# Deep Neural Networks to Register and Annotate the Cells of the *C. elegans* Nervous System

**Adam A. Atanas[1], Alicia Kun-Yang Lu[1], Jungsoo Kim[1], Saba Baskoylu[1], Di Kang[1], Talya S. Kramer[1], Eric Bueno[1], Flossie K. Wan[1], Steven W. Flavell[1,*]**

[1]Picower Institute for Learning & Memory, Department of Brain & Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA
[*]Corresponding Author: flavell@mit.edu

**ABSTRACT**

Aligning and annotating the heterogeneous cell types that make up complex cellular tissues remains a major challenge in the analysis of biomedical imaging data. Here, we present a series of deep neural networks that allow for automatic non-rigid registration and cell identification in the context of the nervous system of freely-moving *C. elegans*. A semi-supervised learning approach was used to train a *C. elegans* registration network (BrainAlignNet) that aligns pairs of images of the bending *C. elegans* head with single pixel-level accuracy. When incorporated into an image analysis pipeline, this network can link neuronal identities over time with 99.6% accuracy. A separate network (AutoCellLabeler) was trained to annotate >100 neuronal cell types in the *C. elegans* head based on multi-spectral fluorescence of genetic markers. This network labels >100 different cell types per animal with 98% accuracy, exceeding individual human labeler performance by aggregating knowledge across manually labeled datasets. Finally, we trained a third network (CellDiscoveryNet) to perform unsupervised discovery and labeling of >100 cell types in the *C. elegans* nervous system by analyzing unlabeled multi-spectral imaging data from many animals. The performance of CellDiscoveryNet matched that of trained human labelers. These tools will be useful for a wide range of applications in *C. elegans* research and should be straightforward to generalize to many other applications requiring alignment and annotation of dense heterogeneous cell types in complex tissues.

## INTRODUCTION

Optical imaging of dense cellular tissues is widespread in biomedical research. Recently developed methods to label cells with highly multiplexed fluorescent probes should soon make it feasible to determine the heterogeneous cell types in any given sample[1–3]. However, it remains challenging to extract critical information about cell identity and position from fluorescent imaging data. Aligning images within or across animals that have non-rigid deformations can be inefficient and lack cellular-level accuracy. Additionally, annotating cell types in a given sample can involve time-consuming manual labeling and often only results in coarse labeling of the main cell classes, rather than full annotation of the vast number of defined cellular subtypes.

Deep neural networks provide a promising avenue for aligning and annotating complex images of fluorescently-labeled cells with high levels of efficiency and accuracy[4]. Deep learning has generated high-performance tools to segment cells from background in images[5,6]. In addition, deep learning approaches have proven useful for non-rigid image registration in the context of medical image alignment[7]. However, this has not been as widely applied to align images of fluorescently labeled cells, which requires micron-level accuracy. Automated cell annotation using clustering approaches, for example applied to single-cell RNA sequencing data, has been widely adopted[8]. Recent studies have also shown the feasibility of using deep learning applied on image features[9] or raw imaging data to label major cell classes[8,10,11]. However, these methods are still not sufficiently advanced to label the potentially hundreds of cellular subtypes in images of complex tissues. In addition, fully unsupervised discovery of the many distinct cell types in cellular imaging data remains an unsolved challenge.

There is considerable interest in using these methods to automatically align and annotate cells in the nervous system of *C. elegans*, which consists of 302 uniquely identifiable neurons[12–14]. The optical transparency of the animal enables *in vivo* imaging of fluorescent indicators of neural activity at brain-wide scale.[15,16] Advances in closed-loop tracking made this imaging feasible in freely-moving animals.[17,18] These approaches are being used to map the relationship between brain-wide activity and flexible behavior (reviewed in[19,20]). However, the animal bends and warps its head as it moves, resulting in non-rigid deformations of the densely-packed cells in its nervous system. Fully automating the alignment and annotation of cells in *C. elegans* imaging data would facilitate high-throughput and high-SNR brain-wide calcium imaging. These methods could also be applied to unsolved problems in quantifying reporter gene expression, developmental trajectories, and more.

Previous studies have described methods to align and annotate cells in multi-cellular imaging datasets from *C. elegans* and species with related imaging challenges like *Hydra*. Datasets from freely-moving animals pose an especially challenging case. Methods for aligning cells across timepoints in moving datasets include approaches that link neurons across adjacent timepoints[21–23], as well as approaches that use signal demixing[24], alignment of body position markers using anatomical constraints[25,26], or registration/clustering/matching based on features of the neurons, such as their centroid positions[27–32]. Targeted data augmentation combined with deep learning applied to raw images has recently been used to reduce manual labeling time during cell alignment.[33] Deep learning applied to raw images has also been used to identify specific image features, like multi-cellular structures in *C. elegans*.[34] We have previously applied non-rigid registration to full fluorescent images from brain-wide calcium imaging datasets to perform

2

79    neuron alignment, but performing this complex image alignment via gradient descent is very
80    slow, taking multiple days to process a single animal's data even on a computing cluster[35].  In
81    summary, all of these current methods for neuron alignment are constrained by a tradeoff
82    between alignment accuracy and time spent processing each dataset, either due to manually
83    labeling subsets of neurons or computing the complex alignments that actually yield >95%
84    alignment accuracy.
85
86    For *C. elegans* neuron class annotation, ground-truth measurements of neurons' locations in the
87    head have allowed researchers to develop atlases describing the statistical likelihood of finding a
88    given neuron in a given location[36–42]. Some of these atlases have utilized the NeuroPAL
89    transgene in which four fluorescent proteins are expressed in genetically-defined sets of cells,
90    allowing users to manually determine their identity based on multi-spectral fluorescence and
91    neuron position[40–42]. However, this manual labeling is time-consuming (hours per dataset), and
92    statistical approaches to automate neuron annotation based on manual labeling have still not
93    achieved human-level performance (>95% accuracy).
94
95    Here we describe deep neural networks that solve these alignment and annotation tasks. First, we
96    trained a neural network (BrainAlignNet) that can perform non-rigid registration to align images
97    of the worm's head from different timepoints in freely-moving data. It is >600-fold faster than
98    our previous gradient descent-based approach using elastix[35] and aligns neurons with 99.6%
99    accuracy. Second, we trained a neural network (AutoCellLabeler) that annotates the identity of
100   each *C. elegans* neuron in the head based on multi-spectral NeuroPAL labels. This network
101   achieves 98% accuracy; versions trained on subsets of the fluorescent channels in NeuroPAL
102   also achieve high performance. Finally, we trained a different network (CellDiscoveryNet) that
103   can perform unsupervised discovery and labeling of >100 cell types of the *C. elegans* nervous
104   system by comparing unlabeled NeuroPAL images across animals. Overall, our results reveal
105   how to train neural networks to align and annotate cells in complex cellular imaging data with
106   high performance.
107
108   **RESULTS**
109
110   **BrainAlignNet: a neural network that registers cells in the deforming head of freely-**
111   **moving *C. elegans***
112   When analyzing neuronal calcium imaging data, it is essential to accurately link neurons'
113   identities over time to construct reliable calcium traces. This task is challenging in freely-moving
114   animals where the nervous system is warped on sub-second timescales by animal movement.
115   Therefore, we sought to develop a fast and accurate method to perform non-rigid image
116   registration that can deal with these warping images. Previous studies have described such
117   methods for non-rigid registration of point clouds (e.g. neuron centroid positions)[28–30,43], but, as
118   we describe below, we found that performing full image alignment allows for higher accuracy
119   neuron position alignments.
120
121   To solve this task, we used a previously-described network architecture[44,45] that takes as input a
122   pair of 3-D images (i.e. volumes of fluorescent imaging data of the head of the worm) from
123   different timepoints of the same neural recording (Fig. 1A). The network is tasked with
124   determining how to warp one 3-D image (termed the "moving image") so that it resembles the

125  other 3-D image (termed the "fixed" image). Specifically, the network outputs a dense
126  displacement field (DDF), a pixel-wise coordinate transformation function designed to indicate
127  which points in the moving and fixed images are the same (see Methods). The moving image is
128  then transformed through this DDF to create a warped moving image, which should look like the
129  fixed image. This network was selected because its LocalNet architecture (a modified 3-D U-
130  Net) allows it to do the feature extraction and image reconstruction necessary to solve the task.
131  To train and evaluate the network, we used data from freely-moving animals expressing both
132  pan-neuronal NLS-GCaMP and NLS-tagRFP, but only provided the tagRFP images to the
133  network, since this fluorophore's brightness should remain static over time. Since Euler
134  registration of images (rotation and translation) is simple, we performed Euler registration on the
135  images using a GPU-accelerated grid search prior to inputting them into the network. During
136  training, we also provided the network with the locations of the centroids of matched neurons
137  found in both images, which were available for these training and validation data since we had
138  previously used gradient descent to solve those registration problems ("registration problem"
139  here is defined as a single image pair that needs to be aligned) and link neurons' identities[35]. The
140  centroid locations are only used for network training and are not required for the network to
141  solve registration problems after training. The loss function that the network was tasked with
142  minimizing had three components: (1) image loss: the Local squared zero-Normalized Cross-
143  Correlation (LNCC) of the fixed and warped moving RFP images, which takes on a lower value
144  when the images are more similar; (2) centroid alignment loss: the average of the Euclidean
145  distances between the matched centroid pairs, where lower values indicate better alignment; and
146  (3) regularization loss: a term that increases the overall loss the more that the images are
147  deformed in a non-rigid manner (in particular, penalizing image scaling and scrambling of
148  adjacent pixels; see Methods).
149
150  We trained and validated the network on 5,176 and 1,466 image pairs, respectively, over 300
151  epochs, at which point the validation loss plateaued (Fig. 1B). We then evaluated network
152  performance on a separate set of 447 image pairs reserved for testing that were recorded from
153  five different animals. On average, the network improved the Normalized Cross-Correlation
154  (NCC) from 0.577 in the input image pairs to 0.947 in the registered image pairs – the maximum
155  achievable score is 1 (Fig. 1C shows example of centroid positions; Fig. 1D shows image
156  example; Fig. 1E shows both). The average distance between aligned centroids was 1.45 pixels
157  (Fig. 1F). These results were only modestly different depending on the animal or the exact
158  registration problem being solved (Extended Data Fig. 1A-C).
159
160  To determine which features of the network were critical for its performance, we trained the
161  network under conditions where we omitted either the centroid alignment loss, the regularization
162  loss, or the image loss. In the first case, the network would not be able to learn based on whether
163  the neuron centroids were well-aligned; in the second case, there would be no constraints on the
164  network performing any type of deformation to solve the task; in the third case, the deformations
165  that the network learned to apply could only be learned from the alignment of the centroids, not
166  the raw tagRFP images. Registration performance of each network was evaluated using the NCC
167  and centroid distance, which quantify the quality of tagRFP image alignment and centroid
168  alignment, respectively (Fig. 1F). While the NCC scores were similar for the full network and
169  the no-regularization and no-centroid alignment networks, other performance metrics like
170  centroid distance were significantly impaired by the absence of centroid alignment loss or

171 regularization loss (Fig. 1E-F). This suggests that in the absence of centroid alignment loss or
172 regularization loss, the network learns how to align the tagRFP images, but does so using
173 unnatural deformations that do not reflect how the worm bends. In the case of the no-image loss
174 network, all performance metrics, including both image and centroid alignment, were impaired
175 compared to the full network (Fig. 1F). This suggests that allowing the network to learn how to
176 warp the RFP images also enhances the network's ability to learn how to align the neuron
177 positions (i.e. centroids).
178
179 The finding that the centroid positions were precisely aligned by the full network indicates that
180 the centers of the neurons were correctly registered by the network. However, it does not ensure
181 that all of the pixels that comprise a neuron are being correctly registered, which could be
182 important for subsequent feature extraction from the aligned images. For example, it is formally
183 possible to have perfect RFP image alignment in a context where the pixels from one neuron in
184 the moving RFP image are scrambled to multiple neuron locations in the warped moving RFP
185 image. In fact, we observed this in our first efforts to build such a network, where the loss
186 function was only composed of the image loss. As an additional control to test for this possibility
187 in our trained networks, we examined the network's performance on data from a different strain
188 that expresses pan-neuronal NLS-mNeptune (analogous to the pan-neuronal NLS-tagRFP) and
189 *eat-4*::NLS-GFP, which is expressed in ~40% of the neurons in the *C. elegans* head (Fig. 1G
190 shows example image). If the pixels within the neurons are being correctly registered, then
191 applying image registration to the GFP channel for these image pairs should result in highly
192 correlated images (i.e., a high NCC value close to 1). If the pixels within neurons are being
193 scrambled, then these images should not be well-aligned. We used the DDF that the network
194 learned from pan-neuronal mNeptune data to register the corresponding *eat-4*::NLS-GFP images
195 from the same timepoints and found that this resulted in high-quality GFP image alignment (Fig.
196 1H). In contrast, while the no-centroid alignment and no-regularization networks output a DDF
197 that successfully aligned the RFP images, applying this DDF to corresponding GFP images
198 resulted in poor GFP image registration (Fig. 1H shows that the no-centroid alignment network
199 aligns the RFP channel, but not the GFP channel, in the *eat-4*::NLS-GFP strain). This further
200 suggests that these reduced networks lacking centroid alignment or regularization loss are
201 aligning the RFP images through unnatural image deformations. Altogether, these results suggest
202 that the full **Brain Align**ment Neural Net**work** (**BrainAlignNet**) can perform non-rigid
203 registration on pairs of images from freely-moving brain-wide calcium imaging data.
204
205 The registration problems included in the training, validation, and test data above were pulled
206 from a set of registration problems that we had been able to solve with gradient descent (example
207 images in Extended Data Fig. 1D). These problems did not include the most challenging cases,
208 for example when the two images to be registered had the worm's head bent in opposite
209 directions (though we note that it did include substantial non-rigid deformations). We next asked
210 whether a network trained on arbitrary registration problems, including those that were not
211 solvable with gradient descent (example images in Extended Data Fig. 1E), could obtain high
212 performance. For this test, we also omitted the Euler registration step that we performed in
213 advance of network training, since the goal was to test whether this network architecture could
214 solve any arbitrary *C. elegans* head alignment problem. For this analysis, we used the same loss
215 function as the successful network described above. We also increased the amount of training
216 data from 5,176 to 335,588 registration problems. The network was trained for 300 epochs, at

217  which point the validation loss plateaued. However, the test performance of the network was not
218  high in terms of image alignment or centroid alignment (Extended Data Fig. 1F). This suggests
219  that additional approaches may be necessary to solve these more challenging registration
220  problems. Overall, our results suggest that, provided that there is an appropriate loss function, a
221  deep neural network can perform non-rigid registration problems to align neurons across the *C.*
222  *elegans* head with high speed and accuracy.
223
224  **Integration of BrainAlignNet into a complete calcium imaging processing pipeline**

225  The above results suggest that BrainAlignNet can perform high quality image alignments. These
226  alignments are only valuable if they enable accurate linking of neurons over time. To test
227  whether performance was sufficient for this, we incorporated BrainAlignNet into our existing
228  image analysis pipeline for brain-wide calcium imaging data and compared the results to our
229  previously-described pipeline, which used gradient descent to solve image registration[35]. This
230  image analysis pipeline, the **A**utomated **N**euron **T**racking **S**ystem for **U**nconstrained **N**ematodes
231  (ANTSUN), includes steps for neuron segmentation (via a 3D U-Net), image registration, and
232  linking of neurons' identities (Fig. 2A). Several steps are required to link neurons' identities
233  based on image registration. First, image registration defines a coordinate transformation
234  between the two images, which is then applied to the segmented neuron ROIs, warping them into
235  a common coordinate frame. To link neurons' identities over time, we then build a N-by-N
236  matrix (where N is the number of all segmented neuron ROIs at all timepoints in a given
237  recording) with the following structure: (1) Enter zero if the ROIs were in an image pair that was
238  not registered (we do not attempt to solve all registration problems, as this is unnecessary); (2)
239  Enter zero if the ROIs were from a registered image pair, but the registration-warped ROI did not
240  overlap with the fixed ROI; and (3) Otherwise, enter a heuristic value indicating the confidence
241  that the ROIs are the same neurons based on several ROI features. These features include
242  similarity of ROI positions and sizes, similarity of red channel brightness, registration quality
243  (computed as NCC of the red channel images), a penalty for overly nonlinear registration
244  transformations, and a penalty if ROIs were displaced over large distances during alignment.
245  Finally, custom hierarchical clustering is applied to the matrix to generate clusters consisting of
246  the ROIs that reflect the same neuron recorded at different timepoints. Calcium traces are then
247  constructed from all of these timepoints, normalizing the GCaMP signal to the tagRFP signal
248  (Fig. 2B-D shows example GCaMP dataset and GFP control datasets). We term the ANTSUN
249  pipeline with gradient descent registration ANTSUN 1.4[35,46] and the version with BrainAlignNet
250  registration ANTSUN 2.0 (Fig. 2A).
251
252  We ran a series of control datasets through both versions of ANTSUN to benchmark their results.
253  The first was from the previously-described animals with pan-neuronal NLS-mNeptune and *eat-*
254  *4*::NLS-GFP. The resulting GFP traces from these recordings allow us to quantify the number of
255  timepoints where the neuron identities are not accurately linked together into a single trace (Fig.
256  2B shows example dataset). Specifically, in this strain, this type of error can be easily detected
257  since it can result in a low-intensity GFP neuron (*eat-4-*) suddenly having a high-intensity value
258  when the trace mistakenly incorporates data from a high-intensity neuron (*eat-4+*), or vice versa.
259  We computed this error rate, taking into account the overall similarity of GFP intensities (i.e.
260  since we can only observe errors when GFP- and GFP+ neurons are combined into the same
261  trace). For both versions of ANTSUN, the error rates were <0.5%, suggesting that >99.5% of
262  timepoints reflect correctly linked neurons (Fig. 2E).

263

264  We next estimated the SNR of the data collected from ANTSUN 2.0, as compared to ANTSUN
265  1.4. Here, we processed data from three pan-neuronal GCaMP animals and compared them to
266  three animals expressing pan-neuronal GFP, in place of GCaMP. The relative signal fluctuations
267  in GCaMP traces versus GFP traces (the GFP traces should ideally be flat) can provide an
268  indication of the entire recording/processing pipeline's SNR (Fig. 2C-D show examples). Results
269  were similar for ANTSUN 1.4 and 2.0, which indicates that incorporating BrainAlignNet did not
270  impair the SNR of the data (Fig. 2F). ANTSUN 2.0 also successfully extracted traces from a
271  similar number of neurons (Fig. 2G). However, while ANTSUN 1.4 requires 250 CPU days per
272  dataset for registration, ANTSUN 2.0 only requires 9 GPU hours, reflecting a >600-fold increase
273  in computation speed (Fig. 2H). These results suggest that ANTSUN 2.0, which uses
274  BrainAlignNet, provides a massive speed improvement in extracting neural data from GCaMP
275  recordings without compromising the SNR or accuracy of the data.

276

277  **AutoCellLabeler: a neural network that automatically annotates >100 neuron classes in the**
278  ***C. elegans* head from multi-spectral fluorescence**

279  We next turned our attention to annotating the identities of the recorded neurons in these brain-
280  wide calcium imaging data. *C. elegans* neurons have fairly stereotyped positions in the heads of
281  adult animals, though fully accurate inference of neural identity from position alone has not been
282  shown to be possible. Fluorescent reporter gene expression using well-defined genetic drivers
283  can provide additional information to assist with neuron annotation. The NeuroPAL strain is
284  especially useful in this regard. It expresses pan-neuronal NLS-tagRFP, but also has expression
285  of NLS-mTagBFP2, NLS-CyOFP1, and NLS-mNeptune2.5 under a set of well-chosen genetic
286  drivers (example image in Fig. 3A)[40]. With proper training, humans can manually label the
287  identities of most neurons in this strain using neuron position and multi-spectral fluorescence.
288  For most of the brain-wide recordings collected using our calcium imaging platform, we used a
289  previously characterized strain with a pan-neuronal NLS-GCaMP7F transgene crossed into
290  NeuroPAL[35] . While freely-moving recordings were conducted with only NLS-GCaMP and
291  NLS-tagRFP data acquisition, animals were immobilized at the end of each recording in order to
292  capture multi-spectral fluorescence. Humans could manually label many neurons' identities in
293  these multi-spectral images, and the image registration approaches described above could map
294  the ROIs in the immobilized data to ROIs in the freely-moving recordings to match neuron
295  identity to GCaMP traces.

296

297  Manual annotation of NeuroPAL images is time-consuming. First, to perform accurate labeling,
298  the individual needs substantial amounts of training. Even after being fully trained, labeling all
299  the ROIs in one NeuroPAL animal can take 3-5 hours. In addition, different individuals have
300  different degrees of knowledge or confidence in labeling certain cell classes. For these reasons, it
301  was desirable to automate NeuroPAL labeling, using datasets that had previously been labeled by
302  a panel of human labelers. In particular, the labels that they provided with a high degree of
303  confidence in their accuracy would be most useful for training an automated labeling network.
304  Previous studies have developed statistical approaches for semi-automated labeling to label
305  neural identity from NeuroPAL images, but the maximum precision that we are aware of is 90%
306  without manual correction[40].

