Automated neuron tracking inside moving and deforming animals using deep learning and targeted augmentation

Core Francisco Park^{1,+}, Mahsa Barzegar Keshteli^{2,+}, Kseniia Korchagina^{2,+}, Ariane Delrocq², Vladislav Susoy¹, Corinne L. Jones³, Aravinthan D. T. Samuel¹, and Sahand Jamal Rahi^{2,*}

¹Department of Physics and Center for Brain Science, Harvard University, Cambridge, MA, USA, ²Laboratory of the Physics of Biological Systems, Institute of Physics, École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland, ³Swiss Data Science Center, École polytechnique fédérale de Lausanne (EPFL), Lausanne, Switzerland, ^{*}sahand.rahi@epfl.ch, ⁺these authors contributed equally to this work

Advances in functional brain imaging now allow sustained rapid 1 3D visualization of large numbers of neurons inside behaving ani-2 mals. To decode circuit activity, imaged neurons must be individu-3 ally segmented and tracked. This is particularly challenging when 5 the brain itself moves and deforms inside a flexible body. The field has lacked general methods for solving this problem effectively. To 6 address this need, we developed a method based on a convolutional 7 neural network (CNN) with specific enhancements which we apply 8 to freely moving Caenorhabditis elegans. For a traditional CNN 9 to track neurons across images of a brain with different postures, 10 the CNN must be trained with ground truth (GT) annotations of 11 similar postures. When these postures are diverse, an adequate 12 number of GT annotations can be prohibitively large to generate 13 manually. We introduce 'targeted augmentation', a method to au-14 tomatically synthesize reliable annotations from a few manual an-15 notations. Our method effectively learns the internal deformations 16 of the brain. The learned deformations are used to synthesize an-17 notations for new postures by deforming the manual annotations 18 of similar postures in GT images. The technique is germane to 19 3D images, which are generally more difficult to analyze than 2D 20 images. The synthetic annotations, which are added to diversify 21 training datasets, drastically reduce manual annotation and proof-22 reading. Our method is effective both when neurons are repre-23 sented as individual points or as 3D volumes. We provide a GUI 24 that incorporates targeted augmentation in an end-to-end pipeline, 25 from manual GT annotation of a few images to final proofreading 26 of all images. We apply the method to simultaneously measure ac-27 tivity in the second-layer interneurons in C. elegans: RIA, RIB, 28 and RIM, including the RIA neurite. We find that these neurons 29 show rich behaviors, including switching entrainment on and off 30 dynamically when the animal is exposed to periodic odor pulses. 31 (300 words) 32

33 Introduction

Whole-brain imaging with single-cell resolution is widely used 34 to study the neural circuits for behavior in many organisms in-35 cluding C. elegans, Drosophila, zebrafish, and hydra [1–4]. Fast 36 3D microscopes - light-sheet, spinning disk confocal, light-37 field, and multifocus microscopes - are expanding the range of 38 brain imaging with fluorescent genetically-encoded calcium in-39 dicators to many animals and their behaviors [5-10]. It is often 40 preferable to study brain dynamics in animals as they perform 41 behaviors without restraint: immobilization can change brain 42 activity [11] and many behaviors such as mating or predation 43 only occur in moving animals in natural contexts [12, 13]. 44

Analyzing whole-brain imaging datasets requires solving two
problems. One problem is segmentation – the pixels for each
neuron must be separated and identified in each image volume.
Another problem is tracking – a given neuron must be correctly
identified in every image volume. These problems are espe-

cially challenging in animals with flexible bodies like nema-50 todes or hydra. Their neurons are small, densely packed, and 51 follow complex trajectories as the brain deforms during behav-52 ior. Genetically-encoded labeling of neurons creates additional 53 challenges in imaging and analysis. Conditions for microscopy 54 can vary with different animals and different experiments. Ex-55 pression patterns of fluorescent reporters can vary from animal 56 to animal. When imaging with dim reporters inside moving an-57 imals, neuronal signals may be intermittent and neurons may 58 drop out of view at some time points. 59

