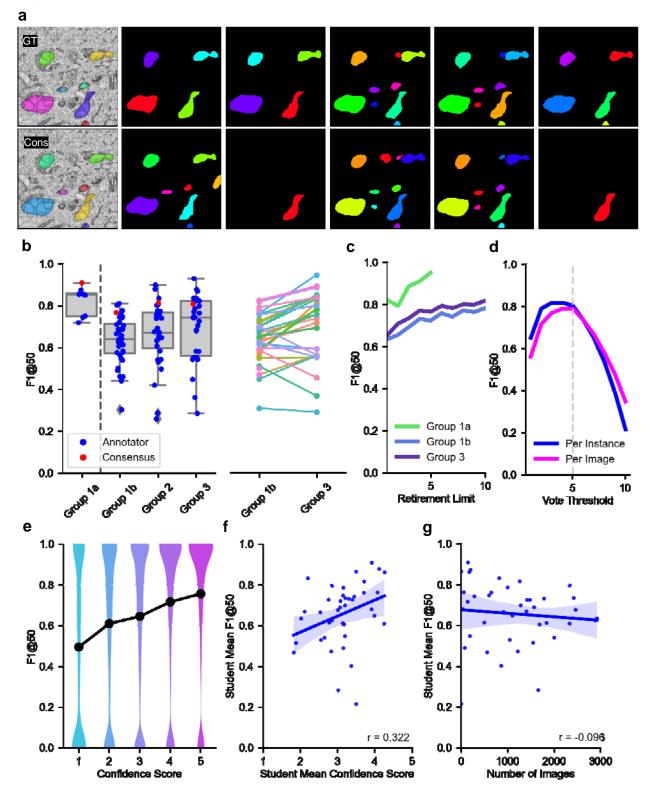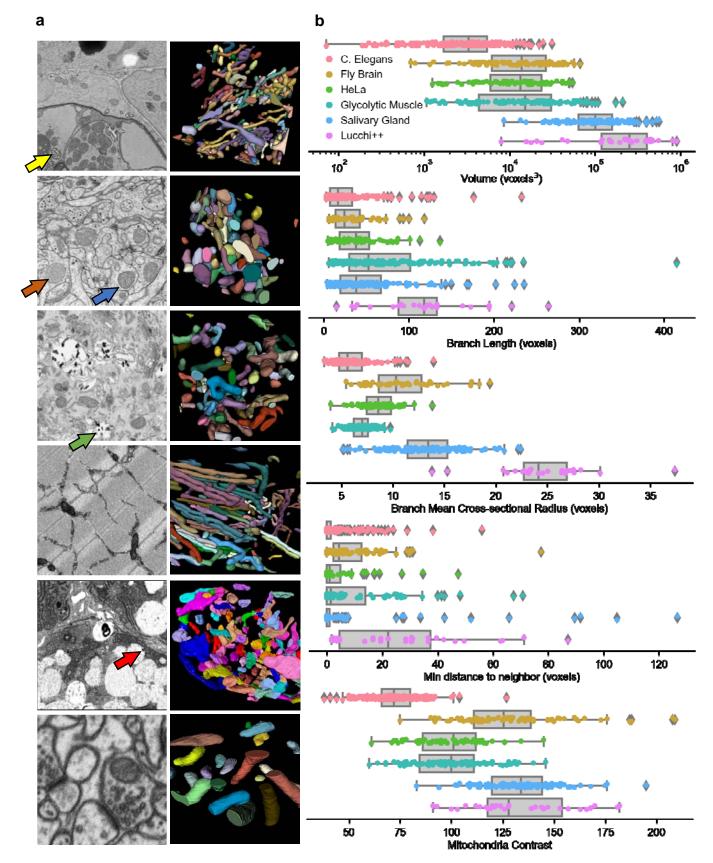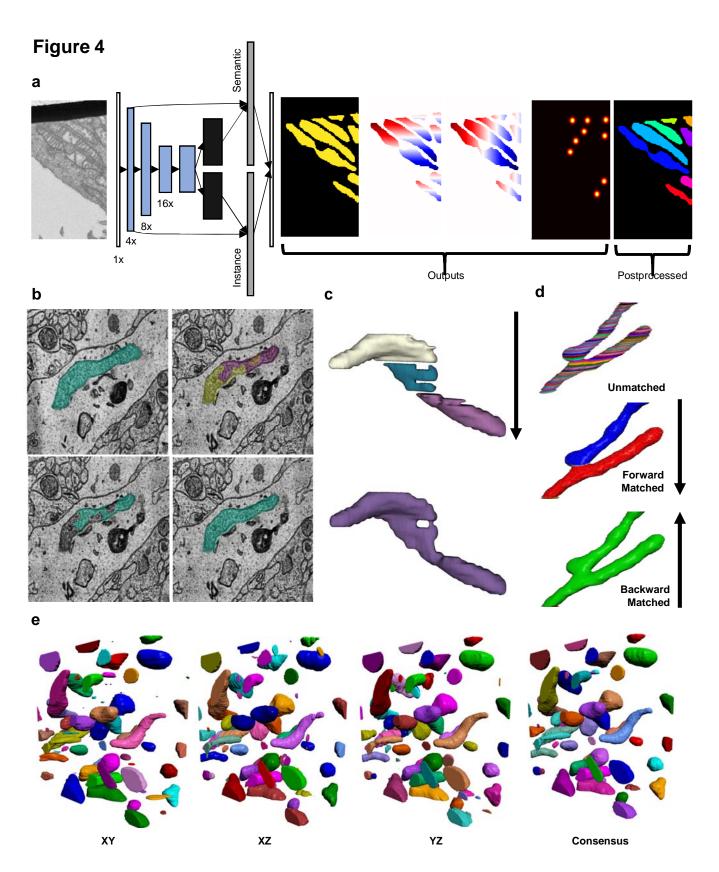# Figure 2

# Figure 3

## Figure 4

# Figure 5