307

308   We trained a 3-D U-Net[47] to label the *C. elegans* neuron classes in a given NeuroPAL 3-D
309   image. As input, the network received four fluorescent 3-D images from the head of each worm:
310   pan-neuronal NLS-tagRFP, plus the NLS-mTagBFP2, NLS-CyOFP1, and NLS-mNeptune2.5
311   images that label stereotyped subsets of neurons (Fig. 3A). During training, the network also
312   received the human-annotated labels of which pixels belong to which neurons. Humans provided
313   ROI-level labels and the boundaries of each ROI were determined using a previously-described
314   neuron segmentation network[35] trained to label all neurons in a given image (agnostic to their
315   identity). Finally, during training the network also received an array indicating the relative
316   weight to assign each pixel during training (Fig. 3B). This was incorporated into a pixel-
317   weighted cross-entropy loss function (lower values indicate more accurate labeling of each
318   pixel), summing across the pixels in a weighted manner. Pixel weighting was adjusted as
319   follows: (1) background was given extremely low weight; (2) ROIs that humans were not able to
320   label were given low weight; (3) all other ROIs received higher weight, proportional to the
321   subjective confidence that the human had in assigning the label to the ROI and the rarity of the
322   label. Regarding this latter point, neurons that were less frequently labeled by human annotation
323   received higher weight so that the network could potentially learn how to classify these neurons
324   from fewer labeled examples.
325
326   We trained the network over 300 epochs using a training set of 81 annotated images and a
327   validation set of 10 images (Fig. 3C). Because the size of the training set was fairly small, we
328   augmented the training data using both standard image augmentations (rotation, flipping, adding
329   gaussian noise, etc.) and a custom augmentation where the images were warped in a manner to
330   approximate worm head bending (see Methods). Because this ***Automatic Cell Labeling Network***
331   (**AutoCellLabeler**) labels individual pixels, it was necessary to convert these pixel-wise
332   classifications into ROI-level classifications. AutoCellLabeler outputs its confidence in its label
333   for each pixel, and we noted that the network's confidence for a given ROI was highest near the
334   center of the ROI (Fig. 3D). Therefore, to determine ROI-level labels, we took a weighted
335   average of the pixel-wise labels within an ROI, weighing the center pixels more strongly. The
336   overall confidence of these pixel scores was also used to compute a ROI-level confidence score,
337   reflecting the network's confidence that it labeled the ROI correctly. Finally, after all ROIs were
338   assigned a label, heuristics were applied to identify and delete problematic labels. Labels were
339   deleted if (1) the network already labeled another ROI as that label with higher confidence; (2)
340   the label was present too infrequently in the network's training data; (3) the network labeled that
341   ROI as something other than a neuron (e.g. a gut granule or glial cell, which we supplied as valid
342   labels during training); or (4) the network confidently predicted different parts of the ROI as
343   different labels.
344
345   We evaluated the performance of the network on 11 separate datasets that were reserved for
346   testing. We assessed the accuracy of AutoCellLabeler on the subset of ROIs with high-
347   confidence human labels (subjective confidence scores of 4 or 5, on a scale from 1-5). On these
348   neurons, average network confidence was 96.8% and its accuracy was 97.1%. We furthermore
349   observed that the network was more confident in its correct labels (average confidence 97.3%)
350   than its incorrect labels (average confidence 80.7%; Fig. 3E). More generally, AutoCellLabeler
351   confidence was highly correlated with its accuracy (Fig. 3F). Indeed, excluding the neurons
352   where the network assigns low (<75%) confidence increased its accuracy to 98.1% (Extended
353   Data Fig. 2A displays the full accuracy-recall tradeoff curve). Under this confidence threshold

354   cutoff, AutoCellLabeler still assigned a label to 90.6% of all ROIs with high-confidence human
355   labels, so we elected to delete the low-confidence (<75%) labels from the set of valid network
356   output labels (see Extended Data Fig. 2A for rationale for the 75% cutoff value).
357
358   We also examined model performance on data where humans had either low confidence or did
359   not assign a neuron label. In these cases, it was harder to estimate the ground truth. Overall,
360   model confidence was much lower for neurons that humans labeled with low confidence (87.3%)
361   or did not assign a label (81.3%). The concurrence of AutoCellLabeler relative to low-
362   confidence human labels was also lower (84.1%; we note that this is not truly a measure of
363   accuracy since these 'ground-truth' labels had low confidence). Indeed, overall the network's
364   concurrence versus human labels scaled with the confidence of the human label (Fig. 3G).
365
366   We carefully examined the subset of ROIs where the network had high confidence (>75%), but
367   humans had either low-confidence or entered no label at all. This was quite a large set of ROIs:
368   AutoCellLabeler identified significantly more high confidence neurons (119/animal) than the
369   original human labelers (83/animal), and this could conceivably reflect a highly accurate pool of
370   labels exceeding human performance. To determine whether this was the case, we obtained new
371   human labels (by different human labelers) for a random subset of these neurons. Whereas some
372   human labels remained low-confidence, others were now labeled with high confidence (20.9% of
373   this group of ROIs). The new human labelers also labeled neurons that were originally labeled
374   with high confidence so that we could compare the network's performance on relabeled data
375   where the original data was unlabeled, low confidence, or high confidence. AutoCellLabeler's
376   performance on all three groups was similar (88%, 86.1%, and 92.1%, respectively), which was
377   comparable to the accuracy of humans relabeling data relative to the original high-confidence
378   labels (92.3%). The slightly lower accuracy on these re-labeled data is likely due to the human
379   labeling of the original training, validation, and testing data being highly vetted and thoroughly
380   double-checked, whereas the re-labeling that we performed just for this analysis was done in a
381   single pass. Overall, these analyses indicate that the high-confidence network labels (119/animal)
382   have similar accuracy regardless of whether the original data had been labeled by humans as un-
383   labelable, low confidence, or high confidence. This indicates that AutoCellLabeler can
384   confidently label more neurons per dataset than individual human labelers.
385
386   We also split out model performance by cell type. This largely revealed similar trends. Model
387   labeling accuracy and confidence were variable among the neuron types, with highest accuracy
388   and confidence for the cell types where there were higher confidence human labels and a higher
389   frequency of human labels (Fig. 3K). For the labels where there were high confidence network
390   and human labels, we generated a confusion matrix to see if AutoCellLabeler's mistakes had
391   recurring trends (Extended Data Fig. 2B). While mistakes of this type were very rare, we
392   observed that the ones that occurred could mostly be categorized as either mislabeling a gut
393   granule as the neuron RMG, or mislabeling the dorsal/ventral categorization of the neurons IL1
394   and IL2 (e.g.: mislabeling IL2D as IL2). Together, these categories accounted for 50% of all
395   AutoCellLabeler's mistakes. We also observed that across cell types, AutoCellLabeler's
396   confidence was highly correlated with human confidence (Extended Data Fig. 2C), suggesting
397   that the main limitations of model accuracy are due to human labeling accuracy and confidence.
398

9

399 To provide better insights into which network features were critical for its performance, we
400 trained additional networks lacking some of AutoCellLabeler's key features. To evaluate these
401 networks, we considered both the number of high confidence labels assigned by AutoCellLabeler
402 and the accuracy of those labels measured against high-confidence human labels. Surprisingly, a
403 network that was trained with only standard image augmentations (i.e. lacking the custom
404 augmentation to bend the images in a manner that approximates a worm head bend) had similar
405 performance (Fig. 3I). However, a network that was trained without a pixel-weighting scheme
406 (i.e. where all pixels were weighted equally) provided far fewer high-confidence labels. This
407 suggests that devising strategies for pixel weighting is critical for model performance, though our
408 custom augmentation was not important. Interestingly, all trained networks had similar accuracy
409 (Fig. 3J) on their high-confidence labels, suggesting that the network architecture in all cases is
410 able to accurately assess its confidence.
411

412 **Automated annotation of *C. elegans* neurons from fewer fluorescent labels and in different**
413 **strains**

414 We examined whether the full group of fluorophores were critical for AutoCellLabeler
415 performance. This is a relevant question because (i) it is laborious to make, inject, and annotate a
416 large number of plasmids driving fluorophore expression, and (ii) the large number of plasmids
417 in the NeuroPAL strain has been noted to adversely impact the animals' growth and
418 behavior[35,40,48]. To test whether fewer fluorescent labels could still facilitate automatic labeling,
419 we trained four additional networks: one that only received the pan-neuronal tagRFP image as
420 input, and three that received pan-neuronal tagRFP plus a single other fluorescent channel
421 (CyOFP, tag-mBFP2, or mNeptune). As we still had the ground-truth labels based on humans
422 viewing the full set of fluorophores, the supervised labels were identical to those supplied to the
423 full network.
424

425 We evaluated the performance of these models by quantifying the number of high-confidence
426 labels that each network provided in each testing dataset (Fig. 4A) and the accuracy of these
427 labels measured against high-confidence human labels (Fig. 4B). We found that all four
428 networks had attenuated performance relative to the full AutoCellLabeler network, which was
429 almost entirely explainable by these networks having lower confidence in their labels, since
430 network accuracy was always consistent with its confidence (Extended Data Fig. 3A). This
431 means that labels from any version of the network can be treated equivalently, where the
432 confidence of a given label can be taken as an indication of its accuracy. Additionally, of the four
433 attenuated networks, the tagRFP+CyOFP network performance (107 neurons per animal labeled
434 at 97.4% accuracy) was quite close to the full network in its performance. Given that there
435 are >20 mTagBFP2 and mNeptune plasmids in the full NeuroPAL strain, these results raise the
436 possibility that a smaller set of carefully chosen plasmids could permit training of a network with
437 equal performance to the full network that we trained here.
438

439 We did not expect the tagRFP-only network to perform well, since the task of labeling tagRFP-
440 only images is nearly impossible for humans. Surprisingly, this network still exhibited relatively
441 high performance, with an average of 94 high-confidence neurons per animal and 94.8%
442 accuracy on those neurons. On most neuron classes, it behaved nearly as well as the full network,
443 though there are 10-20 neuron classes that it is much worse at labeling, such as ASG, IL1, and
444 RMG (Fig. 4C). Since this network only requires the red channel fluorescence, it could in theory

10

445  be used directly on freely-moving data, which has only GCaMP and tagRFP channel data.
446  Potentially, network performance could be increased by evaluating it on many different
447  timepoints from the freely-moving data to allow it to see the worm in many different postures.
448  Since the tagRFP-only network was trained only on high-SNR images collected from
449  immobilized animals, we first checked that the network was able to generalize outside its
450  training distribution to single images with lower SNR (example images in Extended Data Fig.
451  3B). It was able to label 79 high-confidence neurons per animal at 95.2% accuracy on the lower
452  SNR images (Fig. 4A-B, right). We then investigated whether allowing the network to access
453  different postures of the same animal improved its accuracy. Specifically, we evaluated the
454  tagRFP-only network on 100 randomly-selected timepoints in the freely-moving data of each
455  animal (example images in Extended Data Fig. 3B). We then related these 100 network labels to
456  the human labels, which could be easily determined, since ANTSUN registers free-moving
457  images back to the immobilized NeuroPAL images that had been labeled by humans. We
458  averaged the 100 network labels to obtain the most likely network label for each neuron, as well
459  as the average confidence for that label. To properly compare network versions, we determined
460  how many neurons could be labeled at any given target labeling accuracy – for example, how
461  many neurons the network can label and still achieve 95% accuracy (Fig. 4D; changing the
462  threshold network confidence value to include a given label allowed us to determine these full
463  curves). This analysis revealed that averaging network labels across the 100 timepoints improved
464  network performance, though only modestly. These results suggest that single color labels can be
465  used to train networks to a high level of performance, but additional fluorescence channels
466  further improve performance.
467
468  The strong performance of the tagRFP-only network on out-of-domain lower SNR images
469  suggest an impressive ability of the AutoCellLabeler network to generalize across different
470  modalities of data. This raised the possibility that it may be possible to use this network
471  architecture to build a foundation model of *C. elegans* neuron annotation that works across
472  strains and imaging conditions. As a first step to explore this, we investigated to what extent the
473  tagRFP-only network could generalize to other strains of *C. elegans* besides the NeuroPAL
474  strain. We used our previously-described SWF415 strain, which contains pan-neuronal NLS-
475  GCaMP7F, pan-neuronal NLS-mNeptune2.5, and sparse tagRFP expression[35]. Notably, the pan-
476  neuronal promoter utilized in this strain for NLS-mNeptune expression (P*rimb-1*) is distinct from
477  the pan-neuronal promoter that drives NLS-tagRFP expression in NeuroPAL (a synthetic
478  promoter). Since humans do not know how to label neurons in SWF415, we did a more limited
479  analysis by analyzing network labels for a subset of neurons that have highly reliable activity
480  dynamics with respect to behavior (AVA, AVE, RIM, and AIB encode reverse locomotion; RIB,
481  AVB, RID, and RME encode forward locomotion; SMDD encodes dorsal head curvature; and
482  SMDV and RIV encode ventral head curvature)[35,49–56]. Specifically, we asked whether neurons
483  labeled with high confidence by the network had the behavior encoding properties typical of the
484  neuron, assessed via analysis of the GCaMP traces from that neuron. Our previously-described
485  CePNEM model[35] was used to determine whether each labeled neuron encoded forward/reverse
486  locomotion or dorsal/ventral head curvature. The network provided high-confidence labels for an
487  average of 7.4/21 of these neurons per animal, and the encoding properties of these neurons
488  matched expectations 68% of the time (randomly labeled neurons had a match of 19%).
489  However, it was possible for the network to (i) incorrectly label a neuron as another neuron that
490  happened to have the same encoding; or (ii) correctly label a neuron that CePNEM lacked

11

491    statistical power to declare an encoding for. We accounted for these effects via simulations (see
492    Methods), which estimated that the actual labeling accuracy of the network on SWF415 was
493    69% (Fig. 4E). This is substantially lower than this network's accuracy on similar images from
494    the NeuroPAL strain (i.e. the strain used to train the network), where an average of 12.5 of these
495    neurons per animal were labeled with 97.1% accuracy. Nevertheless, this analysis indicates that
496    AutoCellLabeler has a reasonable ability to generalize to strains with different genetic drivers
497    and fluorophores, suggesting that in the future it may be worthwhile to pursue building a
498    foundation model that labels *C. elegans* neurons across many strains.
499
500    **A neural network (CellDiscoveryNet) that facilitates unsupervised discovery of >100 cell**
501    **types by aligning data across animals**
502
503    Annotation of cell types via supervised learning is fundamentally limited by prior knowledge and
504    humans' ability to label multi-spectral imaging data. In principle, unsupervised approaches that
505    can automatically identify stereotyped cell types would be preferable. Thus, we next sought to
506    train a neural network to perform unsupervised discovery of the cell types of *C. elegans* nervous
507    system (Fig. 5A). If successful, these approaches could be useful for labeling of mutant
508    genotypes, new permutations of NeuroPAL, or even related species. In addition, such an
509    approach would be useful in more complex animals that do not yet have complete catalogs of
510    cell types.
511
512    To facilitate unsupervised cell type discovery, we trained a network to register different animals'
513    multi-spectral NeuroPAL imaging data to one another. Successful alignment of cells across all
514    recorded animals would amount to unsupervised cell type annotation, since the cells that align
515    across animals would be the same cell type identified in different animals. The architecture of
516    this network was similar to BrainAlignNet, but the training data here consisted of pairs of 4-color
517    NeuroPAL images from two different animals and the network was tasked with aligning all four
518    fluorescent channels (Fig. 5B). No cell type positions (i.e. centroids) or neuronal identities were
519    provided to the network during training. Regularization and augmentation were similar to that of
520    BrainAlignNet (see Methods). Training and validation data were comprised of 91 animals'
521    datasets, which gave rise to 3285 unique pairs for alignment; 11 animals were withheld for
522    testing (the same test set as for AutoCellLabeler). The validation loss plateaued after 600 epochs
523    (Fig. 5C) and we characterized the network that had the minimum validation loss (at epoch 596).
524    In the analyses below, we characterize performance on training data and withheld testing data,
525    describing any differences. We note that, in contrast to the networks described above, high
526    performance on training data is still useful in this case, since the only criterion for success in
527    unsupervised learning is successful alignment (i.e. even if all data need to be used for training to
528    do so). Strong performance on testing data is still more desirable though, since it is less efficient
529    to train different networks over and over as new data are incorporated into the full dataset.
530
531    We first characterized the ability of this Unsupervised **Cell Discovery Net**work
532    (**CellDiscoveryNet**) to align images across different animals. Image alignment was reasonably
533    high for all four fluorescent NeuroPAL channels with a median NCC of 0.80 overall (Fig. 5D).
534    Alignment accuracy was nearly equivalent in training and testing data (Fig. 5D). We also
535    examined how well the centroid positions of defined cell types were aligned, utilizing our prior
536    knowledge of neurons' locations – i.e. the human labels (Fig. 5E). We computed this metric only

12

537     on cell types that were identified with high confidence in both of the images of a given
538     registration problem. The median centroid distance was 7.2 pixels, with similar performance on
539     training and testing data. This was initially rather disappointing, as it suggested that the majority
540     of neurons were not being placed at their correct locations. However, we observed two important
541     properties of the centroid alignments. First, the distribution of centroid distances was bimodal –
542     the $20^{th}$ percentile centroid distance was only 1.4 pixels, which corresponds to a correct neuron
543     alignment. Second, the median centroid distance decreased to 3.3 for registration problems with
544     high (> $90^{th}$ percentile = 0.85) NCC scores on the images. Together, these observations suggest
545     that CellDiscoveryNet correctly aligns neurons some of the time.
546
547     We next sought to differentiate the neuron alignments where CellDiscoveryNet was correct from
548     those where it was incorrect. Effectively, we wanted to treat CellDiscoveryNet as a hypothesis
549     generator for which neurons might be the same, and then algorithmically separate good
550     hypotheses from bad ones, stitching together the accurate hypotheses into a full set of neuron
551     alignments. To accomplish this, we adapted our ANTSUN pipeline (described in Fig. 2) to use
552     CellDiscoveryNet instead of BrainAlignNet. This modified ANTSUN 2U (**U**nsupervised) takes
553     as input multi-spectral data from many animals instead of monochrome images from different
554     time points of the same animal. This approach then allows us to effectively cluster neurons that
555     might be the same neuron found in different animals. Thus, we ran CellDiscoveryNet on pairs of
556     images and used the resulting DDFs to align the corresponding segmented neuron ROIs. We then
557     constructed a N-by-N matrix where N is all segmented neurons detected across all of the
558     NeuroPAL images (i.e. all neurons in all animals). Entries in the matrix are zero if the two
559     neurons were in an image pair that was never registered or if the two neurons did not overlap at
560     all in the registered image pair. Otherwise, a heuristic value indicating the likelihood that the
561     neurons are the same was entered into the matrix. This heuristic included the same information
562     as in ANTSUN 2.0 (described above), such as registration quality and ROI position similarity.
563     The only difference was that the heuristic for tagRFP brightness similarity was replaced with a
564     heuristic for 4-channel color similarity (see Methods). Custom hierarchical clustering of the rows
565     of this matrix then identified groups of ROIs hypothesized to be the same cell type identified in
566     different animals.
567
568     To determine the performance of this unsupervised cell type discovery approach, we quantified
569     both the number of cell types that were discovered (i.e. number of clusters) and the accuracy of
570     cell type labeling within each cluster. Here, accuracy was computed by first determining the
571     most frequent neuron label for each cell type, based on the human labels. We then determined
572     the number of correct versus incorrect detections of this cell type for all cells that fell within the
573     cluster, where a correct detection was defined to be when the human label for that cell matched
574     the most frequent label for that cell's cluster. The number of cell types identified and the labeling
575     accuracy are directly related: more permissive clustering identifies more cell types, but at the
576     cost of lower accuracy. A full curve revealing this tradeoff is shown in Fig. 5F (the parameter $w_7$
577     controls the restrictiveness of the clustering; see Methods). Based on this curve, we selected the
578     clustering parameter $w_7 = 10^{-9}$ that identified 125 cell types with 93% labeling accuracy. Not
579     every cell type is detected in every animal. On the testing data, the CellDiscoveryNet-powered
580     ANTSUN 2U roughly matched human-level performance in terms of accuracy and number of
581     neurons labeled per animal (Fig. 5G-H). However, it fell slightly short of AutoCellLabeler (Fig.

582    5G-H). Overall, this analysis reveals that CellDiscoveryNet facilitates unsupervised cell type
583    discovery with a high level of performance, matching trained human labelers.
584
585    We examined whether the accuracy of cell identification was different across cell types or across
586    animals. Fig. 5I shows the accuracy of labeling for each cell type (see Extended Data Figure 4A
587    for per-animal accuracy). Indeed, results were mixed: some cell types had highly accurate
588    detections across animals (eg: OLQD and RME), whereas a smaller subset of cell types were
589    detected with lower accuracy (eg: AIZ and ASG), and yet other cell types were harder to assess
590    accuracy due to a smaller number of human labels (eg: AIM and I4). In addition, there were five
591    clusters which did not contain a sufficient number of human-labeled ROIs to be given a cell type
592    label (<3 cells in these clusters had matching human labels; these are labeled "NEW 1" through
593    "NEW 5"). To examine which neurons these might correspond to, we examined the high-
594    confidence AutoCellLabeler labels for ROIs in these clusters. This produced enough labels to
595    categorize four of these five clusters as SAAD, SMBD, VB02, and VB02. The repeated VB02
596    label is likely an indication of under-clustering (ie: the two VB02 clusters should have been
597    merged into the same cluster). The identity of the fifth cluster was unclear, as the ROIs in that
598    cluster were not well labeled by either humans or AutoCellLabeler.
599
600    Finally, we examined whether CellDiscoveryNet was able to label cells not detected via
601    AutoCellLabeler. Specifically, we determined the fraction of the cells detected by
602    CellDiscoveryNet that were labeled by AutoCellLabeler, which was 86%. The new unsupervised
603    detections (the remaining 14%) included: new labels for cells that were otherwise well-labeled
604    by AutoCellLabeler (e.g.: M3); the detection and labeling of several cell types that were
605    uncommonly labeled by AutoCellLabeler (e.g.: RMEV); and the previously-mentioned cell type
606    that could not be identified. This suggests that the unsupervised approach that we describe here is
607    able to provide cell annotations that were not possible via human labeling or AutoCellLabeler.
608
609
610    **DISCUSSION**
611
612    Aligning and annotating the cell types that make up complex tissues remains a key challenge in
613    computational image analysis. We trained a series of deep neural networks that allow for
614    automated non-rigid registration and neuron identification in the context of brain-wide calcium
615    imaging in freely-moving *C. elegans*. This provides an appealing test case for the development
616    of such tools. *C. elegans* movement creates major challenges with tissue deformation and the
617    animal has >100 defined cell types in its nervous system. We describe BrainAlignNet, which can
618    perform non-rigid registration of the neurons of the *C. elegans* head, allowing for 99.6%
619    accuracy in aligning individual neurons. We also describe AutoCellLabeler, which can
620    automatically label >100 neuronal cell types with 98% accuracy, exceeding the performance of
621    individual human labelers by aggregating their knowledge. Finally, CellDiscoveryNet aligns data
622    across animals to perform unsupervised discovery of stereotyped cell types, identifying >100 cell
623    types of the *C. elegans* nervous system from unlabeled data. These tools should be useful for a
624    wide range of applications in *C. elegans* and should be straightforward to generalize to analyses
625    of other complex tissues.
626

627 Our newly-described network for freely-moving worm registration on average aligns neurons
628 with single pixel-level accuracy. Incorporating the network into a full image processing pipeline
629 indicates that it allows us to link neurons across time with 99.6% accuracy. Training a network to
630 achieve this high performance highlighted a series of general challenges. For example, our
631 attempt to train the network in a fully unsupervised manner (i.e. to simply align two images with
632 no further information) failed. While the resulting networks aligned RFP images of testing data
633 nearly perfectly, it turned out that the image transformations underlying this registration reflected
634 a scrambling of pixels and that the network was not warping the images in the manner that the
635 animal actually bends. We note that it was only possible to detect this failure mode through
636 unique control datasets that we had available to us, namely a strain that also had GFP in a sparse
637 subset of neurons and prior knowledge of ROI locations in the images. A semi-supervised
638 training procedure that utilized information about ROI locations ultimately prevented this failure
639 mode. While this approach is quite feasible for our use case, other types of datasets may not have
640 additional features such as ROI centroids to serve as supervised labels. It is possible that image
641 augmentation[33] might be able to assist in such cases.