A solution to segmenting and tracking neurons during brain-60 wide imaging has stringent requirements. In every neuron, ev-61 ery signal at every time-point might contain useful information. 62 Recording circuit activity requires accurate and reliable signal 63 analysis throughout all image volumes. Reliable analysis of im-64 age volumes often involves laborious manual annotation. In a 65 recent study of the mating circuit of male C. elegans, we re-66 quired 200 hours to manually annotate 76 neurons in image vol-67 umes recorded at 5 Hz in each 10 min experiment [12]. Auto-68 mated annotation is needed to accelerate the field of brain-wide 69 imaging. Because automated annotation can always generate er-70 rors, it must be followed by comprehensive proofreading. How-71 ever, if automated annotation were sufficiently fast and reliable, 72 the time required to manually proofread an entire dataset could 73 become less than the time required for full manual annotation. 74

One way to simplify the problem of neuron segmentation is to 75 restrict fluorescent labeling to cell nuclei. In each image vol-76 ume, cell nuclei form a constellation of non-overlapping and 77 nearly spherical volumes. With their stereotyped appearance, 78 cell nuclei are, in principle, suitable for segmentation by stan-79 dard methods in image processing. One method involves iden-80 tifying objects that have a certain size and convex shape along 81 all axes [14]. After objects have been segmented, they need to 82 be tracked across image volumes. One approach to the tracking 83 problem is to seek the optimal alignment of different 3D brain 84 images using methods from point-set registration [15]. These 85 techniques have been applied to C. elegans by treating cell nu-86 clei as points in space, and then searching for and registering the 87 local constellation of points surrounding each point in an image 88 volume [16]. However, these methods have not been extended 89 to segmenting and tracking neurons as 3D volumes. 90

Another approach to tracking neurons is to identify all segmented neurons in each image volume and assign them a unique label computationally. This is possible in animals like *C. elegans* where all neurons have unique identities, and can be done by building probabilistic models of the positions of all neurons using brain atlases that describe a given posture. For each brain image, one calculates the most likely distribution of unique neu-



Figure 1. 3D volumes from typical multi-neuron or whole-brain recordings of freely moving C. elegans worms after maximum intensity projection along the z-direction. A) The tail of a male worm with pan-neuronal nuclear GCaMP6s and red fluorescent proteins is shown in the presence of another, unlabeled hermaphrodite worm, indicated by arrows in the two top panels. In the bottom right panel, the head of the male worm, which is indicated by an arrow, appears in the field of view, next to the tail that is to be tracked. B) The interneuron strain carries a mixture of nuclear and cytosolic markers. The arrows indicate thin neurites.

rons [17–19]. However, this approach is not easily extended to 98 moving, deforming brains, and has only been applied to immo-99 bilized animals. [17] 100

A third way to track neurons is to characterize the collective 101 3D motion trajectories of neurons over time. In the rapidly de-102 forming body of the hydra, an effective particle motion tracking 103 algorithm has been developed that calculates the most likely set 104 of collective movements of visible neurons [20]. 105

Alternatively, neuron segmentation and tracking during brain-106 wide imaging can be viewed as problems in pattern recogni-107 tion, for which deep neural networks are ideally suited. [21] 108 Deep neural networks have been used to perform point-set reg-109 istration when tracking neurons by training the networks to find 110 the most likely alignment of 3D images. Deep neural networks 111 have also been used to learn the most likely trajectory in collec-112 tive motion tracking [22, 23]. These approaches facilitated the 113 114 tracking of neurons, but began with accurately segmented images. In addition to requiring high-quality segmentations, 3Dee-115 CellTracker [22] was tested only on worms that were compu-116 tationally straightened using additional low-magnification im-117 ages. fDLC [23] aims to allow deformations, however, it is un-118 clear how the method, which relies on point clouds, would be 119 applied to tracking neurons that are represented as 3D shapes. 120

We sought a comprehensive pipeline that begins with raw un-121 segmented brain-wide recordings and ends with proofreading 122 of all tracked and segmented neurons, each neuron represented 123 either as a key-point at its nucleus or as a 3D shape. To be gener-124 alizable, we did not want the pipeline to use information beyond 125 the brain images themselves, such as low-magnification images 126 127 of brain or body posture that are used in some approaches [16, 23]. To do this, we developed a deep learning method that is 128 both significantly faster than full manual annotation, and that 129 functions end-to-end by simultaneously segmenting and track-130 ing neurons in freely moving C. elegans. 131