643 Another limitation was that even with the semi-supervised approach, we were only able to train
644 networks to register images from reasonably well initialized conditions. Specifically, we
645 provided Euler-registered image pairs that were selected to have moderately similar head
646 curvature (though we note that these examples still had fairly dramatic non-rigid deformations;
647 see Figure 1). Solving this problem was sufficient to fully align neurons from freely-moving *C.*
648 *elegans* brain-wide calcium imaging, since clustering could effectively be used to link identity
649 across all timepoints even if our image registration only aligned a subset of the image pairs. Our
650 attempts to train a network to register all timepoints to one another was unsuccessful, though a
651 variety of approaches could conceivably improve upon this moving forward.

653 The AutoCellLabeler network that we describe here now automates a task that previously
654 required several hours of manual labeling per dataset. It achieves 98% accuracy in cell
655 identification and labels more neurons per dataset than individual human labelers. This
656 performance required a pixel weighting scheme where the network was trained to be especially
657 sensitive to high-confidence labels of neurons that were not ubiquitously labeled by all human
658 labelers. In other words, the network could aggregate knowledge across human labelers and
659 example animals to achieve high performance. While the high performance of AutoCellLabeler
660 is extremely useful from a practical point of view, we note that AutoCellLabeler still cannot label
661 all ROIs in a given image, which would be the highest level of desirable performance. Our
662 analyses suggest that it is currently bounded by human labeling of training data, which in turn is
663 bounded by our NeuroPAL image quality and the ambiguity of labeling certain neurons in the
664 NeuroPAL strain.

666 While improvements in human labeling could improve performance of the network, this analysis
667 also highlighted that it would be highly desirable to perform fully unsupervised cell labeling,
668 where the cell types could be inferred and labeled in multispectral images even without any
669 human labeling. To accomplish this, we developed CellDiscoveryNet, which aligns NeuroPAL
670 images across animals. Together with a custom clustering approach, this enabled us to identify
671 125 neuron classes, labeling them with 93% accuracy in a completely unsupervised manner. This
672 approach could be very useful within the *C. elegans* system, since it is extremely time

15

673 consuming to perform human labeling and it is conceivable that the NeuroPAL labels may
674 change in different genotypes or if the NeuroPAL transgene is modified. Beyond *C. elegans*,
675 these unsupervised approaches should be useful, since the vast majority of tissues in larger
676 animals do not yet have a full catalog of cell types and, therefore, would greatly benefit from
677 unsupervised discovery. In this spirit, other recent studies have started to develop approaches for
678 unsupervised labeling of imaging data[11,57,58], though these efforts were not aimed at identifying
679 the full set of cellular subtypes (>100) in individual images, which was the chief objective of
680 CellDiscoveryNet.
681
682 We also note the importance of our post-processing clustering approach to improving the
683 robustness of neural networks in solving these image registration problems. BrainAlignNet and
684 especially CellDiscoveryNet will sometimes generate incorrect solutions to individual neuron
685 mappings between a single pair of images. Relying solely on the network output for an
686 individual pair of images would thus be prone to inaccuracy. However, by treating the networks
687 as hypothesis generators across many images and using clustering to consider more likely
688 hypotheses first, we can generate highly accurate linkages across all images. We speculate that
689 this strategy might generalize across disciplines to many problems where it is possible to use
690 deep neural networks to generate large numbers of hypotheses whose likelihoods can be
691 heuristically evaluated.
692
693 We trained alternative versions of AutoCellLabeler with subsets of the spectral information in
694 NeuroPAL, which provides some insights into the possibility of performing high-accuracy neural
695 identification in strains with less fluorophores than NeuroPAL. On the one hand, all networks
696 that were trained with fewer than the full set of 4 fluorescent channels exhibited poorer
697 performance. However, it is notable that the network trained with only pan-neuronal RFP still
698 achieved 95% accuracy in labeling 94 neurons per image. It is important to note that this is the
699 performance of a network that was only trained to evaluate a single static image. It is
700 conceivable that there could be an improvement in performance if the network were trained on
701 pan-neuronal RFP images from all freely-moving timepoints, since this might allow the network
702 to infer identity based on the full range of movement and deformations that a given neuron
703 exhibits, which is quite stereotyped[36–40]. The fact that AutoCellLabeler exhibited surprisingly
704 good out-of-domain performance on images with different SNRs and on different strains
705 suggests that it may also be possible to improve performance across strains by building a
706 foundation model similar to AutoCellLabeler that has been specifically engineered to solve the
707 general task of labeling the cell types of the *C. elegans* brain in a wide range of images and
708 strains (data are starting to be aggregated into data repositories[42]). Future efforts should be able
709 to build upon the tools described here to lead to these types of improvements.
710
711 It should also be possible to combine the tools that we describe here to great effect. For example,
712 the unsupervised cell labels from CellDiscoveryNet could be used to train AutoCellLabeler in
713 order to obtain more unsupervised labels at higher accuracy. Moreover, this process could
714 potentially be multiplexed to achieve better cell annotation from TagRFP only. For example,
715 multiple multi-spectral transgenic lines like NeuroPAL could be subject to CellDiscoveryNet
716 labeling and fed into parallel AutoCellLabeler variants that only use the red channel for
717 prediction. These networks could then be combined to potentially achieve high-performance cell
718 identification from TagRFP only, utilizing knowledge gained from multiple multi-spectral lines.

16

719

720 These approaches for registering and annotating cells in dense tissues should be straightforward
721 to generalize to other species. For example, variants of BrainAlignNet could be trained to
722 facilitate alignment of tissue sections or to register imaging data onto a common anatomical
723 reference atlas. Our results suggest that training these networks on subsets of data with labeled
724 feature points, such as cell centroids (i.e. the semi-supervised approach we use here), will
725 facilitate more accurate solutions that, after training, can still be applied to datasets without any
726 labeled feature points. In addition, variants of AutoCellLabeler could be trained on any multi-
727 color cellular imaging data with manual labels. A pixel-wise labeling approach, together with
728 appropriate pixel weighting during training, should be generally useful to build models for
729 automatic cell labeling in a range of different tissues and animals. Finally, models similar to
730 CellDiscoveryNet could be broadly useful to identify previously uncharacterized cell types in
731 many tissues. It is conceivable that hybrid or iterative versions of AutoCellLabeler and
732 CellDiscoveryNet could lead to even higher performance cell type discovery and labeling.

733

740

741 **AUTHOR CONTRIBUTIONS**
742 Conceptualization, A.A.A., J.K., S.W.F. Methodology, A.A.A., J.K., A.KY.L. Software, A.A.A.,
743 J.K., A.KY.L.  Formal analysis, A.A.A., A.KY.L. Investigation, A.A.A., J.K., A.KY.L.. T.S.K.,
744 S.B., E.B., F.K.W., D.K. Writing – Original Draft, A.A.A. and S.W.F. Writing – Review &
745 Editing, A.A.A., J.K., A.KY.L., and S.W.F. Funding Acquisition, S.W.F.

746

747 **DECLARATION OF INTERESTS**
748 The authors have no competing interests to declare
749
750
751 **MATERIALS AND METHODS**
752
753 ***C. elegans* Strains and Genetics**
754
755 All data were collected from one-day old adult hermaphrodite *C. elegans* animals raised at 22C
756 on standard nematode growth medium (NGM) plates.
757
758 For the GCaMP-expressing animals without NeuroPAL, two transgenes were present: (1)
759 *flvIs17*: *tag-168::NLS-GCaMP7F + NLS-tagRFPt* expressed under a small set of cell-specific
760 promoters: *gcy-28.d, ceh-36, inx-1, mod-1, tph-1(short), gcy-5, gcy-7*; and (2) *flvIs18*: *tag-*
761 *168::NLS-mNeptune2.5*. This resulting strain SWF415, has been previously characterized[35].
762

For the GCaMP-expressing animals with NeuroPAL, two transgenes were present in the strain: (1) *flvIs17*: described above; and (2) *otIs670*: low-brightness NeuroPAL. This resulting strain, named SWF702, has been previously characterized[35].

The animals with *eat-4::NLS-GFP* and *tag-168::NLS-GFP* were also previously described[35]. As is described in the strain list, *tag-168::NLS-mNeptune2.5* was also co-injected with each of these plasmids to generate the two strains: SWF360 (*eat-4::NLS-GFP; tag-168::NLS-mNeptune2.5*) and SWF467 (*tag-168::NLS-GFP; tag-168::NLS-mNeptune2.5*).

We provide here a list of these four strains:

**SWF415** *flvIs17[tag-168::NLS-GCaMP7F, gcy-28.d::NLS-tag-RFPt, ceh-36:NLS-tag-RFPt, inx-1::tag-RFPt, mod-1::tag-RFPt, tph-1(short)::NLS-tag-RFPt, gcy-5::NLS-tag-RFPt, gcy-7::NLS-tag-RFPt]; flvIs18[tag-168::NLS-mNeptune2.5]; lite-1(ce314); gur-3(ok2245)*

**SWF702** *flvIs17; otIs670 [low-brightness NeuroPAL]; lite-1(ce314); gur-3(ok2245)*

**SWF360** *flvEx450[eat-4::NLS-GFP, tag-168::NLS-mNeptune2.5]; lite-1(ce314); gur-3(ok2245)*

**SWF467** *flvEx451[tag-168::NLS-GFP, tag-168::NLS-mNeptune2.5]; lite-1(ce314); gur-3(ok2245)*

**<u>Microscope and Recording Conditions</u>**

Data used to train and evaluate the models include previously-published datasets[35,46,59] and newly-collected data. These animals were recorded under similar recording conditions to those described in our previous study[35]. There were two types of datasets collected, relevant to this study: freely-moving GCaMP/TagRFP data, and immobilized NeuroPAL data.

Briefly, all neural data (free-moving and NeuroPAL) were acquired on a dual light-path microscope that was previously described[35]. The light path used to image GCaMP, mNeptune, and the fluorophores in NeuroPAL at single cell resolution is an Andor spinning disk confocal system with Nikon ECLIPSE Ti microscope. Light supplied from a 150 mW 488 nm laser, 50 mW 560 nm laser, 100 mW 405 nm laser, or 140 mW 637 nm laser passes through a 5000 rpm Yokogawa CSU-X1 spinning disk unit with a Borealis upgrade (with a dual-camera configuration). A 40x water immersion objective (CFI APO LWD 40X WI 1.15 NA LAMBDA S, Nikon) with an objective piezo (P-726 PIFOC, Physik Instrumente (PI)) was used to image the volume of the worm's head (a Newport NP0140SG objective piezo was used in a subset of the recordings). A custom quad dichroic mirror directed light emitted from the specimen to two separate sCMOS cameras (Zyla 4.2 PLUS sCMOS, Andor), which had in-line emission filters (525/50 for GCaMP/GFP, and 570 longpass for tagRFP/mNeptune in freely-moving recordings; NeuroPAL filters described below). Data was collected at a volume rate of 1.7 Hz (1.4 Hz for the datasets acquired with the Newport piezo).

For recordings, L4 worms were picked 18-22 hours before the imaging experiment to a new NGM agar plate seeded with OP50 to ensure that we recorded one day-old adult animals. Animals were recorded a thin, flat NGM agar pad (2.5cm x 1.8cm x 0.8mm). On the 4 corners of

18

809 the agar pad, we placed a single layer of microbeads with a diameter of 80um to alleviate the
810 pressure of the coverslip (#1.5) on the worm. Animals were transferred to the agar pad in a drop
811 of M9, after which the coverslip was added.
812      For NeuroPAL data collection, animals were immobilized via cooling, after which multi-
813 spectral information was captured. For cooling, the slide was mounted with a thermoelectric
814 cooling element attached to it, set to cool the agar temperature to 1 °C. A closed-loop
815 temperature controller (TEC200C, Thorlabs) with a micro-thermistor (SC30F103A, Amphenol)
816 embedded in the agar kept the agar temperature at the 1 °C set point. Once the temperature
817 reached the set point, we waited 5 minutes for the worm to be fully immobilized before imaging.
818      We obtained a series of images from each recorded animal while the animal was
819 immobilized (this has been previously described[35]):
820      (1-3) Spectrally isolated images of mTagBFP2, CyOFP1, and mNeptune2.5. We excited
821 CyOFP1 using the 488nm laser at 32% intensity under a 585/40 bandpass filter. mNeptune2.5
822 was recorded next using a 637nm laser at 48% intensity under a 655LP-TRF filter, in order to not
823 contaminate this recording with TagRFP-T emission. Finally, mTagBFP2 was isolated using a
824 405nm laser at 27% intensity under a 447/60 bandpass filter.
825      (4) An image with TagRFP-T, CyOFP1, and mNeptune2.5 (all of the "red" markers) in
826 one channel, and gCaMP7f in the other channel. As described in our previous study, this image
827 was used for neuronal segmentation and registration to both the freely moving recording and
828 individually isolated marker images. We excited TagRFP-T and mNeptune2.5 via 561nm laser at
829 15% intensity and CyOFP1 and gCaMP6f via 488nm laser at 17% intensity. TagRFP-T,
830 mNeptune2.5, and CyOFP1 were imaged with a 570LP filter and gCaMP6f was isolated using a
831 525/50 bandpass filter.
832      All isolated images were recorded for 60 timepoints. We increased the signal to noise
833 ratio for each of the images by first registering all timepoints within a recording to one another
834 and then averaging the transformed images. For manual labeling of these datasets, we created a
835 composite, 3-dimensional RGB image by setting the mTagBFP2 image to blue, CyOFP1 image
836 to green, and mNeptune2.5 image to red as done by Yemini et al. (2021) and manually adjusting
837 the intensity of each channel to optimally match their manual.
838
839 **Availability of Code**
840
841 All code is freely and publicly available (use main/master branches unless otherwise specified):
842
843 • BrainAlignNet: https://github.com/flavell-lab/BrainAlignNet and
844   https://github.com/flavell-lab/DeepReg (main branch)

845 • GPU-accelerated Euler registration: https://github.com/flavell-lab/euler_gpu

846 • ANTSUN 2.0: https://github.com/flavell-lab/ANTSUN (branch v2.1.0); see also
847   https://github.com/flavell-lab/flv-c-setup and https://github.com/flavell-
848   lab/FlavellPkg.jl/blob/master/src/ANTSUN.jl for auxiliary package installation.

849 • AutoCellLabeler: https://github.com/flavell-lab/pytorch-3dunet and
850   https://github.com/flavell-lab/AutoCellLabeler

851 • CellDiscoveryNet: https://github.com/flavell-lab/DeepReg (multicolor branch)

852 • ANTSUN 2U: https://github.com/flavell-lab/ANTSUN-Unsupervised

19

853
854 **BrainAlignNet**
855
856 **Network architecture**
857 BrainAlignNet's architecture is derived from the DeepReg software package, which uses a
858 variation of a 3-D U-Net architecture termed a LocalNet[44,45]. BrainAlignNet first has a
859 concatenation layer that concatenates the moving and fixed images together along a new, channel
860 dimension. The resulting $284 \times 120 \times 64 \times 2$ image is then passed as input to the LocalNet,
861 which outputs a $284 \times 120 \times 64 \times 3$ dense displacement field (DDF). The DDF defines a
862 coordinate transformation from fixed image coordinates to moving image coordinates, relative to
863 the fixed image coordinate system. So, for instance, if $DDF[x, y, z] = (\Delta x, \Delta y, \Delta z)$, it means that
864 the coordinates $(x, y, z)$ in the fixed image are mapped to the coordinates $(x + \Delta x, y + \Delta y, z +$
865 $\Delta z)$ in the moving image. The network has a final warping layer that applies the DDF to
866 transform the moving image into a predicted fixed image whose pixel at location $(x, y, z)$
867 contains the moving image pixel at location $(x, y, z) + DDF[x, y, z]$. It also has another final
868 warping layer that transforms the fixed image centroids $(x, y, z)$ into predicted moving image
869 centroids $(x, y, z) + DDF[x, y, z]$. The network's loss function causes it to seek to minimize the
870 difference between its predictions and the corresponding input data.
871
872 The LocalNet is at its core a 3-D U-Net with an additional output layer that receives inputs from
873 multiple output levels. In more detail, it has 3 input levels and 3 output levels, with $16 \cdot 2^i$
874 feature channels at the $i$th level for $i \in \{0,1,2\}$. It contains an encoder block mapping the input to
875 level 0, followed by two more encoder blocks mapping input level $i$ to level $i + 1$ for $i \in \{0,1\}$.
876 Each of these three encoder blocks contains a convolutional block, a residual convolutional
877 block, and a $2 \times 2 \times 2$ max-pool layer. The convolutional block consists of a 3-D convolutional
878 layer with kernel size 3 that doubles the number of feature channels, followed by a batch
879 normalization layer, followed by a ReLU activation function. The residual convolutional block
880 consists of two convolutional blocks in sequence, except that the input (to the residual
881 convolutional block) is added to the output of the second convolutional block right before its
882 ReLU activation function. The bottom block comes after the encoder block at level 2, mapping
883 input level 2 to output level 2. It has the same architecture as a single convolutional block;
884 notably, it does not contain the max-pool layer.
885
886 There are three decoder blocks receiving inputs from the three encoder blocks described above.
887 The first two decoder blocks map output level $i + 1$ to output level $i$ for $i \in \{1,0\}$; the third one
888 maps output level 0 to the preliminary output with the same $(x, y, z)$ dimensions as the input.
889 Each decoding block consists of an upsampling block, a skip-connection layer, a convolutional
890 block, and a residual convolutional block. The upsampling block contains a transposed 3D
891 convolutional layer with kernel size 3 that halves the number of feature channels and an image
892 resizing layer (run independently on the upsampling block's input) using bilinear interpolation to
893 double each dimension of the image. The output of the resizing layer is then split into two equal
894 pieces along the channel axis and summed, and then added to the output of the transposed
895 convolutional layer. The skip-connection layer appends the output of the mirrored encoder block
896 $i$ (for the third decoder block, this corresponds the first encoder block) right before that encoder
897 block's max pool layer. The skip-connection layer appends this output to the channel dimension,
898 doubling its size. The convolutional and residual convolutional blocks are identical to those in

20

899 the encoding block, except that the convolutional block halves the number of input channels
900 instead of doubling it.
901
902 Finally, there is the output layer. It takes as input the output of the bottom block, as well as the
903 output of every decoder block. To each of these inputs, it applies a 3D convolutional layer that
904 outputs exactly 3 channels, followed by an upsampling layer that uses bilinear interpolation to
905 increase the dimensions to the size of the original input images. It then averages together all of
906 these images to compute the final $284 \times 120 \times 64 \times 3$ DDF.
907
### Preprocessing
909 To train and validate a registration network that aligns neurons across time series in freely-
910 moving *C. elegans*, we took several steps to prepare the calcium imaging datasets with images
911 and their corresponding centroids. The preprocessing procedure consisted of (i) selecting two
912 different time points from a single video (fixed and moving time points) at which to obtain RFP
913 images (all images given to the network are from the red channel, which contains the signal from
914 NLS-TagRFP) and neuron centroids; (ii) cropping all RFP images to a consistent size; (iii)
915 performing Euler registration (translation and rotation) to align neurons from the image at the
916 moving time point (moving image) to the image at the fixed time point (fixed image); (iv)
917 creating image centroids for the network, which consist of matched lists of centroid positions of
918 all the neurons in both the fixed and moving images.
919
### *(i) Selection of registration problems.*
921 We refer to the task of solving the transformation function that aligns neurons from the moving
922 image to the fixed image as a registration problem. We selected our registration problems based
923 on previously constructed[35] image registration graphs using ANTSUN 1.4. In these registration
924 graphs, the time points of a single calcium imaging recording served as vertices. An edge
925 between two time points indicates a registration problem that we will attempt to solve. Edges
926 were preferentially created between time points with higher worm posture similarities.
927   In ANTSUN 1.4, we selected approximately 13,000 pairs of time points (fixed and
928 moving) per video that had sufficiently high worm posture similarity. These registration
929 problems were solved by gradient descent using our old image processing pipeline, and
930 ANTSUN clustering yielded linked neuron ROIs across frames that were the basis of
931 constructing calcium traces[35]. To train BrainAlignNet here, we randomly sampled about 100
932 problems across a total of 57 animals, ultimately compiling 5,176 registration problems for
933 training (some registration problems were discarded during subsequent preprocessing steps). To
934 prepare the validation datasets, we sampled 1,466 problems across 22 animals. Testing data was
935 447 problems from 5 animals.
936
### *(ii) Cropping.*
938 The registration network requires all 3D image volumes in training, validation, and testing to be
939 of the same size. Therefore, a crucial step in preprocessing was to crop or pad the images along
940 the *x*, *y*, *z* dimensions to a consistent size of (284, 120, 64). Before reshaping the images, we first
941 subtracted the median pixel value from each image (both fixed and moving) and set the negative
942 pixels to zero. Then, we either cropped or padded with zeros around the centers of mass of these
943 images to make the *x* dimension 284, the *y* dimension 120, and the *z* dimension 64.
944

945     *(iii) Euler registration.*

946     Through experimentation with various settings of the network, we have found that it is difficult

947     for the network to learn large rotations and translations at the same time as smaller nonlinear

948     deformations. Euler registration is far more computationally tractable than nonlinear

949     deformation, so we solved Euler registration for the images before providing them to the

950     network. In Euler registration, we rotate or translate the moving images by a certain amount,

951     aiming to maximize their normalized cross-correlation (NCC) with the fixed image. The optimal

952     parameters of translation and rotation that resulted in the highest NCC were determined using a

953     brute-force, GPU-accelerated parameter grid search. To further accelerate the grid search, we

954     projected the fixed and moving images onto the $x$-$y$ plane using a maximum-intensity projection

955     along the $z$-axis. We also downsampled the fixed and moving images by a factor of 4 after the $z$

956     maximal projection. The best parameters identified for transforming the projected images were

957     then applied to each z-slice to transform the entire 3D image. This approach was feasible because

958     the vast majority of worm movement occurs along the $x$-$y$ axes.

959

960     *(iv) Creating image centroids.*

961     We obtained the neuronal ROI images for both the fixed and moving RFP images, designating

962     them as the fixed and moving ROI images respectively. The full sets of ROIs in each image were

963     obtained using ANTSUN 1.4's image segmentation and watershedding functions. ROI images

964     were them constructed as follows. Each pixel in an ROI image contains an index value: 0 for

965     background, or a positive integer for a neuron. All pixels belonging to a specific neuron have the

966     same index, and pixels belonging to any other neuron have a different index. Since the ROI

967     images are created independently at each time point, their neuronal indices are not *a priori*

968     consistent across time points. Therefore, we used previous runs of ANTSUN 1.4 to link the ROI

969     identities across time points, and generated new ROI images with consistent indices across time

970     points – for example, all pixels with value 6 in one time point correspond to the same neuron as

971     pixels with value 6 in any other time point. We deleted any ROIs with indices that were not

972     present in both the moving and fixed images.

973

974        We then cropped these ROI images to the same size and subjected them to Euler

975     transformations using the same parameters as their corresponding fixed and moving RFP images.

976     Next, we computed the centroids of each neuron index in the resulting moving and fixed ROI

977     images. The centroid was defined to be the mean $x$, $y$, and $z$ coordinates of all pixels of a given

978     ROI. We stored these centroids as two lists of equal length (typically, around 110). Note that

979     these lists are now the matched positions of neurons in the fixed and moving images.

980

981        Since the network expects image centroids to be of the same size, all neuronal centroids

982     in the fixed and moving images were padded and aggregated into arrays of shape (200, 3),

983     ensuring the same ordering of neurons. The extra entries that do not contain neurons are filled

984     with (-1, -1, -1) to make the total number of neurons equal to 200. We designate the neuronal

985     centroid positions in the fixed and moving ROI images as fixed and moving centroids,

986     respectively.

987

988     **Loss functions**

989     Our main custom modifications to the DeepReg network focus on the design of the loss function.

990     In particular, we implemented a new supervised centroid alignment loss component and new

991  regularization loss sub-components. Overall, the loss function consists of three major
992  components:

- **Image loss** $L_I$ captures the difference between the warped moving image and the ground-truth fixed image.
- **Centroid alignment loss** $L_C$ is a supervised portion of the loss function. Given pre-labeled centroids corresponding to ground-truth information about neuron positions in the fixed and moving images, this loss component captures the difference between the predicted moving centroids and the ground-truth moving centroids.
- **Regularization loss** $L_R$ captures the prior that the "simplest" DDF that achieves the desired transform outcome is the best. For example, it's implausible that a pair of neurons that start close together end up on opposite sides of the worm, so a DDF that generates such a transformation would have a high value of regularization loss.

1004  The total loss is then computed as $Loss = w_I L_i + w_C L_C + w_R L_R$. We set $w_I = 1$, $w_C = 0.1$, and
1005  $w_R = 1$.

***(i) Image loss.***
The image loss is the negative of the local squared zero-normalized cross-correlation (LNCC) between the fixed and warped moving RFP images. We designate the fixed image as $X_{true}$ and the warped moving image as $X_{pred}$. Define $E(X)$ as a function that computes the discrete expectation of image $X$ within a sliding cube of side length $n=16$:

$$E(X)[x,y,z] = \frac{1}{n^3} \sum_{i=x}^{x+n-1} \sum_{j=y}^{y+n-1} \sum_{k=z}^{z+n-1} X[i,j,k]$$

We then can compute the discrete sliding variance as

$$V(X) = E(X^2) - E(X)^2$$

The image loss (i.e., negative LNCC) is then defined as

$$L_I = -\text{LNCC} = -\frac{\left[E(X_{true} \circ X_{pred}) - E(X_{true}) \circ E(X_{pred})\right]^2}{V(X_{true}) \circ V(X_{pred}) + \epsilon}$$

***(ii) Centroid alignment loss.***
The centroid alignment loss is calculated as the negative of the sum of the Euclidean distances between the moving centroids and the network's predicted moving centroids, averaged across the number of centroids available. We designate the ground-truth and network predicted centroids as $N \times 3$ matrices $y_{true}$ and $y_{pred}$ respectively, where $N$ is the number of centroids, and the $i$th row of each matrix represents the coordinates of neuron $i$'s centroid. Centroid alignment loss in the overall loss function is then expressed as follows:

23

1033
$$L_C = \frac{1}{N} \sum_{i=0}^{N-1} \sqrt{\sum_{d=0,1,2} \left( y_{true}[i,d] - y_{pred}[i,d] \right)^2}$$

1034 *(iii) Regularization loss.*
1035 Our regularization loss function consists of four terms that seek to penalize DDFs that do not
1036 correspond to possible physical motion of the worm. Of these terms, gradient norm is unchanged
1037 from its previous implementation in the DeepReg package, while the other three components are
1038 our additions:

1039

1040 • **Gradient norm loss** $L_{Grad}$ penalizes transformations for being nonuniform.
1041 • **Difference norm loss** $L_{Diff}$ penalizes transformations for moving pixels too far.
1042 • **Axis difference norm loss** $L_{AxisDiff}$ penalizes transformations for moving pixels too far
1043 along the $z$-dimension, which is less plausible than movement along the $x$- and $y$-
1044 dimensions in our recordings.
1045 • **Nonrigid penalty loss** $L_{Nonrigid}$ penalizes transformations for being nonrigid (i.e., not
1046 translation and rotation). (Note that unlike the gradient norm loss, this loss function will
1047 not penalize DDFs that apply rigid-body rotations.)

1048

1049 We then set $L_R = 0.02\, L_{Grad} + 0.005 L_{Diff} + 0.001\, L_{AxisDiff} + 0.02\, L_{Nonrigid}$

1050


1051 *Gradient Norm.* The gradient norm computes the average gradient of the DDF by
1052 summing up the central finite difference of the DDF as the approximation of derivatives along
1053 the $x$, $y$, and $z$ axes. Specifically, we first approximate the partial derivatives for $m \in \{0, 1, 2\}$ as
1054 follows:

1055
$$\frac{\partial D_m}{\partial x} \approx \frac{1}{2} \left( D[2:X,\ 1:Y-1,\ 1:Z-1,\ m] - D[0:X-2,\ 1:Y-1,\ 1:Z-1,\ m] \right)$$

1056
$$\frac{\partial D_m}{\partial y} \approx \frac{1}{2} \left( D[1:X-1,\ 2:Y,\ 1:Z-1,\ m] - D[1:X-1,\ 0:Y-2,\ 1:Z-1,\ m] \right)$$

1057
$$\frac{\partial D_m}{\partial z} \approx \frac{1}{2} \left( D[1:X-1,\ 1:Y-1,\ 2:Z,\ m] - D[1:X-1,\ 1:Y-1,\ 0:Z-2,\ m] \right)$$


1058 These results are then stacked to obtain $\frac{\partial D}{\partial x}, \frac{\partial D}{\partial y}$, and $\frac{\partial D}{\partial z}$. The gradient norm is calculated as the
1059 squared sum of these derivatives, averaged across all elements:


1060
$$L_{Grad} = \frac{1}{3(X-2)(Y-2)(Z-2)} \sum_{i=0}^{X-3} \sum_{j=0}^{Y-3} \sum_{k=0}^{Z-3} \sum_{m=0}^{2} \left[ \left( \frac{\partial D}{\partial x} \right)^2 + \left( \frac{\partial D}{\partial y} \right)^2 + \left( \frac{\partial D}{\partial z} \right)^2 \right]_{i,j,k,m}$$


1061 *Difference Norm.* The difference norm computes the average squared displacement of a
1062 pixel under the DDF $D$:

1063
$$L_{Diff} = \frac{1}{3XYZ} \sum_{i=0}^{X-1} \sum_{j=0}^{Y-1} \sum_{k=0}^{Z-1} \sum_{m=0}^{2} (D[i,j,k,m])^2$$

1064 where $X, Y, Z$ are the sizes of the image along the $x$, $y$, and $z$ axes respectively.

24

1065

1066

1067        *Axis Difference Norm.* Axis difference norm of the DDF $D$ calculates the average squared
1068 displacement of a pixel along the $z$-axis:

$$D_z = D[:,:,:,2] \qquad L_{AxisDiff} = \frac{1}{XYZ} \sum_{i=0}^{X-1} \sum_{i=0}^{Y-1} \sum_{k=0}^{Z-1} (D_z[i,j,k])^2$$

1069

1070        *Nonrigid penalty.* This term penalizes nonrigid transformations of the neurons by
1071 utilizing the gradient information of the DDF. Unlike the approach used in computing the
1072 gradient norm, where global rotations would have nonzero gradient, here we are interested in
1073 penalizing specifically nonrigid transforms. We accomplish this by constructing a reference
1074 DDF, denoted as $D_{ref}$, which warps the entire image to the origin: $D_{ref}[x,y,z,:] =$
1075 $[-x,-y,-z]$. Then the difference DDF $D_{diff} = D - D_{ref}$ has the property that the magnitude of
1076 its gradient is rotation-invariant. We can then compute $\frac{\partial D_{diff}}{\partial x}$, $\frac{\partial D_{diff}}{\partial y}$, and $\frac{\partial D_{diff}}{\partial z}$ as for the
1077 gradient norm and define the gradient magnitude:

$$M = \left(\frac{\partial D_{diff}}{\partial x}\right)^2 + \left(\frac{\partial D_{diff}}{\partial y}\right)^2 + \left(\frac{\partial D_{diff}}{\partial z}\right)^2$$

1078

1079 Under any rigid-body transform, $M = 1$. Thus, the nonrigid penalty is calculated as

$$L_{Nonrigid} = \frac{1}{3(X-2)(Y-2)(Z-2)} \sum_{i=0}^{X-3} \sum_{j=0}^{Y-3} \sum_{k=0}^{Z-3} \sum_{m=0}^{2} \left| M + \frac{1}{M} - 2 \right|_{i,j,k,m}$$

1080

1081 In this way, rigid-body transforms will have 0 loss while any nonrigid transform will have a
1082 positive loss.

1083

1084 **Data augmentation**
1085 During training, input data was subject to augmentation. We used random affine transformations
1086 for augmentation. Each transformation was generated by perturbing the corner points of a cube
1087 by random amounts, and computing the affine transformation resulting in that perturbation. The
1088 same transformation was then applied to the moving image, fixed image, moving centroids, and
1089 fixed centroids.

1090

1091 **Optimizer**
1092 BrainAlignNet was trained using the Adam optimizer with a learning rate of $10^{-4}$.

1093

1094 **Configuration file**

1095 The full configuration file we used during network training is available at
1096 https://github.com/flavell-lab/BrainAlignNet/tree/main/configs

1097

1098 **<u>Automatic Neuron Tracking System for Unconstrained Nematodes (ANTSUN) 2.0</u>**

1099 We integrated BrainAlignNet into our previously-described ANTSUN pipeline[35,46] (also applied
1100 in[59]). Briefly, the pipeline: (i) performs some image pre-processing such as shear-correction and

25

1101    cropping; (ii) segments the images into neuron ROIs via a 3D U-Net; (iii) finds time points
1102    where the worm postures are similar; (iv) performs image registration to define a coordinate
1103    mapping between these time points; (v) applies that coordinate mapping to the ROIs; (vi)
1104    constructs an ROI similarity matrix storing how likely different ROIs are to correspond to the
1105    same neuron; (vii) clusters that matrix to extract neuron identity; (viii) maps the linked ROIs
1106    onto the GCaMP data to extract neural traces; and (ix) performs some postprocessing such as
1107    background-subtraction and bleach correction to extract neural traces.

1108    The differences in ANTSUN 2.0 compared with our previously-published version of this
1109    pipeline, ANTSUN 1.4, are that in ANTSUN 2.0 we use BrainAlignNet to perform image
1110    registration rather than the gradient descent-based elastix, and we modified the heuristic function
1111    used to construct the ROI similarity matrix. We only replaced the freely-moving registration with
1112    BrainAlignNet; the immobilized registrations, channel alignment registration, and freely-moving
1113    to immobilized registration are still performed with elastix. These remaining elastix-based
1114    registrations are much less computationally expensive, taking only about 2% of the total
1115    computation time of the original ANTSUN 1.4 pipeline. They will also likely be replaced with
1116    BrainAlignNet in a future release of ANTSUN, after further diagnostics and controls are run.

1117    The heuristic function used to compute the ROI similarity matrix was updated to add additional
1118    terms specific to BrainAlignNet, including regularization and an additional ROI displacement
1119    term that serves to implement our prior that ROIs which moved less far in the registration are
1120    more likely to be correctly registered. Letting $i$ and $j$ be two different ROIs in our recording at
1121    time points $t_i$ (moving) and $t_j$ (fixed), the full expression for the ROI similarity matrix is:

1122
$$M_{ij} = R_{t_i t_j} \frac{1}{1 + w_1 d_i} \, q_{t_i t_j}^{w_2} r_{ij}^{w_3} e^{-\left(w_4 a_{ij} + w_5 c_{ij} + w_6 n_{t_i t_j}\right)}$$

1123                                            Where:

1124        $R_{t_i t_j}$ is 1 if there exists a registration mapping $t_i$ to $t_j$, and 0 otherwise.

1125    $d_i$ is the displacement of the centroid of ROI $i$ under the DDF registration between $t_i$ and $t_j$.

1126    $q_{t_i t_j}$ is the registration quality, computed as the NCC of warped moving image $t_i$ vs fixed image
1127                                            $t_j$.

1128    $r_{ij}$ is the fractional overlap of warped moving ROI $i$ and fixed ROI $j$ (intersection / max size).

1129    $a_{ij}$ is the absolute difference in marker channel activity (i.e. tagRFP brightness) between ROIs $i$
1130            and $j$, normalized to mean activity at the corresponding timepoints $t_i$ and $t_j$.

1131        $c_{ij}$ is the distance between the centroids of warped moving ROI $i$ and fixed ROI $j$.

1132        $n_{t_i t_j}$ is the (unweighted) nonrigid penalty loss of the DDF registration from $t_i$ to $t_j$.

1133                    $w_i$ are weights controlling how important each variable is.

26

1134    Additionally, the matrix is forced to be symmetrical by setting $M_{ji} = M_{ij}$ whenever $M_{ji} = 0$ and
1135    $M_{ij} \neq 0$. It is also sparse since $R_{t_i t_j}$ and $r_{ij}$ are usually 0. Finally, there are two additional
1136    hyperparameters in the clustering algorithm, $w_7$ and $w_8$. $w_7$ controls the minimum height the
1137    clustering algorithm will reach (effectively, $w_7$ is a cap on how low $M_{ij}$ values can get, or how
1138    low the heuristic value can fall before determining that the ROIs are not the same neuron) and
1139    $w_8$ controls the acceptable collision fraction (a collision is defined by a cluster containing
1140    multiple ROIs from the same timepoint, which should not happen since each neuron should
1141    correspond to only one ROI at each time point).

1142    We determined the weights $w_i$ by performing a grid search through 2,912 different combinations
1143    of weights on three *eat-4*::NLS-GFP datasets. To evaluate the outcome of each combination, we
1144    computed the error rate (rate of incorrect neuron linkages) and number of detected neurons. The
1145    error rate was computed as previously described[35]: since the strain *eat-4*::NLS-GFP expresses
1146    GFP in some but not all neurons, we can quantify registration errors as instances where a GFP-
1147    positive neuron lacked GFP in a time point and vice versa, as these correspond to neuron
1148    mismatches. We then selected the combination of parameters that maximize the number of
1149    detected neurons while minimizing the error rate. One *eat-4*::NLS-GFP dataset (the one shown in
1150    Figure 2) was used as a withheld testing animal to determine this optimal set of parameters. The
1151    pan-neuronal GFP and pan-neuronal GCaMP animals were not included in this parameter search.

1152    The values of the parameters we used were:

1153    $$w_1 = 2$$

1154    $$w_2 = 25$$

1155    $$w_3 = 1$$

1156    $$w_4 = 3$$

1157    $$w_5 = 1$$

1158    $$w_6 = 1$$

1159    $$w_7 = 0.0001$$

1160    $$w_8 = 0.05$$

1161

## AutoCellLabeler

*Network Architecture*

1164    AutoCellLabeler uses a 3-D U-Net architecture[35,47], with input dimensions $4 \times 64 \times 120 \times 284$
1165    (fluorophore channel, $z$, $y$, $x$) and output dimensions $185 \times 64 \times 120 \times 284$ (label channel, $z$, $y$,
1166    $x$). The 3D U-Net has 4 input levels and 4 output levels, with $64 \cdot 2^i$ feature channels at the $i$th
1167    level for $i \in \{0,1,2,3\}$.

1168    There is an encoder block that maps an input image to the $0^{th}$ input level, followed by three
1169    additional encoder blocks that map input level $i$ to input level $i + 1$ for $i \in \{0,1,2\}$. Each encoder
1170    block consists of two convolutional blocks followed by a $2 \times 2 \times 2$ max pool layer, with the
1171    exception of the first encoder layer which does not have the max pool layer. The first
1172    convolutional block in each encoder increases the number of channels by a factor of 2 and the
1173    second leaves it unchanged.

1174    Each convolutional block consists of a GroupNorm layer with group size 16 (except for the first
1175    convolutional layer in the first encoder, which has group size 1), followed by a 3D convolutional
1176    layer with kernel size 3 and the appropriate number of input and output channels, followed by a
1177    ReLU activation layer.

1178    After the encoder, the 3D U-Net then has three decoder blocks mapping output level $i + 1$ and
1179    input level $i$ to output level $i$ for $i \in \{0,1,2\}$. Output level 3 is defined to be the same as input
1180    level 3. Each decoder layer consists of an $2 \times 2 \times 2$ upsampling layer which upsamples output
1181    level $i$ via interpolation, followed by a concatenation layer that concatenates it to input level $i -$
1182    1 along the channel axis, followed by two convolutional blocks. The first convolutional block
1183    decreases the number of channels by a factor of 2 and the second convolutional block leaves the
1184    number of channels unchanged. After the final decoder layer, a $1 \times 1$ convolutional layer is
1185    applied to increase the number of output channels to the desired 185.

1186    *Training Inputs*

1187    We trained the AutoCellLabeler network on a set of 81 human-annotated NeuroPAL images,
1188    with 10 images withheld for validation and another 11 withheld for testing. Each training dataset
1189    contained three components: image, label, and weight. The images were $4 \times 64 \times 120 \times 284$ ,
1190    with the first dimension corresponding to channel: we spectrally isolated each of the four
1191    fluorescent proteins NLS-mNeptune 2.5, NLS-CyOFP1, NLS-mTagBFP2, and NLS-tagRFP
1192    using our previously described imaging setup[35], described in detail above. The training images
1193    were then created by registering all of the images to the NLS-tagRFP image as described above,
1194    cropping all of them to $64 \times 120 \times 284$ dimensions $(z, y, x)$, and then stacking them along the
1195    channel axis to be $4 \times 64 \times 120 \times 284$ (in the reverse order that they were for the
1196    BrainAlignNet).

1197    To create the labels, we ran our segmentation U-Net on each such image to generate ROIs
1198    corresponding to neurons in these images. Humans then manually annotated the images and
1199    assigned a label and a confidence to these ROIs. These confidence values ranged from 1-5, with
1200    5 being the maximum. For network training, only confidence-1 labels were excluded while all
1201    labels from confidence 2 through 5 were included. We then made a list $\ell$ of length 185: the
1202    background label, and all 184 labels that were ever assigned in any of the human-annotated
1203    images. This list contained all neurons expected to be in the *C. elegans* head with the exceptions
1204    of ADFR, AVFR, RMHL, RMHR, and SABD, as these neurons were not labeled in any dataset.
1205    The list also contained six other possible classes corresponding to neurons in the anterior portion
1206    of the ventral cord: VA01, VB01, VB02, VD01, DD01, and DB02, as well as the classes "glia"
1207    and "granule" to denote non-neuronal objects that fluoresce (and might be labeled with an ROI),

1208 and the class "RMH?" as the human labelers were never able to disambiguate whether their
1209 "RMH" labels corresponded to RMHL or RMHR.

1210 Due to a data processing glitch, labels for 2 of the 81 training datasets were imported incorrectly;
1211 validation and testing datasets were unaffected. This resulted in those datasets effectively having
1212 random labels during training. We are currently re-training all versions of the AutoCellLabeler
1213 network and expect their performance to modestly increase once this is rectified.

1214 For each image, the human labels were transformed into matrices $L$ with dimensions
1215 $185 \times 64 \times 120 \times 284$ via one-hot encoding, so that $L[n, z, y, x]$ denotes whether the pixel at
1216 position $(x, y, z)$ has label $\ell[n]$. Specifically, we set $L[n, z, y, x]$ for $n > 0$ to be 1 if the pixel at
1217 position $(x, y, z)$ corresponded to an ROI that the human labeled as $\ell[n]$, and 0 otherwise. For
1218 example, the fourth element of $\ell$ was I2L (i.e., $\ell[3] = $ "I2L"), so $L[3, z, y, x]$ would be 1 in the
1219 ROI labeled as I2L and 0 everywhere else. The first label (i.e., $n = 0$) corresponded to the
1220 background, which was 1 if all other channels were 0, and 0 otherwise.

1221 Finally, we create a weight matrix $W$ of dimensions $64 \times 120 \times 284$ (in the code, this matrix
1222 has dimensions $185 \times 64 \times 120 \times 284$, but the loss function is mathematically equivalent to the
1223 version presented here). The entries of $W$ are determined by the following set of rules for
1224 weighting each corresponding pixel in the human label matrix $L$:

1225 • $W[z, y, x] = 1$ for all $x, y, z$ with the background label, i.e. $L[0, z, y, x] = 1$
1226 • $W[z, y, x] = \frac{130}{N(l_r)} f(c_r)$ if there is an ROI at $(x, y, z)$ with label $l_r$ that has confidence $c_r$.
1227   Here $N(l_r)$ is the number of ROIs across all datasets (train, validation and testing) with
1228   the label $l_r$. This makes neurons with fewer labels more heavily weighted in training.
1229   Additionally, $f$ is a function that weighs labels based on human confidence score $c_r$,
1230   where $c_r \in \{2, 3, 4, 5\}$. Specifically, $f(2) = 50$, $f(3) = 600$, $f(4) = 900$, and $f(5) =$
1231   1000. The number 130 was the maximum number of times that any neuronal label (e.g.:
1232   not "granule" or "glia") was detected across all of the training datasets.
1233 For the "no weight" network described in Figure 4, all entries of this matrix were set to 1.

1234 *Loss function*

1235 The loss function is pixel-wise weighted cross-entropy loss. This is computed as:

$$1236 \quad \textbf{Loss} = -\frac{1}{d_x d_y d_z K} \sum_{n=0}^{K-1} \sum_{x=0}^{d_x-1} \sum_{y=0}^{d_y-1} \sum_{z=0}^{d_z-1} W[z, y, x] L[n, z, y, x] \log \left( \frac{e^{P[n,z,y,x]}}{\sum_{m=0}^{K-1} e^{P[m,z,y,x]}} \right)$$

1237

1238 Here $(d_z, d_y, d_x)$ are the image dimensions (64, 120, 284), $K$ is the number of total labels (i.e.,
1239 the length of $\ell$), and $(n, x, y, z)$ are indices within label and image dimensions. $W$ and $L$ are as
1240 defined above, and $P$ is the prediction (output) of the network. In this way, the network has a
1241 lower loss if $P[n, z, y, x]$ is high when $L[n, z, y, x] = 1$ (ie: the network got the label right), as

1242    then the softmax $\log\left(\frac{e^{P[n,z,y,x]}}{\sum_{m=0}^{K-1} e^{P[m,z,y,x]}}\right)$ term will be close to 0 and therefore multiply $L[n, z, y, x]$

1243    by a small (negative) number, resulting in an overall small (positive) loss. The $W[z, y, x]$ term

1244    makes it so the network cares more about pixels and labels with high weight – in particular, it

1245    cares more about foreground labels $n > 0$ and about higher-confidence and rarer labels.