Brain-wide imaging is subject to substantial experiment-to-132 experiment variations. To be robust, a convolutional neural net-133 work (CNN) must be separately trained using ground-truth data 134 that comes from each experiment. Generating annotated train-135 ing data for every experiment is labor intensive. We thus sought 136 to minimize the amount of manually annotated data required 137 for each experiment. We developed 'targeted augmentation', a 138 means of synthesizing large amounts of annotated training data 139 for the many different shapes and postures of the brain in a given 140 experiment. Our pipeline learns a set of internal deformations 141 of a freely moving C. elegans. The pipeline uses these deforma-142 tions to generate synthetic image volumes and annotations from 143 a small number of original image volumes and their manual an-144 notations. When a CNN is trained with both manual annotations 145 and synthetic annotations for each experiment, its accuracy and 146 reliability dramatically improve. 147

We further increased the accuracy and reduced the computa-148 tional burden of the CNN by developing a new architecture 149 to deploy targeted augmentation. Our low error rate for auto-150 mated annotation substantially reduces the time required for fi-151 nal proofreading. We implement all steps in the pipeline - man-152 ual annotation, analysis, and proofreading - in an easy-to-use 153 graphical user interface (GUI). The GUI also implements addi-154 tional machine learning methods to speed up the initial manual 155 annotation. 156

By reducing the amount of manually annotated training data required to train a CNN to reliably perform both neuron segmentation and tracking across time, as well as reducing the amount 159 of proofreading required to remove errors, our pipeline achieves a 3x-fold increase in analysis throughput in comparison to full manual annotation for the most challenging brain imaging problem in *C. elegans* to date from ref. [12].

157

158

160

161

162

163

Because our method works both for key-point as well as volu-164 metric segmentation and tracking, we used our pipeline to ana-165 lyze recordings of freely moving C. elegans with labeled neu-166

rons that need to be analyzed as 3D volumes, which cannot be 167 easily accomplished with previous methods. Specifically, we 168 investigated the second-layer interneurons in C. elegans: RIA, 169 RIB, and RIM. This set of neurons is thought to be important 170 for sensorimotor integration [24, 25] and each has been found to 171 be associated with different functions during C. elegans chemo-172 taxis: RIA shows compartmentalized calcium activity [26, 27], 173 in which different segments of the RIA neurite encode dorsal 174 and ventral head movements, respectively, but the soma does 175 not show prominent calcium activity. Thus, RIA has to be seg-176 mented and tracked as a 3D volume. RIB activity promotes 177 forward locomotion [2, 28-30]. RIM depolarization extends 178 reversals, while hyperpolarization extends the forward motor 179 state [31]. We asked whether these interneurons also repre-180 sented olfactory stimuli in addition to their tight link to chemo-181 tactic behavioral output. We applied periodic inputs to behaving 182 animals in the form of odor pulses to dissect fundamental net-183 work and circuit properties [32–35]. We observed that second-184 layer interneurons switched between entrained responses to the 185 sensory input or non-entrained activity, indicating switches be-186 tween states where these interneurons couple or decouple from 187 sensory information. 188

In summary, our main achievements are 1) a 'targeted augmen-189 tation' method that reduces the need for manual annotation to 190 create training datasets for a CNN that solves the segmentation 191 and tracking problem in brain imaging, 2) a new CNN archi-192 tecture that is optimized for this application, 3) a generalizable 193 pipeline, implemented in a graphical user interface (GUI), that 194 is able to track neurons as either key-points or 3D shapes us-195 ing only information contained in the brain images themselves 196 (no pre-training or additional recordings necessary), and 4) the 197 application of our method to track both the nuclei of the second-198 layer interneurons RIB and RIM as well as the whole 3D vol-199 ume of RIA, including its (thin) neurite, showing the complex 200 coupling of these neurons to sensory information. Furthermore, 201 we made the GUI and 4D image datasets for testing and fur-202 ther method development freely available, see 'Code and data 203 availability'. 204