1246    *Evaluation metric*

1247    The evaluation metric is weighted mean intersection-over-union (IoU) across channels. Let $A$ be

1248    the network's argmax label matrix. Specifically, $A[n, z, y, x] = 1$ when $P[n, z, y, x] =$

1249    $\max_m P[m, z, y, x]$ and $A[n, z, y, x] = 0$ otherwise. Then the evaluation metric is defined as:

1250

$$\textbf{MeanIoU} \approx \frac{1}{K} \sum_{n=0}^{K-1} \frac{\sum_{x=0}^{d_x-1} \sum_{y=0}^{d_y-1} \sum_{z=0}^{d_z-1} W[z,y,x] \cdot L[n.z,y,x] \cdot A[n,z,y,x]}{\sum_{x=0}^{d_x-1} \sum_{y=0}^{d_y-1} \sum_{z=0}^{d_z-1} W[z,y,x] \cdot \max(L[n,z,y,x], A[n,z,y,x])}$$

1251    In this manner, if the network is always correct, $A = L$, the numerator and denominator will be

1252    equal, and the evaluation score will be 1. Similarly, if the network is always wrong, the

1253    evaluation score will be 0. (In the code, this metric is slightly different from the version

1254    presented here due to additional complexity with the $W$ matrix having a nonuniform extra

1255    dimension, but they act very similarly.)

1256    *Optimizer*

1257    The network was optimized with the Adam optimizer with a learning rate of $10^{-4}$.

1258    *Data augmentation*

1259    The following data augmentations are performed on the training data. One augmentation is

1260    generated for each iteration, in the following order. The same augmentation is applied to the

1261    image, label, and weight matrices, except that contrast adjustment and noise are not used for the

1262    label and weight matrices. Missing pixels are set to the median of the image, or to 0 for the label

1263    and weight matrices. Interpolation is linear for the images and nearest-neighbors for label and

1264    weight. Full parameter settings such as strength or range of each augmentation are given in the

1265    parameter file (see below).

1266    - **Rotation.** The rotations in the $xy$ plane and $yz$ plane are much larger than the rotation in

1267        the $xz$ plane because the worm is oriented to lay roughly along the $x$ axis, and the physics

1268        of the coverslip are such that it cannot rotate about the $y$ axis.

1269    - **Translation.** The image is translated.

1270    - **Scaling.** The image is scaled.

1271    - **Shearing.** The image is sheared.

1272    - **B-Spline Deformation.** Evenly-spaced control points are chosen and a random

1273        piecewise-cubic B-Spline deformation is generated between them. Additionally, a second

1274        B-Spline deformation with the same control points is generated that focuses on

1275        deformations in the $xy$ plane designed to resemble worm bending. The two transforms are

1276        added and then executed.

30

1277    -    **Rotation by multiples of 90 degrees.** The image is rotated.
1278    -    **Contrast adjustment.** Each channel is adjusted separately.
1279    -    **Gaussian blur.** Gaussian blur is added to the image, in a gradient along the *z*-axis. The
1280         gradient is intended to mimic the optical effect of the image becoming blurrier farther
1281         away from the objective.
1282    -    **Gaussian noise.** Added to the image, with each pixel being sampled independently.
1283    -    **Poisson noise.** Added to the image, with each pixel being sampled independently.
1284    *"Less aug" network training*

1285    We trained a version of the network with some of our custom augmentations disabled, to see
1286    how important they were to the overall performance, compared with the other more standard data
1287    augmentations. The specific augmentations that were disabled were:

1288    -    The second B-Spline deformation focusing on deformations in the *xy* plane
1289    -    Contrast adjustment
1290    -    Gaussian blur
1291    *Parameter file*

1292    The full parameter files are available at:

1293    https://github.com/flavell-lab/pytorch-3dunet/tree/master/AutoCellLabeler_parameters

1294    They include augmentation hyperparameters and various other settings not listed here. There is a
1295    different parameter file for each version of the network, though in most cases the differences are
1296    simply the number of input channels. If a user installs the pytorch-3dunet package from that
1297    GitHub repository and replace the paths to the training and validation data with their locations on
1298    your computer, they can train it with the exact settings we used here. Training will require a
1299    GPU with at least 48GB of VRAM.

1300    *Evaluation*

1301    During evaluation, an additional softmax layer is applied to convert network output into
1302    probabilities. Let $I$ be the input image and let $P$ be the network's output (after the softmax layer).
1303    Then at every pixel $(x, y, z)$, the network's output array $P[n, z, y, x]$ represents the probability
1304    that this pixel has label $\ell[n]$.

1305    *ROI image creation*

1306    To convert the output into labels, we first ran our previously-described neuron segmentation
1307    network[35,46] on the tagRFP channel of the NeuroPAL image. Specifically, since this
1308    segmentation network was trained on lower-SNR freely-moving data, we ran it on a lower-SNR
1309    copy of the tagRFP channel. (This copy was one of the 60 images we averaged together to get
1310    the higher-SNR image fed to AutoCellLabeler.)

1311    The segmentation network and subsequent watershed post-processing[35] were then used to
1312    generate a matrix $R$ with dimensions $284 \times 120 \times 64$ (same as the original tagRFP image).
1313    Each pixel in $R$ contains an index, either 0 for background or a positive integer indicating a
1314    specific neuron. The segmentation network and watershed algorithms were designed such that all

31

1315 pixels belonging to a specific neuron have the same index, and pixels belonging to any other
1316 neuron have different indices. We define an ROI $R_i = \{(x, y, z) \mid R[x, y, z] = i\}$.

*ROI label assignment*

1318 We now wish to use AutoCellLabeler to assign a label to ROI $R_i$. To do this, we first generate a
1319 mask matrix $M_i$ with the same dimensions as $R$, defined by:

1320 • $M_i[x, y, z] = 0$ if $R[x, y, z] \neq i$
1321 • $M_i[x, y, z] = 0.01$ if $R[x, y, z] = i$ and there exists $(X, Y, Z)$ face-adjacent to $(x, y, z)$
1322   such that $R[X, Y, Z] \neq i$.
1323 • $M_i[x, y, z] = 1$ otherwise.
1324 Here, the 0.01 entries are provided to the edges of the ROI so as to weight the central pixels of
1325 each ROI more heavily when determining the neuron's identity.

1326 Finally, we define a prediction matrix $D$ that allows us to determine the label of each ROI and
1327 the corresponding confidence of each label. Letting $V$ be the number of distinct nonzero values
1328 in $R$ (ie: the number of ROIs) and $K = 185$ be the number of possible labels (as before), we
1329 define a $V \times K$ prediction matrix $D$ whose $(i, n)$th entry represents the probability that ROI $R_i$
1330 has label $n$ as follows:

1331
$$D[i, n] = \frac{\sum_{xyz} M[i, x, y, z] P[n, z, y, x]}{\sum_{xyz} M[i, x, y, z]}$$

1332 Here the sums are taken over all pixels in the image.

1333 Note that because of the additional softmax layer, we have $\sum_n D[i, n] = 1$ for all $i$. From this, we
1334 can then define the label index of ROI $R_i$ to be $n_i = \text{argmax}_n D[i, n]$. From this, we can define
1335 its label to be $\ell[n_i]$, and the confidence of that label to be $D[i, n_i]$.

*ROI Label Postprocessing*

1337 After all ROIs are assigned a label, they are sorted by confidence in descending order. The ROIs
1338 are iterated through in this order, with each ROI being assigned its most likely label and the set
1339 of all assigned labels being tracked. If an ROI $R_i$ has its most likely label $l_i$ already assigned to a
1340 different ROI $R_j$, the distance between the centroids of ROIs $R_i$ and $R_j$ is computed. If this
1341 distance is small enough, the collision is likely due to over-segmentation by the segmentation U-
1342 Net (i.e., ROIs $R_i$ and $R_j$ are actually the same neuron). In this case, they are assigned the same
1343 label. Otherwise, the collision is likely due to a mistake on the part of AutoCellLabeler, and the
1344 label for ROI $R_i$ is deleted. (i.e. the higher-confidence label for ROI $R_j$ is kept and the lower-
1345 confidence label $R_i$ is discarded.)

1346 Additionally, ROIs are checked for under-segmentation. This rarely happens when the
1347 segmentation U-Net incorrectly merges two neurons into the same ROI. This is assessed by
1348 checking how many pixels in the ROI $R_i$ have predictions other than the full ROI label index $n_i$.
1349 Specifically, we count the number of pixels with $P[n, x, y, z] > 0.75$ within $R_i$ for some $n \neq n_i$.
1350 If there exists at least 10 pixels with label $n \neq n_i$, or 20% of the pixels in the ROI are labeled as

32

1351 $n \neq n_i$ in this way, it is plausible that the ROI contains parts of another neuron. In this case, the
1352 label for that ROI is deleted.

1353 Most neuron classes in the *C. elegans* brain are bilaterally symmetric and have two distinct cell
1354 bodies on the left and right part of the animal. These are genetically identical and therefore have
1355 exactly the same shape and color, which can often make it difficult to distinguish between them.
1356 For most applications, it is also usually unnecessary to distinguish between them since they
1357 typically have nearly-identical activity and function. In some cases, AutoCellLabeler can be
1358 confident in the neuron class but uncertain about the L/R subclass, assigning a probability of
1359 >10% to both L and R subclasses. In this case, we do not assign a specific subclass, instead
1360 assigning a label only for the main class with the sum of its confidence for either of the two
1361 subclasses. We note that this is only done for the L/R subclass – other neurons can also have D/V
1362 subclasses, but these are typically functionally distinct, so we require the network to
1363 disambiguate D/V for all neuron classes.

1364 Finally, certain neuron classes were present few times in our manually-labeled data, making it
1365 more likely for the network to mislabel them due to lack of training data, and simultaneously
1366 making it difficult for us to assess its performance on these neuron classes due to the lack of
1367 testing data where they were labeled. We deleted any AutoCellLabeler labels corresponding to
1368 one of these classes, which were ADF, AFD, AVF, AVG, DB02, DD01, RIF, RIG, RMF, RMH,
1369 SAB, SABV, SIAD, SIBD, VA01, and VD01. Additionally, there are other fluorescent cell types
1370 in the worm's head. AutoCellLabeler was trained to label them as either "glia" or "granule", to
1371 avoid mislabeling them as neurons, and any AutoCellLabeler labels of "glia" or "granule" were
1372 deleted to ensure all of our analyses are based on actual neuron labels.

1373 Altogether, these postprocessing heuristics resulted in deleting network labels for only 6.3% of
1374 ROIs with confidence 4 or greater human neuron labels (ie: not "granule" or "glia").

1375

1376 **CePNEM Simulation Analysis (Figure 4E)**

1377 To assess performance of our AutoCellLabeler network on the SWF415 strain, we could not
1378 compare its labels to human labels since humans do not know how to label neurons in this strain.
1379 Therefore, we used functional information about neuron activity patterns to assess accuracy of
1380 the network. We used our previously-described CePNEM model to do this[35]. Briefly, CePNEM
1381 fits a single neural activity trace to the animal's behavior to extract parameters about how that
1382 neuron represents information about the animal's behavior. CePNEM fits a posterior distribution
1383 for each parameter, and statistical tests run on that posterior are used to determine encoding of
1384 behavior. For example, if nearly all parameter sets in the CePNEM posterior for a given neuron
1385 have the property that they predict the neuron's activity is higher when the animal is reversing,
1386 then CePNEM would assign a reversal encoding to that neuron.

1387 By doing this analysis in NeuroPAL animals where the identity of each neural trace is known,
1388 we have previously created an atlas of neural encoding of behavior[35]. This atlas revealed a set of
1389 neurons that have consistent encodings across animals: AVA, AVE, RIM, and AIB encode

33

1390  reverse locomotion; RIB, AVB, RID, and RME encode forward locomotion; SMDD encodes
1391  dorsal head curvature; and SMDV and RIV encode ventral head curvature. Based on this prior
1392  knowledge, we decided to quantify the fraction $f$ of labeled neurons with the expected activity-
1393  behavior coupling. For example, if AutoCellLabeler labeled 10 neurons as AVA and 7 of them
1394  encoded reverse locomotion when fit by CePNEM, this fraction would be 0.7.

1395  However, this fraction is not necessarily an accurate estimate of AutoCellLabeler's accuracy. For
1396  example, it might have been possible for AutoCellLabeler to mislabel a neuron as AVA that
1397  happened to encode reverse locomotion in that animal, thus making the incorrect label appear
1398  accurate. On the other hand, CePNEM is limited by statistical power, and can sometimes fail to
1399  detect the appropriate encoding. This could make a correct label appear inaccurate.

1400  To correct for these factors, we ran a simulation analysis to try to estimate the fraction $p$ of labels
1401  that were correct. To do this, we iterated through every one of AutoCellLabeler's labels that was
1402  one of the consistent-encoding neuron classes (i.e. one of the neurons listed above). In each
1403  simulation, we assign labels to neurons in the following manner: with probability $p_{sim}$ (i.e., the
1404  fraction of labels estimated by our simulation to be correct), the label was reassigned to a random
1405  neuron that was given that label by a human in a NeuroPAL animal (at confidence 3 or greater);
1406  with probability $1 - p_{sim}$, the label was reassigned to a random neuron in the same (SWF415)
1407  animal. In this way, the simulation controls for both of the possible inaccuracies outlined above.
1408  Then the fraction $f_{sim}$ of labeled neurons with the expected encoding was computed for each
1409  simulation. 1000 simulation trials were run for each value of $p_{sim}$, which ranged from 0 to 100 –
1410  the mean and standard deviation of these trials are shown in Figure 4E. We then computed the
1411  probability $p_{sim}$ for which $f_{sim}$ was in closest agreement to $f$, which was 69% (dashed vertical
1412  line). This is our estimate for the ground-truth correct label probability $p$.

1413

1414  **CellDiscoveryNet**

1415  *Network Architecture*

1416  The architecture of CellDiscoveryNet uses the same LocalNet backbone from DeepReg that
1417  BrainAlignNet uses, with the following modifications to the architecture and training procedure
1418  (these modifications are currently in the multicolor branch):

1419  - The input images to CellDiscoveryNet are $284 \times 120 \times 64 \times 4$ instead of
1420    $284 \times 120 \times 64$.
1421  - The image concatenation layer in CellDiscoveryNet concatenates the moving and fixed
1422    images along the existing channel dimension instead of adding a new channel dimension.
1423    Effectively, this means that the output of that layer (and input to the LocalNet backbone)
1424    is now $284 \times 120 \times 64 \times 8$ instead of $284 \times 120 \times 64 \times 2$.
1425  - The affine data augmentation procedure was adjusted to first construct a 3D affine
1426    transformation, then independently apply that same transformation to each channel in the
1427    4D input images.

1428    -    In the output warping layer, the DDF is now applied independently to each channel of the
1429         moving image to create the predicted fixed image.

1430    *Loss function*

1431    The loss function in CellDiscoveryNet is a weighted sum of image loss and regularization loss.
1432    As this is an entirely unsupervised learning procedure, label loss is not used.

1433    The image loss component has weight 1 and uses GNCC instead of LNCC used by
1434    BrainAlignNet. The GNCC loss is computed as:

$$\mathbf{GNCC} = \frac{1}{C}\sum_{c=1}^{C}\frac{\sum_{x=0}^{d_x-1}\sum_{y=0}^{d_y-1}\sum_{z=0}^{d_z-1}(F[x,y,z]-\mu_{Fc})(P[x,y,z]-\mu_{Pc})}{d_xd_yd_z\sigma_{Fc}^2\sigma_{Pc}^2}$$

1436    Where $(d_x, d_y, d_z, C)$ are the dimensions of the image, $F$ is the fixed image, $P$ is the predicted
1437    fixed image (i.e.: DDF-transformed moving image), $C$ is the number of channels, $\mu_{Ic}$ is the mean
1438    of image $I$ in channel $c$, and $\sigma_{Ic}^2$ is the variance of image $I$ in channel $c$.

1439    The regularization loss terms are as before for BrainAlignNet, except with weights 0 for the axis
1440    difference norm, 0.05 for the gradient norm, 0.05 for the nonrigid penalty, and 0.0025 for the
1441    difference norm.

1442    *Training data*

1443    CellDiscoveryNet was trained on 3,240 pairs of $284 \times 120 \times 64 \times 4$ images, comprising every
1444    possible pair combination of 81 distinct images. These were the same 81 images used to train
1445    AutoCellLabeler. Each pair consisted of a moving image and a fixed image. Both images were
1446    pre-processed by setting the dynamic range of pixel intensities to [0, 1], independently for each
1447    channel.

1448    Each moving image was additionally pre-processed by using our previously-described GPU-
1449    accelerated Euler registration to coarsely align it to the corresponding fixed image. This
1450    registration was run on the NLS-tagRFP channel, and the Euler transform fit to that channel was
1451    then independently applied to each other channels to generate the full transformed moving
1452    image.

1453    There were 45 validation image pairs (from 10 validation images), and 1,866 testing image pairs.
1454    The testing image pairs added 11 additional images, and consisted of all pairs not present in
1455    either the training or validation data. (So, for example, a registration problem between an image
1456    in the training data and an image in the validation data would count as a testing image pair, since
1457    the network never saw that image pair in training or validation.) The split of images in the
1458    validation and testing data was identical to that for AutoCellLabeler.

1459    The network was trained for 600 epochs with the Adam optimizer and a learning rate of $10^{-4}$.
1460    Full training parameter settings are available at https://github.com/flavell-
1461    lab/DeepReg/blob/multicolor/CellDiscoveryNet/train_config.yaml

1462

35

1463 **ANTSUN 2U**

1464 To convert the CellDiscoveryNet registration outputs into neuron labels across animals, we
1465 created a modified version of our ANTSUN image processing pipeline. We skipped the pre-
1466 processing steps since the images were already pre-processed, used 102 four-channel images
1467 instead of 1600 one-channel images, set the registration graph to be the complete graph (except
1468 each pair of images is only registered once and not once in each direction), substituted
1469 BrainAlignNet with CellDiscoveryNet for the nonrigid registration step, and skipped the trace
1470 extraction steps of ANTSUN (stopping after it computed linked neuron IDs).

1471 We also modified the heuristic function in the matrix that was subject to clustering to better
1472 account for the nature of this multi-spectral data. Specifically, we removed the marker channel
1473 brightness heuristic $a_{ij}$ since brightness of neurons relative to the mean ROI is not likely to be
1474 well conserved across different animals. We replaced it with a more problem-specific heuristic:
1475 color. Specifically, the color $C_i$ of an ROI $R_i$ was defined as the 4-vector of its brightness in each
1476 of the four channels, normalized to the average brightness of that ROI across the four channels.
1477 We then define

1478
$$a_{ij} = \sum_{k=1}^{4} |C_{ik} - C_{jk}|$$

1479 where $i, j$ each indicates an ROI label.

1480 In this way, $a_{ij}$ will be small if the ROIs have similar colors and large if they have different
1481 colors. We use this new color-based $a_{ij}$ in the same way in the heuristic function that we used
1482 the original brightness-based $a_{ij}$, except that we set its weight $w_4 = 7$.

1483 We did not run hyperparameter search on any of the other weight parameters $w_i$ for this dataset
1484 to avoid overfitting to the 102 animals included in it, instead leaving them all at their default
1485 values from the original ANTSUN 2.0 pipeline (with the one exception of $w_8$ which we set to 0
1486 here in light of having much fewer animals than we did timepoints). We hypothesize that
1487 performance may increase even further upon hyperparameter search, though this would likely
1488 require considerably more data for testing. The only exception was that we varied parameter $w_7$,
1489 which controls the precision vs recall tradeoff. Larger values of $w_7$ result in more, but less
1490 accurate, detected clusters; each cluster corresponds to a single neuron class label. We elected to
1491 use a value of $w_7 = 10^{-9}$ for all displayed results; the full tradeoff curve is available in Figure
1492 5f.

1493 *Accuracy metric*

1494 By construction, clusters in ANTSUN 2U should correspond to individual neuron classes. To
1495 compute its accuracy, we checked whether clusters indeed only correspond to single neuron
1496 classes. Let $L(r, a)$ be the function mapping ROI $r$ in animal $a$ to its human label (ignoring L/R),
1497 and let $C_i$ be a set of $(r, a)$ values belonging to the same cluster. We can then define $L_i$ to be the

36

1498    set of labels in $C_i$: $L_i = \{L(r, a) | (r, a) \in C_i, L(r, a) \neq \text{UNKNOWN}\}$. Then let $F_i$ be the most
1499    frequent label in $L_i$. We can then define the accuracy of ANTSUN 2U as follows:

1500
$$\textbf{Accuracy} = \frac{\sum_{i \in S} \sum_{l \in L_i} \delta_{lF_i}}{\sum_{i \in S} |L_i|}$$

1501    Here $|L_i|$ is the number of elements in the set $L_i$, $S$ is the set of all clusters with $|L_i| > 2$ (which
1502    included all but one cluster in our data with $w_7 = 10^{-9}$), and $\delta$ is the Kronecker delta function.

1503
1504

1505    **REFERENCES**
1506    1.   Alon, S., Goodwin, D.R., Sinha, A., Wassie, A.T., Chen, F., Daugharthy, E.R., Bando, Y.,
1507       Kajita, A., Xue, A.G., Marrett, K., et al. (2021). Expansion sequencing: Spatially precise in
1508       situ transcriptomics in intact biological systems. Science *371*, eaax2656.
1509       https://doi.org/10.1126/science.aax2656.

1510    2.   Chen, K.H., Boettiger, A.N., Moffitt, J.R., Wang, S., and Zhuang, X. (2015). RNA imaging.
1511       Spatially resolved, highly multiplexed RNA profiling in single cells. Science *348*, aaa6090.
1512       https://doi.org/10.1126/science.aaa6090.

1513    3.   Ke, R., Mignardi, M., Pacureanu, A., Svedlund, J., Botling, J., Wählby, C., and Nilsson, M.
1514       (2013). In situ sequencing for RNA analysis in preserved tissue and cells. Nat. Methods *10*,
1515       857–860. https://doi.org/10.1038/nmeth.2563.

1516    4.   Moen, E., Bannon, D., Kudo, T., Graf, W., Covert, M., and Van Valen, D. (2019). Deep
1517       learning for cellular image analysis. Nat. Methods *16*, 1233–1246.
1518       https://doi.org/10.1038/s41592-019-0403-1.

1519    5.   Stirling, D.R., Swain-Bowden, M.J., Lucas, A.M., Carpenter, A.E., Cimini, B.A., and
1520       Goodman, A. (2021). CellProfiler 4: improvements in speed, utility and usability. BMC
1521       Bioinformatics *22*, 433. https://doi.org/10.1186/s12859-021-04344-9.

1522    6.   Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead,
1523       S., Berg, A.C., Lo, W.-Y., et al. (2023). Segment Anything. Preprint at arXiv,
1524       https://doi.org/10.48550/arXiv.2304.02643 https://doi.org/10.48550/arXiv.2304.02643.

1525    7.   Zou, J., Gao, B., Song, Y., and Qin, J. (2022). A review of deep learning-based deformable
1526       medical image registration. Front. Oncol. *12*, 1047215.
1527       https://doi.org/10.3389/fonc.2022.1047215.

1528    8.   Zidane, M., Makky, A., Bruhns, M., Rochwarger, A., Babaei, S., Claassen, M., and Schürch,
1529       C.M. (2023). A review on deep learning applications in highly multiplexed tissue imaging
1530       data analysis. Front. Bioinforma. *3*, 1159381. https://doi.org/10.3389/fbinf.2023.1159381.

1531    9.   Geuenich, M.J., Hou, J., Lee, S., Ayub, S., Jackson, H.W., and Campbell, K.R. (2021).
1532       Automated assignment of cell identity from single-cell multiplexed imaging and proteomic
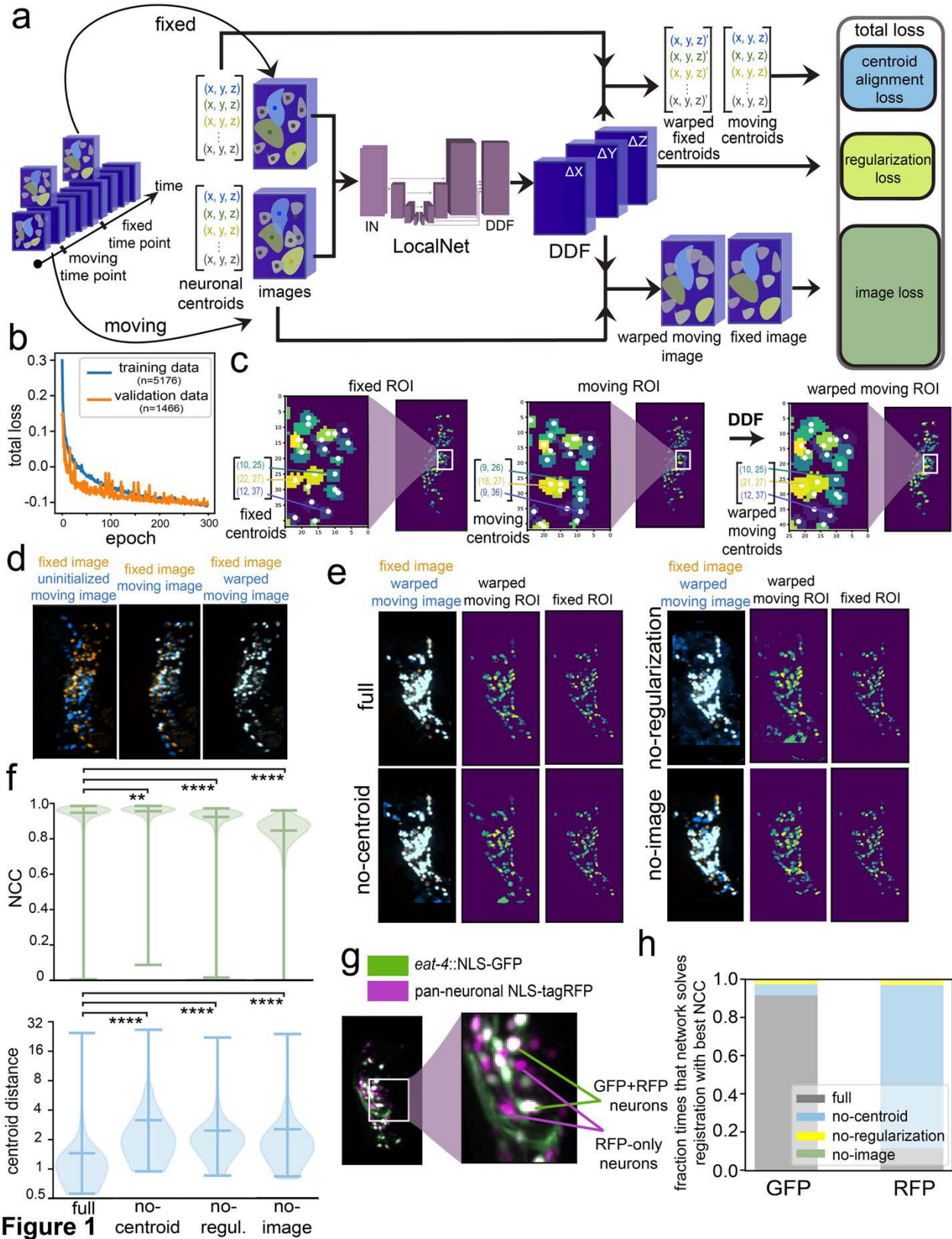1533       data. Cell Syst. *12*, 1173-1186.e5. https://doi.org/10.1016/j.cels.2021.08.012.

37

10. Amitay, Y., Bussi, Y., Feinstein, B., Bagon, S., Milo, I., and Keren, L. (2023). CellSighter: a neural network to classify cells in highly multiplexed images. Nat. Commun. *14*, 4302. https://doi.org/10.1038/s41467-023-40066-7.

11. Brbić, M., Cao, K., Hickey, J.W., Tan, Y., Snyder, M.P., Nolan, G.P., and Leskovec, J. (2022). Annotation of spatially resolved single-cell data with STELLAR. Nat. Methods *19*, 1411–1418. https://doi.org/10.1038/s41592-022-01651-8.

12. White, J.G., Southgate, E., Thomson, J.N., and Brenner, S. (1986). The structure of the nervous system of the nematode Caenorhabditis elegans. Philos. Trans. R. Soc. Lond. B. Biol. Sci. *314*, 1–340. https://doi.org/10.1098/rstb.1986.0056.

13. Witvliet, D., Mulcahy, B., Mitchell, J.K., Meirovitch, Y., Berger, D.R., Wu, Y., Liu, Y., Koh, W.X., Parvathala, R., Holmyard, D., et al. (2021). Connectomes across development reveal principles of brain maturation. Nature *596*, 257–261. https://doi.org/10.1038/s41586-021-03778-8.

14. Cook, S.J., Jarrell, T.A., Brittin, C.A., Wang, Y., Bloniarz, A.E., Yakovlev, M.A., Nguyen, K.C.Q., Tang, L.T.-H., Bayer, E.A., Duerr, J.S., et al. (2019). Whole-animal connectomes of both Caenorhabditis elegans sexes. Nature *571*, 63–71. https://doi.org/10.1038/s41586-019-1352-7.

15. Prevedel, R., Yoon, Y.-G., Hoffmann, M., Pak, N., Wetzstein, G., Kato, S., Schrödel, T., Raskar, R., Zimmer, M., Boyden, E.S., et al. (2014). Simultaneous whole-animal 3D imaging of neuronal activity using light-field microscopy. Nat. Methods *11*, 727–730. https://doi.org/10.1038/nmeth.2964.

16. Schrödel, T., Prevedel, R., Aumayr, K., Zimmer, M., and Vaziri, A. (2013). Brain-wide 3D imaging of neuronal activity in Caenorhabditis elegans with sculpted light. Nat. Methods *10*, 1013–1020. https://doi.org/10.1038/nmeth.2637.

17. Nguyen, J.P., Shipley, F.B., Linder, A.N., Plummer, G.S., Liu, M., Setru, S.U., Shaevitz, J.W., and Leifer, A.M. (2016). Whole-brain calcium imaging with cellular resolution in freely behaving Caenorhabditis elegans. Proc. Natl. Acad. Sci. U. S. A. *113*, E1074-1081. https://doi.org/10.1073/pnas.1507110112.

18. Venkatachalam, V., Ji, N., Wang, X., Clark, C., Mitchell, J.K., Klein, M., Tabone, C.J., Florman, J., Ji, H., Greenwood, J., et al. (2016). Pan-neuronal imaging in roaming Caenorhabditis elegans. Proc. Natl. Acad. Sci. U. S. A. *113*, E1082-1088. https://doi.org/10.1073/pnas.1507109113.

19. Flavell, S.W., and Gordus, A. (2022). Dynamic functional connectivity in the static connectome of Caenorhabditis elegans. Curr. Opin. Neurobiol. *73*, 102515. https://doi.org/10.1016/j.conb.2021.12.002.

20. Kramer, T.S., and Flavell, S.W. (2024). Building and integrating brain-wide maps of nervous system function in invertebrates. Curr. Opin. Neurobiol. *86*, 102868. https://doi.org/10.1016/j.conb.2024.102868.

1572 21. Lagache, T., Hanson, A., Pérez-Ortega, J.E., Fairhall, A., and Yuste, R. (2021). Tracking
1573   calcium dynamics from individual neurons in behaving animals. PLoS Comput. Biol. *17*,
1574   e1009432. https://doi.org/10.1371/journal.pcbi.1009432.

1575 22. Hanson, A., Reme, R., Telerman, N., Yamamoto, W., Olivo-Marin, J.-C., Lagache, T., and
1576   Yuste, R. (2024). Automatic monitoring of neural activity with single-cell resolution in
1577   behaving Hydra. Sci. Rep. *14*, 5083. https://doi.org/10.1038/s41598-024-55608-2.

1578 23. Wen, C., Miura, T., Voleti, V., Yamaguchi, K., Tsutsumi, M., Yamamoto, K., Otomo, K.,
1579   Fujie, Y., Teramoto, T., Ishihara, T., et al. (2021). 3DeeCellTracker, a deep learning-based
1580   pipeline for segmenting and tracking cells in 3D time lapse images. eLife *10*, e59187.
1581   https://doi.org/10.7554/eLife.59187.

1582 24. Nejatbakhsh, A., Varol, E., Yemini, E., Venkatachalam, V., Lin, A., Samuel, A.D.T., and
1583   Paninski, L. (2020). Extracting neural signals from semi-immobilized animals with
1584   deformable non-negative matrix factorization. Preprint at bioRxiv,
1585   https://doi.org/10.1101/2020.07.07.192120 https://doi.org/10.1101/2020.07.07.192120.

1586 25. Christensen, R.P., Bokinsky, A., Santella, A., Wu, Y., Marquina-Solis, J., Guo, M.,
1587   Kovacevic, I., Kumar, A., Winter, P.W., Tashakkori, N., et al. (2015). Untwisting the
1588   Caenorhabditis elegans embryo. eLife *4*, e10070. https://doi.org/10.7554/eLife.10070.

1589 26. Ardiel, E.L., Lauziere, A., Xu, S., Harvey, B.J., Christensen, R.P., Nurrish, S., Kaplan, J.M.,
1590   and Shroff, H. (2022). Stereotyped behavioral maturation and rhythmic quiescence in C.
1591   elegans embryos. eLife *11*, e76836. https://doi.org/10.7554/eLife.76836.

1592 27. Ryu, J., Nejatbakhsh, A., Torkashvand, M., Gangadharan, S., Seyedolmohadesin, M., Kim,
1593   J., Paninski, L., and Venkatachalam, V. (2024). Versatile multiple object tracking in sparse
1594   2D/3D videos via deformable image registration. PLoS Comput. Biol. *20*, e1012075.
1595   https://doi.org/10.1371/journal.pcbi.1012075.

1596 28. Nguyen, J.P., Linder, A.N., Plummer, G.S., Shaevitz, J.W., and Leifer, A.M. (2017).
1597   Automatically tracking neurons in a moving and deforming brain. PLoS Comput. Biol. *13*,
1598   e1005517. https://doi.org/10.1371/journal.pcbi.1005517.

1599 29. Yu, X., Creamer, M.S., Randi, F., Sharma, A.K., Linderman, S.W., and Leifer, A.M. (2021).
1600   Fast deep neural correspondence for tracking and identifying neurons in C. elegans using
1601   semi-synthetic training. eLife *10*, e66410. https://doi.org/10.7554/eLife.66410.

1602 30. Wu, Y., Wu, S., Wang, X., Lang, C., Zhang, Q., Wen, Q., and Xu, T. (2022). Rapid detection
1603   and recognition of whole brain activity in a freely behaving Caenorhabditis elegans. PLoS
1604   Comput. Biol. *18*, e1010594. https://doi.org/10.1371/journal.pcbi.1010594.

1605 31. Nejatbakhsh, A., and Varol, E. (2021). Neuron Matching in C. elegans With Robust
1606   Approximate Linear Regression Without Correspondence. In, pp. 2837–2846.

1607 32. Deng, H., Yu, J., and Venkatachalam, V. (2024). Neuron tracking in C. elegans through
1608   automated anchor neuron localization and segmentation. In Computational Optical Imaging

1609     and Artificial Intelligence in Biomedical Sciences (SPIE), pp. 84–95.
1610     https://doi.org/10.1117/12.3001982.

1611  33. Park, C.F., Barzegar-Keshteli, M., Korchagina, K., Delrocq, A., Susoy, V., Jones, C.L.,
1612     Samuel, A.D.T., and Rahi, S.J. (2024). Automated neuron tracking inside moving and
1613     deforming C. elegans using deep learning and targeted augmentation. Nat. Methods *21*, 142–
1614     149. https://doi.org/10.1038/s41592-023-02096-3.

1615  34. Wang, D., Lu, Z., Xu, Y., Wang, Z.I., Santella, A., and Bao, Z. (2019). Cellular structure
1616     image classification with small targeted training samples. IEEE Access Pract. Innov. Open
1617     Solut. *7*, 148967–148974. https://doi.org/10.1109/access.2019.2940161.

1618  35. Atanas, A.A., Kim, J., Wang, Z., Bueno, E., Becker, M., Kang, D., Park, J., Kramer, T.S.,
1619     Wan, F.K., Baskoylu, S., et al. (2023). Brain-wide representations of behavior spanning
1620     multiple timescales and states in C. elegans. Cell *186*, 4134-4151.e31.
1621     https://doi.org/10.1016/j.cell.2023.07.035.

1622  36. Toyoshima, Y., Wu, S., Kanamori, M., Sato, H., Jang, M.S., Oe, S., Murakami, Y.,
1623     Teramoto, T., Park, C., Iwasaki, Y., et al. (2020). Neuron ID dataset facilitates neuronal
1624     annotation for whole-brain activity imaging of C. elegans. BMC Biol. *18*, 30.
1625     https://doi.org/10.1186/s12915-020-0745-2.

1626  37. Bubnis, G., Ban, S., DiFranco, M.D., and Kato, S. (2019). A probabilistic atlas for cell
1627     identification. Preprint at arXiv, https://doi.org/10.48550/arXiv.1903.09227
1628     https://doi.org/10.48550/arXiv.1903.09227.

1629  38. Chaudhary, S., Lee, S.A., Li, Y., Patel, D.S., and Lu, H. (2021). Graphical-model framework
1630     for automated annotation of cell identities in dense cellular images. eLife *10*, e60321.
1631     https://doi.org/10.7554/eLife.60321.

1632  39. Skuhersky, M., Wu, T., Yemini, E., Nejatbakhsh, A., Boyden, E., and Tegmark, M. (2022).
1633     Toward a more accurate 3D atlas of C. elegans neurons. BMC Bioinformatics *23*, 195.
1634     https://doi.org/10.1186/s12859-022-04738-3.

1635  40. Yemini, E., Lin, A., Nejatbakhsh, A., Varol, E., Sun, R., Mena, G.E., Samuel, A.D.T.,
1636     Paninski, L., Venkatachalam, V., and Hobert, O. (2021). NeuroPAL: A Multicolor Atlas for
1637     Whole-Brain Neuronal Identification in C. elegans. Cell *184*, 272-288.e11.
1638     https://doi.org/10.1016/j.cell.2020.12.012.

1639  41. Varol, E., Nejatbakhsh, A., Hobert, O., Mena, G., Paninski, L., Yemini, E., and Sun, R.
1640     (2020). Statistical atlas of C. elegans neurons (Springer).

1641  42. Sprague, D.Y., Rusch, K., Dunn, R.L., Borchardt, J.M., Bubnis, G., Chiu, G.C., Wen, C.,
1642     Suzuki, R., Chaudhary, S., Dichter, B., et al. (2024). Unifying community-wide whole-brain
1643     imaging datasets enables robust automated neuron identification and reveals determinants of
1644     neuron positioning in C. elegans. Preprint at bioRxiv,
1645     https://doi.org/10.1101/2024.04.28.591397 https://doi.org/10.1101/2024.04.28.591397.

1646 43. Ma, J., Zhao, J., and Yuille, A.L. (2016). Non-Rigid Point Set Registration by Preserving
1647   Global and Local Structures. IEEE Trans. Image Process. *25*, 53–64.
1648   https://doi.org/10.1109/TIP.2015.2467217.

1649 44. Fu, Y., Brown, N.M., Saeed, S.U., Casamitjana, A., Baum, Z.M. c, Delaunay, R., Yang, Q.,
1650   Grimwood, A., Min, Z., Blumberg, S.B., et al. (2020). DeepReg: a deep learning toolkit for
1651   medical image registration. J. Open Source Softw. *5*, 2705.
1652   https://doi.org/10.21105/joss.02705.

1653 45. Hu, Y., Modat, M., Gibson, E., Ghavami, N., Bonmati, E., Moore, C.M., Emberton, M.,
1654   Noble, J.A., Barratt, D.C., and Vercauteren, T. (2018). Label-driven weakly-supervised
1655   learning for multimodal deformable image registration. In 2018 IEEE 15th International
1656   Symposium on Biomedical Imaging (ISBI 2018), pp. 1070–1074.
1657   https://doi.org/10.1109/ISBI.2018.8363756.

1658 46. Pradhan, S., Madan, G.K., Kang, D., Bueno, E., Atanas, A.A., Kramer, T.S., Dag, U., Lage,
1659   J.D., Gomes, M.A., Lu, A.K.-Y., et al. (2024). Pathogen infection induces sickness behaviors
1660   by recruiting neuromodulatory systems linked to stress and satiety in C. elegans. Preprint at
1661   bioRxiv, https://doi.org/10.1101/2024.01.05.574345
1662   https://doi.org/10.1101/2024.01.05.574345.

1663 47. Wolny, A., Cerrone, L., Vijayan, A., Tofanelli, R., Barro, A.V., Louveaux, M., Wenzl, C.,
1664   Strauss, S., Wilson-Sánchez, D., Lymbouridou, R., et al. (2020). Accurate and versatile 3D
1665   segmentation of plant tissues at cellular resolution. eLife *9*, e57613.
1666   https://doi.org/10.7554/eLife.57613.

1667 48. Sharma, A.K., Randi, F., Kumar, S., Dvali, S., and Leifer, A.M. (2024). TWISP: A
1668   Transgenic Worm for Interrogating Signal Propagation in C. elegans. Genetics, iyae077.
1669   https://doi.org/10.1093/genetics/iyae077.

1670 49. Kato, S., Kaplan, H.S., Schrödel, T., Skora, S., Lindsay, T.H., Yemini, E., Lockery, S., and
1671   Zimmer, M. (2015). Global brain dynamics embed the motor command sequence of
1672   Caenorhabditis elegans. Cell *163*, 656–669. https://doi.org/10.1016/j.cell.2015.09.034.

1673 50. Li, Z., Liu, J., Zheng, M., and Xu, X.Z.S. (2014). Encoding of both analog- and digital-like
1674   behavioral outputs by one C. elegans interneuron. Cell *159*, 751–765.
1675   https://doi.org/10.1016/j.cell.2014.09.056.

1676 51. Chalfie, M., Sulston, J.E., White, J.G., Southgate, E., Thomson, J.N., and Brenner, S. (1985).
1677   The neural circuit for touch sensitivity in Caenorhabditis elegans. J. Neurosci. Off. J. Soc.
1678   Neurosci. *5*, 956–964.

1679 52. Gordus, A., Pokala, N., Levy, S., Flavell, S.W., and Bargmann, C.I. (2015). Feedback from
1680   network states generates variability in a probabilistic olfactory circuit. Cell *161*, 215–227.
1681   https://doi.org/10.1016/j.cell.2015.02.018.

1682 53. Luo, L., Wen, Q., Ren, J., Hendricks, M., Gershow, M., Qin, Y., Greenwood, J., Soucy, E.R.,
1683   Klein, M., Smith-Parker, H.K., et al. (2014). Dynamic encoding of perception, memory, and

41

1684 movement in a C. elegans chemotaxis circuit. Neuron *82*, 1115–1128.
1685 https://doi.org/10.1016/j.neuron.2014.05.010.

1686 54. Wang, Y., Zhang, X., Xin, Q., Hung, W., Florman, J., Huo, J., Xu, T., Xie, Y., Alkema, M.J.,
1687 Zhen, M., et al. (2020). Flexible motor sequence generation during stereotyped escape
1688 responses. eLife *9*, e56942. https://doi.org/10.7554/eLife.56942.

1689 55. Ben Arous, J., Tanizawa, Y., Rabinowitch, I., Chatenay, D., and Schafer, W.R. (2010).
1690 Automated imaging of neuronal activity in freely behaving Caenorhabditis elegans. J.
1691 Neurosci. Methods *187*, 229–234. https://doi.org/10.1016/j.jneumeth.2010.01.011.

1692 56. Lim, M.A., Chitturi, J., Laskova, V., Meng, J., Findeis, D., Wiekenberg, A., Mulcahy, B.,
1693 Luo, L., Li, Y., Lu, Y., et al. (2016). Neuroendocrine modulation sustains the C. elegans
1694 forward motor state. eLife *5*, e19887. https://doi.org/10.7554/eLife.19887.

1695 57. Kim, J., Rustam, S., Mosquera, J.M., Randell, S.H., Shaykhiev, R., Rendeiro, A.F., and
1696 Elemento, O. (2022). Unsupervised discovery of tissue architecture in multiplexed imaging.
1697 Nat. Methods *19*, 1653–1661. https://doi.org/10.1038/s41592-022-01657-2.

1698 58. Liu, C.C., Greenwald, N.F., Kong, A., McCaffrey, E.F., Leow, K.X., Mrdjen, D., Cannon,
1699 B.J., Rumberger, J.L., Varra, S.R., and Angelo, M. (2023). Robust phenotyping of highly
1700 multiplexed tissue imaging data using pixel-level clustering. Nat. Commun. *14*, 4618.
1701 https://doi.org/10.1038/s41467-023-40068-5.

1702 59. Dag, U., Nwabudike, I., Kang, D., Gomes, M.A., Kim, J., Atanas, A.A., Bueno, E., Estrem,
1703 C., Pugliese, S., Wang, Z., et al. (2023). Dissecting the functional organization of the
1704 C. elegans serotonergic system at whole-brain scale. Cell *186*, 2574-2592.e20.
1705 https://doi.org/10.1016/j.cell.2023.04.023.

1706
1707

1708

**Figure 1**

1709

43

1710 **Figure 1. BrainAlignNet can perform non-rigid registration to align the neurons in the *C.***
1711 ***elegans* head**

1712     **(A)** Network training pipeline. The network takes in a pair of images and a pair of centroid
1713          position lists corresponding to the images at two different time points (fixed and moving).
1714          (In the LocalNet diagram, this is represented as "IN". Intermediate cuboids represent
1715          intermediate representations of the images at various stages of network processing. In
1716          reality, the cuboids are four-dimensional, but we represent them with three dimensions
1717          (up/down is $x$, left/right is $y$, in/out is channel, and we omit $z$) for visualization purposes.
1718          Spaces and arrows between cuboids represent network blocks, layers, and information
1719          flow. See Methods for a detailed description of network architectures.) Image pairs were
1720          selected based on the similarity of worm postures (see Methods). The fixed and moving
1721          images were pre-registered using an Euler transformation, translating and rotating the
1722          moving images to maximize their cross-correlation with the fixed images. The fixed and
1723          moving neuron centroid positions were obtained by computing the centers of the same
1724          neurons in both the fixed and moving images as a list of $(x, y, z)$ coordinates. This
1725          information was available since we had previously extracted calcium traces from these
1726          videos using a previous, slow version of our image analysis pipeline. The network
1727          outputs a Dense Displacement Field (DDF), a 4-D tensor that indicates a coordinate
1728          transformation from fixed image coordinates to moving image coordinates. The DDF is
1729          then used to transform the moving images and fixed centroids to resemble the fixed
1730          images and moving centroids. During training, the network is tasked with learning a DDF
1731          that transforms the centroids and images in a way that minimizes the centroid alignment
1732          and image loss, as well as the regularization loss (see Methods). Note that, after training,
1733          only images (not centroids) need to be input into the network to align the images.
1734     **(B)** Network loss curves. The training and validation loss curves show that validation
1735          performance plateaued around 300 epochs of training.
1736     **(C)** Example of registration outcomes on neuronal ROI images. The network-learned DDF
1737          warps the neurons in the moving image ('moving ROIs'). The warped-moving ROIs are
1738          meant to be closer to the fixed ROIs. Each neuron is uniquely colored in the ROI images
1739          to represent its identity. The centroids of these neurons are represented by the white dots.
1740          Here, we take a $z$-slice of the 3-D fixed and moving ROI blocks on the $x$-$y$ plane to show
1741          that the DDF can warp the $x$ and $y$ coordinates of the moving centroids to align with the $x$
1742          and $y$ coordinates of the fixed centroids with one-pixel precision.
1743     **(D)** Example of registration outcomes on tagRFP images. We show the indicated image
1744          blocks as Maximal Intensity Projections (MIPs) along the $z$-axis, overlaying the fixed
1745          image (orange) with different versions of the moving image (blue). While the fixed image
1746          remains untransformed, the uninitialized moving image (left) gets warped by an Euler
1747          transformation (middle) and a network-learned DDF (right) to overlap with the fixed
1748          image.
1749     **(E)** Registration outcomes shown on example tagRFP and ROI images for four different
1750          trained networks. We randomly selected one registration problem from one of the testing
1751          datasets and tasked the trained networks with creating a DDF to warp the moving (RFP)
1752          image and moving ROI onto the fixed (RFP) image and fixed ROI. The full network with
1753          full loss function aligns neurons in both RFP and ROI images almost perfectly. For the
1754          networks trained without the centroid alignment loss, regularization loss, or image loss—
1755          while keeping the rest of the training configurations identical—the resulting DDF is
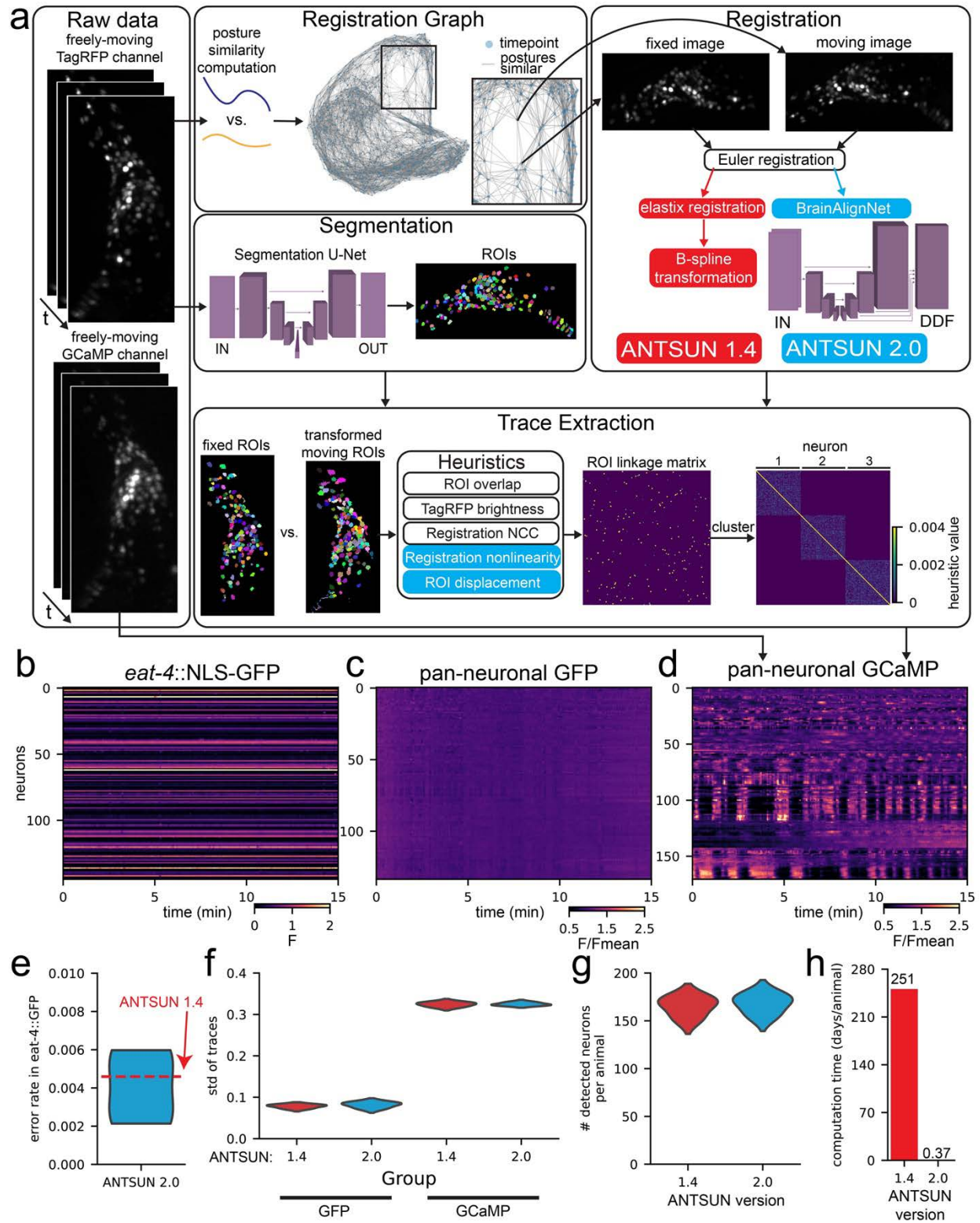
44

1756      unable to fully align the neurons and displays unrealistic deformation (closely inspect the
1757      warped moving ROI images).

1758   **(F)** Evaluation of registration performance on testing datasets under four different network
1759      configurations. Here, we evaluated 80-100 problems per animal for all animals in the
1760      testing data. Two performance metrics are shown. Normalized cross-correlation (NCC,
1761      top) quantifies alignment of the fixed and warped moving RFP images, where a score of
1762      one indicates perfect alignment. Centroid distance (bottom) is measured as the mean
1763      Euclidean distance between the centroids of all neurons in the fixed ROI and the
1764      centroids of their corresponding neurons in the warped moving ROI; a distance of 0
1765      indicates perfect alignment. All violin plots are accompanied by lines indicating the
1766      minimum, mean, and maximum values. **p<0.01, ***p<0.001, ****p<0.0001,
1767      distributions of registration metrics (NCC and centroid distance) were compared pairwise
1768      across all four versions of the network with the Wilcoxon signed rank test.

1769   **(G)** Example image of the head of an animal from a strain that expresses both pan-neuronal
1770      NLS-tagRFP and eat-4::NLS-GFP. The neurons expressing both NLS-tagRFP and eat-
1771      4::NLS-GFP is a subset of all the neurons expressing pan-neuronal NLS-tagRFP.

1772   **(H)** A comparison of the registration qualities of the four trained registration networks: full
1773      network, no-centroid alignment loss, no-regularization loss, no-image loss. Each network
1774      was evaluated on four datasets in which both pan-neuronal NLS-tagRFP and *eat-4*::NLS-
1775      GFP are expressed, examining 3927 registration problems per dataset. For a total of
1776      15,708 registration problems, each network was tasked with registering the tagRFP
1777      images. The resulting DDFs from the tagRFP registrations were also used to register the
1778      *eat-4*::GFP images. For each channel in each problem, we determined which of the four
1779      networks had the highest performance (i.e. highest NCC). Note that the no-centroid
1780      alignment network performs the best of the RFP channel, but not in the GFP channel.
1781      Instead, the full network performs the best in the GFP channel. This suggests that the
1782      network without the centroid alignment loss deforms RFP images in a manner that does
1783      not accurately move the neurons to their correct locations (i.e. scrambles the pixels).

1784

**Figure 2**

46

**Figure 2. BrainAlignNet supports calcium trace extraction with high accuracy and high SNR**

**(A)** Diagram of ANTSUN 1.4 and 2.0, which are two full calcium trace extraction pipelines that only differ with regards to image registration. Raw tagRFP channel data is input into the pipeline, which submits image pairs with similar worm postures for registration using either elastix (ANTSUN 1.4; red) or BrainAlignNet (ANTSUN 2.0; blue). The registration is used to transform neuron ROIs identified by a segmentation U-Net (the cuboid diagram is represented as in Figure 1A). These are input into a heuristic function (ANTSUN 2.0-specific heuristics shown in blue) which defines an ROI linkage matrix. Clustering this matrix then yields neuron identities.

**(B)** Sample dataset from an *eat-4*::NLS-GFP strain, showing ratiometric (GFP/tagRFP) traces without any further normalization. This strain has some GFP+ neurons (bright horizontal lines) as well as some GFP- neurons (dark horizontal lines, which have F~0). Registration artifacts between GFP+ and GFP- neurons would be visible as bright points in GFP- traces or dark points in GFP+ traces.

**(C)** Sample dataset from a pan-neuronal GFP strain, showing F/Fmean fluorescence. Any variation visible here is noise.

**(D)** Sample dataset from a pan-neuronal GCaMP strain, showing F/Fmean fluorescence. Robust calcium dynamics are visible in most neurons.

**(E)** Violin plot of the error rate of ANTSUN 2.0 registration across four *eat-4*::NLS-GFP animals, computed based on mismatches between GFP+ and GFP- neurons in the *eat-4*::NLS-GFP strain. Dashed red line shows the error rate of ANTSUN 1.4. Note that all error rates are <1%.

**(F)** Violin plots of the standard deviation of traces across three animals per strain (pan-neuronal GFP or pan-neuronal GCaMP).

**(G)** Violin plots of the number of detected neurons across three pan-neuronal GCaMP animals for the two different ANTSUN versions (1.4 or 2.0).

**(H)** Computation time to process one animal based on ANTSUN version (1.4 or 2.0). ANTSUN 1.4 was run on a computing cluster that provided an average of 32 CPU cores per registration problem; computation time is the total number of CPU hours used (ie: the time it would have taken to run ANTSUN 1.4 registration locally on a comparable 32-core machine). ANTSUN 2.0 was run locally on NVIDIA A4000, A5500, and A6000 graphics cards.
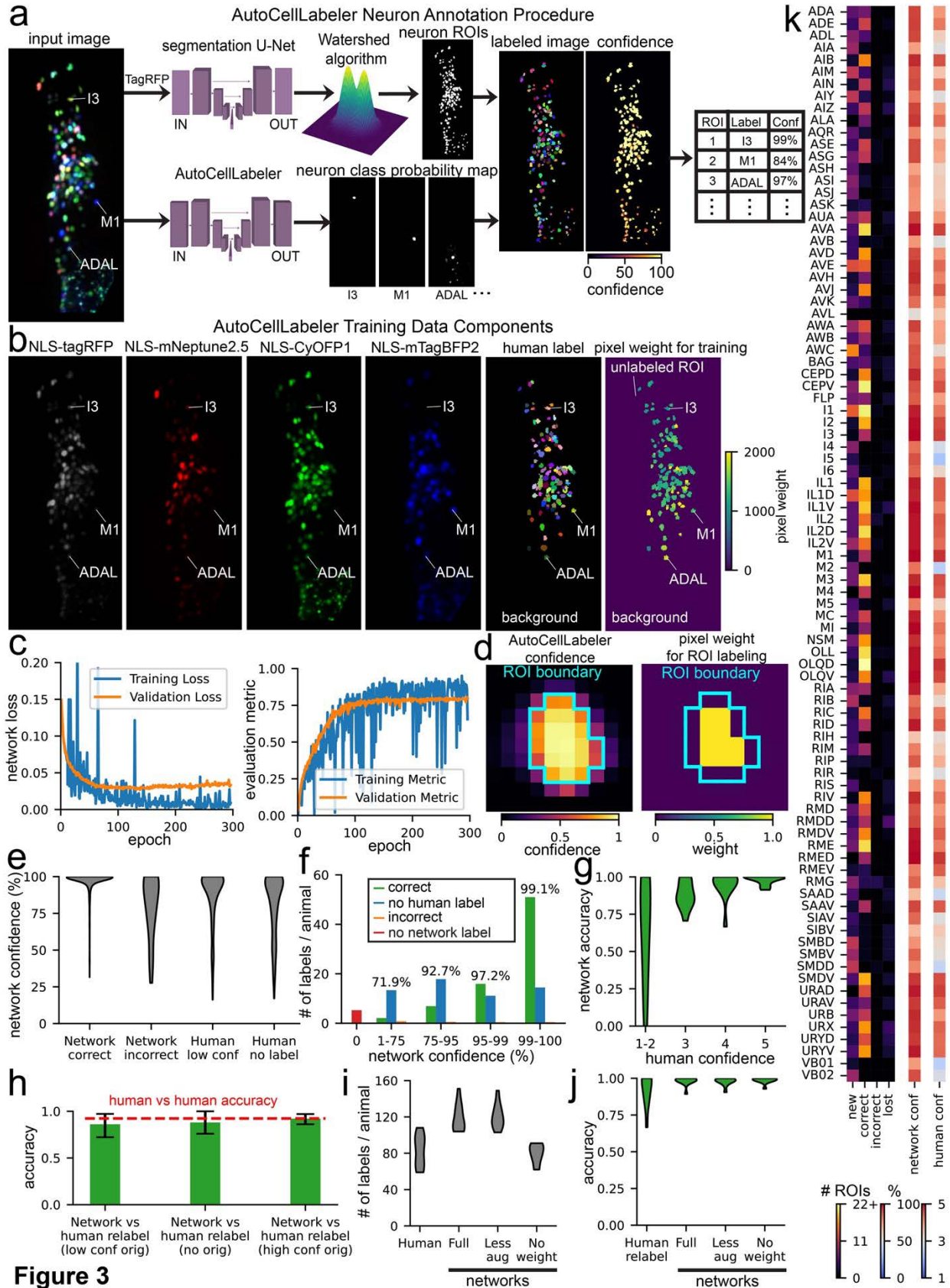
**Figure 3**

1821

48

**Figure 3. The AutoCellLabeler Network can automatically annotate >100 neuronal cell types in the *C. elegans* head**

(A) Procedure by which AutoCellLabeler generates labels for neurons. First, the tagRFP component of a multi-spectral image is passed into a segmentation neural network, which extracts neuron ROIs, labeling each pixel as an arbitrary number with one number per neuron. Then, the full multi-spectral image is input into AutoCellLabeler, which outputs a probability map. This probability map is applied to the ROIs to generate labels and confidence values for those labels. The network cuboid diagrams are represented as in **Figure 1A**.

(B) AutoCellLabeler's training data consists of a set of multi-spectral images (NLS-tagRFP, NLS-mNeptune2.5, NLS-CyOFP1, and NLS-mTagBFP2), human neuron labels, and a pixel weighting matrix based on confidence and frequency of the human labels that controls how much each pixel is weighted in AutoCellLabeler's loss function.

(C) Pixel-weighted cross-entropy loss and pixel-weighted IoU metric scores for training and validation data. Cross-entropy loss captures the discrepancy between predicted and actual class probabilities for each pixel. The IoU metric describes how accurately the predicted labels overlap with the ground truth labels.

(D) During the label extraction procedure, AutoCellLabeler is less confident of its label on pixels near the edge of ROI boundaries. Therefore, we allow the central pixels to have much higher weight when determining the overall ROI label from pixel-level network output.

(E) Distributions of AutoCellLabeler's confidence across test datasets based on the relationship of its label to the human label ("Correct" = agree, "Incorrect" = disagree, "Human low conf" = human had low confidence, "Human no label" – human did not even guess a label for the neuron).

(F) Categorization of neurons in test datasets based on AutoCellLabeler's confidence. Here "Correct" and "Incorrect" are as in **(E)**, but "No human label" also includes low-confidence human labels. Printed percentage values are the accuracy of AutoCellLabeler on the corresponding category, computed as $\frac{correct}{correct+incorrect}$

(G) Distributions of accuracy of AutoCellLabeler's high confidence (>75%) labels on neurons across test datasets based on the confidence of the human labels.

(H) Accuracy of AutoCellLabeler compared with high-confidence labels from new human labelers on neurons in test datasets that were labeled at low confidence, not at all, or at high confidence by the original human labelers. Error bars are bootstrapped 95% confidence intervals. Dashed red line shows accuracy of new human labelers relative to the old human labelers, when both gave high confidence to their labels.

(I) Distributions of number of high-confidence labels per animal over test datasets. High confidence was 4-5 for human labels and >75% for network labels.

(J) Distributions of accuracy of high-confidence labels per animal over test datasets, relative to the original human labels.

(K) Number of ROIs per neuron class labeled at high confidence in test datasets that fall into each category, along with average confidence for all labels for each neuron class in those test datasets. "New" represents ROIs that were labeled by the network as the neuron and were not labeled by the human. "Correct" represents ROIs that were labeled by both AutoCellLabeler and the human as that neuron. "Incorrect" represents ROIs that were labeled by the network as that neuron and were labeled by the human as something else.

49

1868     "Lost" represents ROIs that were labeled by the human as that neuron and were not
1869     labeled by the network. "Network conf" represents the average confidence of the network
1870     for all its labels of that neuron. "Human conf" represents the average confidence of the
1871     human labelers for all their labels of that neuron. Neuron classes with high values in the
1872     "Correct" column and low values in the "Incorrect" column indicate a very high degree
1873     of accuracy in AutoCellLabeler's labels for those classes. If those classes also have a
1874     high value in the "New" column, it could indicate that AutoCellLabeler is able to find the
1875     neuron with high accuracy in animals where humans were unable to label it.
1876
1877

**Figure 4**

1878

51

**Figure 4. Variants of AutoCellLabeler can annotate neurons from fewer fluorescent channels and in different strains**

   **(A)** Distributions of number of high-confidence labels per animal over test datasets for the networks trained on the indicated set of fluorophores. The "tagRFP (on low SNR)" column corresponds to a network that was trained on high-SNR, tagRFP-only data and tested on low-SNR tagRFP data due to shorter exposure times in freely-moving animals.

   **(B)** Distributions of accuracy of high-confidence labels per animal over test datasets for the networks trained on the indicated set of fluorophores. The "tagRFP (on low SNR)" column is as in **(A)**.

   **(C)** Same as **Figure 3K**, except for the tagRFP-only network.

   **(D)** Accuracy vs detection tradeoff for various AutoCellLabeler versions. For each network, we can set a confidence threshold above which we accept labels. By varying this threshold, we can produce a tradeoff between accuracy of accepted labels ($x$-axis) and number of labels per animal ($y$-axis) on test data. Each curve in this plot was generated in this manner. The "tagRFP-only (on low SNR)" values are as in **(A)**. The "tagRFP-only (on freely-moving)" values come from evaluating the tagRFP-only network on 100 randomly-chosen timepoints in the freely-moving (tagRFP) data for each test dataset. The final labels were then computed on each immobilized ROI by averaging together the 100 labels and finding the most likely label. To ensure fair comparison to other networks, only immobilized ROIs that were matched to the freely-moving data were considered for any of the networks in this plot (unlike Extended Data Figure 2A, which used all available ROIs).

   **(E)** Evaluating the performance of tagRFP-only AutoCellLabeler on data from another strain SWF415, where there is pan-neuronal NLS-GCaMP7f and pan-neuronal NLS-mNeptune2.5. Notably, the pan-neuronal promoter used for NLS-mNeptune2.5 differs from the pan-neuronal promoter used for NLS-tagRFP in NeuroPAL. Performance here was quantified by computing the fraction of network labels with the correct expected activity-behavior relationships in the neuron class (y-axis; quantified by whether an encoding model showed significant encoding; see Methods). For example, when the label was the reverse-active AVA neuron, did the corresponding calcium trace show higher activity during reverse? The blue line shows the expected fraction as a function of the true accuracy of the network (x-axis), computed via simulations (see Methods). Orange circle shows the actual fraction when AutoCellLabeler was evaluated on SWF415. Based on this, the dashed line shows estimated true accuracy of this labeling.
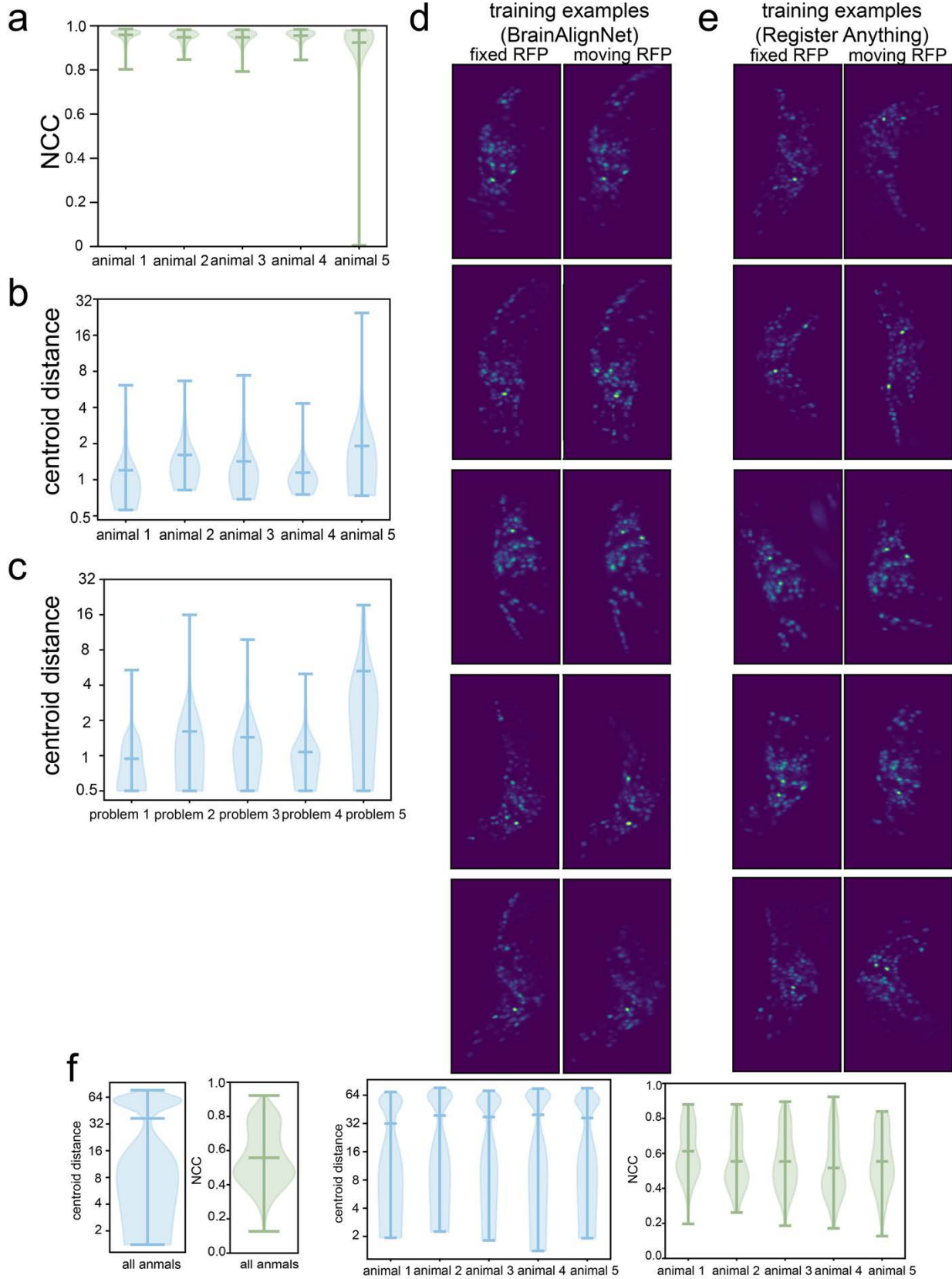
1915
1916    **Figure 5**

53

1917 **Figure 5. CellDiscoveryNet and ANTSUN 2U can perform unsupervised cell type discovery**
1918 **by analyzing data across different *C. elegans* animals**
1919

1920 **(A)** A schematic comparing the approaches of AutoCellLabeler and CellDiscoveryNet.
1921 AutoCellLabeler uses supervised learning, taking as input both images and manual labels
1922 for those images, and learns to label neurons accordingly. CellDiscoveryNet uses
1923 unsupervised learning, and can learn to label neurons after being trained only on images
1924 (with no labels provided).
1925 **(B)** CellDiscoveryNet training pipeline. The network takes as input two multi-spectral
1926 NeuroPAL images from two different animals. It then outputs a Dense Displacement
1927 Field (DDF), which is a coordinate transformation between the two images. It warps the
1928 moving image under this DDF, producing a warped moving image that should ideally
1929 look very similar to the fixed image. The dissimilarity between these images is the image
1930 loss component of the loss function, which is added to the regularization loss that
1931 penalizes non-linear image deformations present in the DDF.
1932 **(C)** Network loss curves. Both training and validation loss curves start to plateau around 600
1933 epochs.
1934 **(D)** Distributions of normalized cross-correlation (NCC) scores comparing the
1935 CellDiscoveryNet predictions (warped moving images) and the fixed images for each
1936 pair of registered images. These NCCs were computed on all four channels
1937 simultaneously, treating the entire image as a single 4D matrix for this purpose. The
1938 "Train" distribution contains the NCC scores for all pairs of images present in
1939 CellDiscoveryNet's training data, while the "Val+Test" distribution contains any pair of
1940 images that was not present in its training data.
1941 **(E)** Distributions of centroid distance scores based on human labels. These are computed
1942 over all (moving, fixed) image pairs on all neurons with high-confidence human labels in
1943 both moving and fixed images. The centroid distance scores represent the Euclidean
1944 distance from the network's prediction for where the neuron was and its correct location
1945 as labeled by the human. Values of a few pixels or less likely roughly indicate that the
1946 neuron was mapped to its correct location, while large values mean the neuron was mis-
1947 registered. The "Train" and "Val+Test" distributions are as in **(D)**. The "High NCC"
1948 distribution is from only (moving, fixed) image pairs where the NCC score was greater
1949 than the 90th percentile of all such NCC scores.
1950 **(F)** Labeling accuracy vs number of linked neurons tradeoff curve. Accuracy is the fraction
1951 of linked ROIs with labels matching their cluster's most frequent label (see Methods).
1952 Number of linked neurons is the total number of distinct clusters; each cluster must
1953 contain an ROI in more than half of the animals to be considered a cluster. The parameter
1954 $w_7$ describes when to terminate the clustering algorithm – higher values mean the
1955 clustering algorithm terminates earlier, resulting in more accurate but fewer detections.
1956 Red dot is the selected value $w_7 = 10^{-9}$ where 125 clusters were detected with 93%
1957 labeling accuracy.
1958 **(G)** Number of neurons labeled per animal in the 11 testing datasets. This plot compares the
1959 number of neurons labeled as follows: human labels with 4-5 confidence,
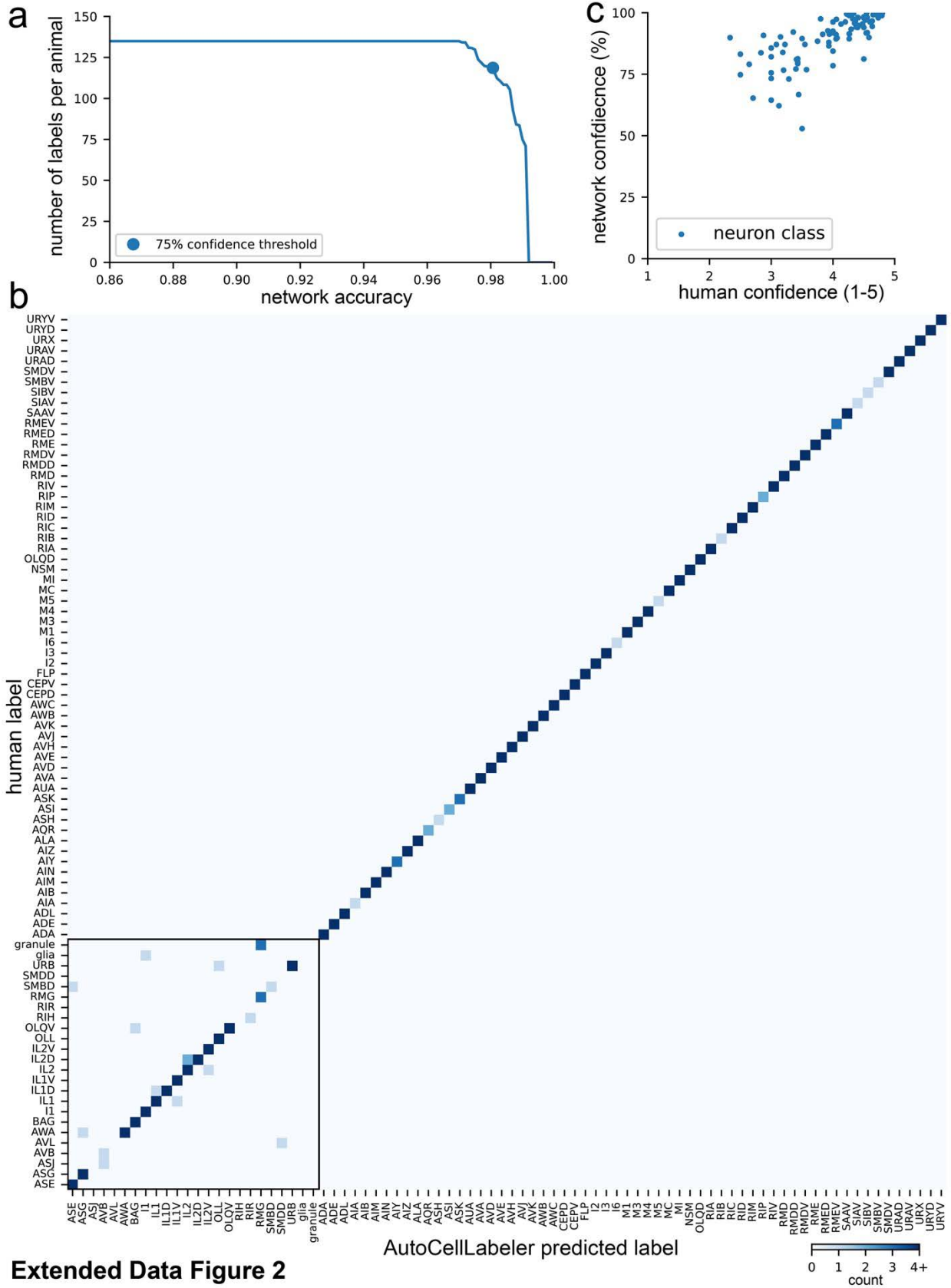
54

1960          AutoCellLabeler labels with 75% or greater confidence, and CellDiscoveryNet with

1961          ANTSUN 2U labels with parameter $w_7 = 10^{-9}$.

1962  **(H)** Accuracy of neuron labels in the 11 testing datasets. This plot defines the original human

1963          confidence 4-5 labels as ground truth. "Human relabel" are confidence 4-5 labels done by

1964          different humans (independently from the first set of human labels). AutoCellLabeler are

1965          confidence 75% or greater labels. CellDiscoveryNet labels were created by running

1966          ANTSUN 2U with $w_7 = 10^{-9}$, and defining the correct label for each cluster to be its

1967          most frequent label.

1968  **(I)** Same as **Figure 4(K)**, except using labels from CellDiscoveryNet with ANTSUN 2U.

1969          The neurons "NEW 1" through "NEW 5" are clusters that were not labeled frequently

1970          enough by humans to be able to determine which neuron class they corresponded to, as

1971          described in the main text.

1972

1973

1974

**Extended Data Figure 1**

56

**Extended Data Figure 1. Example images and performance of network trained to register arbitrary image pairs.**

(A) Performance of image registration in five different animals in the testing set. Normalized Cross-Correlation (NCC) scores of aligned tagRFP images are shown, which indicate the extent of image alignment (best achievable score is 1). 90-100 registration problems examined per animal are shown as violin plots with the overlaying lines indicating minimum, mean, and maximum values.

(B) Performance of image registration in five different animals in the testing set. Centroid distance is the average Euclidean distance between the centroids of matched neurons in each image (best achievable score is 0). 90-100 registration problems examined per animal are shown as violin plots with the overlaying lines indicating minimum, mean, and maximum values.

(C) Performance of image registration in five different registration problems (i.e. image pairs) from one example animal. Centroid distance is the average Euclidean distance between the centroids of matched neurons in that image pair (best achievable score is 0). All the centroid position distances for each registration problem as shown as violin plots with the overlaying lines indicating minimum, mean, and maximum values.

(D) Five example image pairs in the training set for BrainAlignNet. These are maximum intensity projections of the tagRFP channel, showing two different timepoints that were selected to be the fixed and moving images in each of these five registration problems.

(E) Five example image pairs in the training set for the network trained to align arbitrary image pairs, including much more challenging problems. Note that the head bending is more dissimilar for these image pairs, as compared to those in (D). Data are shown as in (D).

(F) Performance of the network trained to register arbitrary image pairs. Quantification is for testing data. We quantify centroid distance (average alignment of neuron centroids) and NCC (image similarity) as in panels (A-C). By both metrics, this network's performance is far worse than that of the BrainAlignNet presented in Fig. 1. The two panels on the right show that results are qualitatively similar for different animals in the testing set.
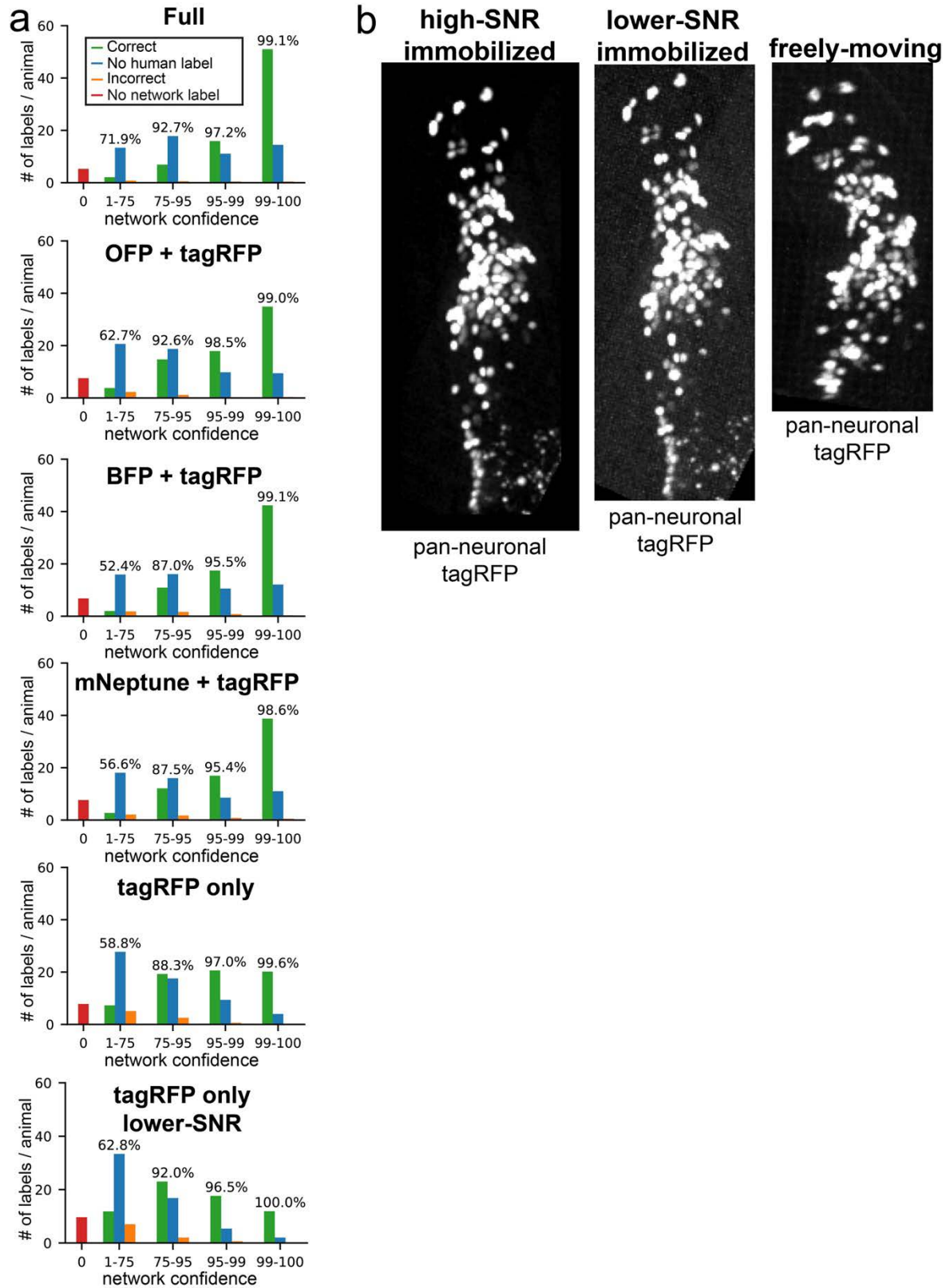
**Extended Data Figure 2**

2022

58

**Extended Data Figure 2. Further characterization of the AutoCellLabeler network**

**(A)** Tradeoff of network labeling accuracy (*x*-axis) and number of neurons labeled (*y*-axis) for the full AutoCellLabeler network. The number of neurons labeled can be varied by adjusting the threshold confidence that the network needs to achieve to label an ROI. By varying this threshold, we were able to generate this curve. This full curve captures the tradeoff and shows the 75% confidence threshold (blue circle) that we selected to use in our analyses.

**(B)** Confusion matrix showing which neurons could potentially be confused for one another by AutoCellLabeler. Note that, except for the diagonal, the matrix is mostly white, reflecting that it is mostly (98%) accurate. Neurons with some inaccuracies were clustered to the lower left (boxed region). Note that with a linear color scale the diagonal would be off-scale bright with correct labels. So we capped the colorbar range at 4 counts so as to not block the ability to see actual confusion entries. For reference, the actual mean value across the diagonal is 9.7.

**(C)** Positive correlation between human and autolabel confidence across the neuronal cell types (each cell type is a blue dot). This plot also highlights that a subset of cells are more difficult for human labelers and, therefore, also for AutoCellLabeler (i.e. the cells that are not clustered in the upper right).
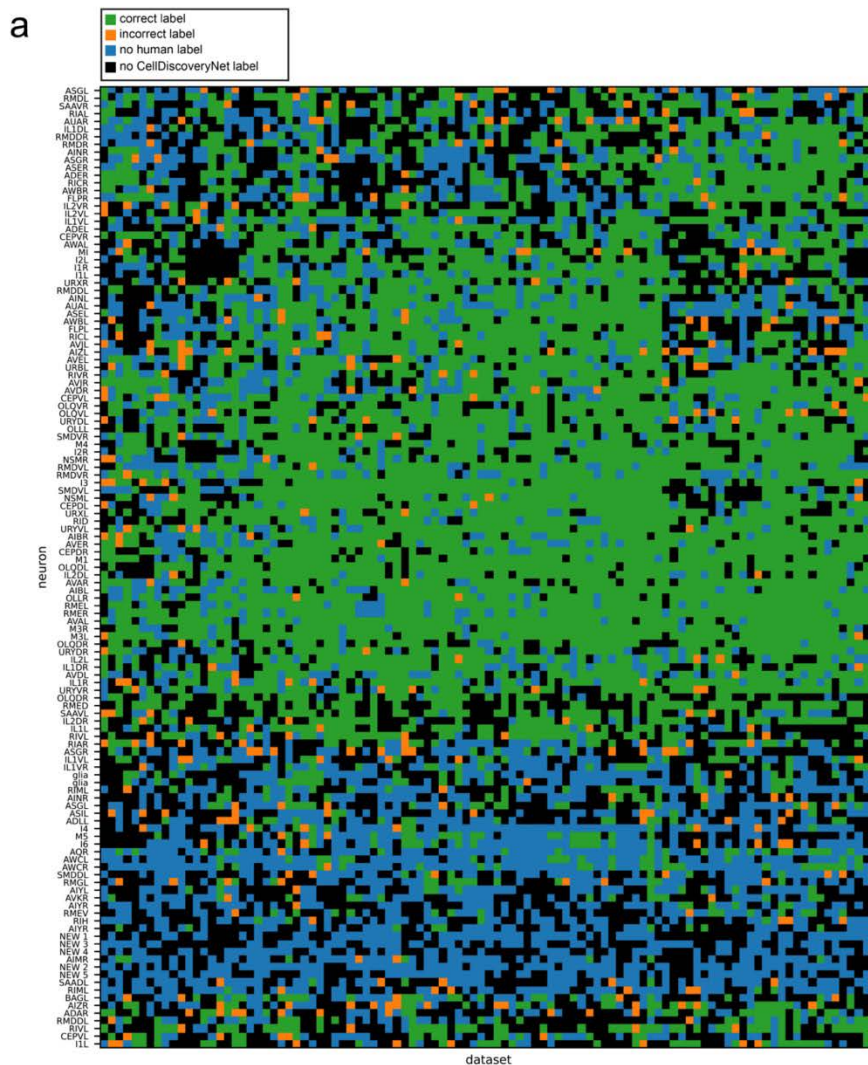
**Extended Data Figure 3**

**Extended Data Figure 3. Further characterization of the different AutoCellLabeler variants.**

   **(A)** These plots show the performance of different indicated cell annotation networks (trained and/or evaluated on different fluorophores, as indicated). Data is displayed to show network performance on different ROIs that it labels with different levels of confidence. Printed percentage values are the accuracy of AutoCellLabeler on the corresponding category, computed as $\frac{correct}{correct+incorrect}$. Note that the lower performing networks (for example, tagRFP-only) are still accurate for their high-confidence labels, and that their decreased accuracy is mostly due to a lower fraction of high-confidence labels (i.e. more cell types where the networks had low confidence in their annotations).

   **(B)** Example maximum intensity projection images of the worm in the tagRFP channel under three different imaging conditions: immobilized high-SNR (created by averaging together 60 immobilized lower-SNR images together, our typical condition for NeuroPAL imaging); immobilized lower-SNR (i.e. one of those 60 images); and freely-moving (which was taken with the same imaging settings as immobilized lower-SNR but in a freely-moving animal)

2061

**Extended Data Figure 4**

2063
2064
2065
2066

2067 **Extended Data Figure 4. Further characterization of CellDiscoveryNet and ANTSUN 2U**
2068 **performance**
2069 **(A)** Matrix of all clusters generated by running ANTSUN 2U. Each row is a distinct cluster
2070 (i.e. inferred cell type), while each column is a distinct animal. Black entries mean that
2071 the given cluster did not include any ROIs in the given animal (ie: ANTSUN 2U failed to
2072 label that cluster in that animal). Non-black entries mean that the cluster contained an
2073 ROI in that animal. Row names correspond to the most frequent human label among
2074 ROIs in the cluster (this was defined by first disambiguating most frequent neuron class,
2075 and then disambiguating L from R). Green entries correspond to cases when the given
2076 ROI's label matched the most frequent class label (row name ignoring L/R), orange
2077 entries correspond to the case when the given ROI's label did not match the most
2078 frequent class label, and blue entries mean that the given ROI did not have a high-
2079 confidence human label. The neurons "NEW 1" through "NEW 5" are clusters that were
2080 not labeled frequently enough by humans to be able to determine which neuron class they
2081 corresponded to, as described in the main text.

2082