205 **Results**

206 Whole-brain recording in behaving animals

It is now possible to measure the activity of an entire C. ele-207 gans brain with cellular resolution during animal behavior us-208 ing fast 3D imaging systems such as spinning-disk confocal mi-209 croscopy adapted for multi-color, multi-neuron, real-time track-210 ing [5, 12, 37, 38]. Brain-wide imaging in C. elegans is usually 211 performed using transgenic strains with panneuronal expression 212 of a nuclear-localized calcium indicator (e.g., GCaMP6s) and 213 red fluorescent protein (e.g., mNeptune). Alternatively, subsets 214 of neurons may be fluorescently labeled throughout their cy-215 tosols to allow recording from cell bodies and neurites. Stable, 216 red fluorescence signals are used to isolate and track neurons. 217 Green fluorescence signals indicate neuronal activity. When an 218 entire brain is captured at single-cell resolution at many vol-219 220 umes per second for minutes, neurons must be accurately segmented and individually tracked over thousands of image vol-221 umes. We sought an end-to-end analysis pipeline that would 222 automate the extraction of neuronal activities from brain-wide 223 imaging in C. elegans. 224

Coarse alignment

The brain of a moving animal exhibits substantial rotations, translations, and deformations. The first step in our image analysis is the coarse global alignment of the 3D images at the scale of the whole brain. Accurate global alignment facilitates both manual annotation and proofreading by reducing differences in neuronal positions at different time points.

225

246

247

248

249

250

We trained a convolutional neural network (CNN) with the U-232 Net architecture [39] to perform a coarse global alignment of 233 each recording. Once trained, the neural network works well 234 on datasets from different animals (see Appendix I). This global 235 alignment CNN (GA-CNN) automatically recognizes the points 236 corresponding to the anterior, posterior, and central axis of the 237 brain. Brain volumes can then be aligned across image vol-238 umes using point-set registration of the anterior, posterior, and 239 axis coordinates. Global alignment is a common problem in 240 image analysis, and alternative methods - such as identifying 241 landmark neurons, multipole matching, *OpenCV* tracking [40] 242 - would also work. At this step, we also reduce imaging noise 243 by applying a Difference-of-Gaussian (DoG) filter to each im-244 age volume. 245

Targeted augmentation

After coarse alignment and noise reduction, image volumes enter a pipeline that performs targeted augmentation of ground truth annotations for machine learning. This pipeline is illustrated in Fig. 2 A.

Ground truth image selection and manual annotation. The first 251 step is to select a small number of image volumes for ground 252 truth (GT) manual annotation. These annotations will be used 253 to train an initial CNN (iCNN) that segments and tracks neu-254 rons across the different postures that the brain can assume. At 255 this step, it is useful to select image volumes corresponding to 256 diverse postures, either individual volumes taken at regular in-257 tervals or a sequence of volumes when the animal exhibits sub-258 stantial movement. 259

There are two useful ways of labeling neurons computation-260 ally for segmentation and tracking. One may label a neuron 261 by a 'key-point', one chosen pixel inside the volume of the neu-262 ron. This is particularly convenient, for example, when only 263 the nuclei of neurons are fluorescently labeled since neighbor-264 ing nuclei are generally well separated and the correspondence 265 between the key-point and the volume of the nucleus is easy 266 to make. One may also label a neuron as a 3D volume, where 267 all pixels in the volume of the neuron are used to identify the 268 neuron. For simplicity, we will assume that we wish to perform 269 key-point tracking in the following, and discuss the relatively 270 small differences with 3D volume segmentation and tracking in 271 the section 'Segmenting and tracking volumetric objects'. 272

Identifying key-points can be partly automated, for example, by
identifying local maxima of fluorescence in each image volume.273When creating the GT manual annotation, non-rigid point-set
registration can be used to track key-points across selected im-
age volumes. Tracking by point-set registration is reliable when
the worm is nearly immobilized but requires substantial manual
correction when the worm is moving (Fig. 1).273

After key-points for all neurons are selected across specific image volumes, they become the set of GT manual annotations.





Figure 2. Illustration of the method. A: Steps involved in tracking. GT = ground truth, ML = machine learning. B: The CNN architecture used for the initial and augmented CNNs. C: The planar embedding of 3D images from the recording allows the similarity between 3D images to be measured by the Cartesian distances between their point representations in the plane. The embedding is performed by compressing the 3D image by an auto-encoder and mapping the latent space representations in the 'bottleneck' layer onto a plane using UMAP [36]. The representation of all (blue), GT (magenta), and target (orange) 3D images from a recording are shown. D: Example of the deformations performed on a GT image to match the target image. Left: maximum intensity projection of 3D images, right: initial CNN annotations of the images, which were used to perform the deformations

Training the iCNN. Using the GT manual annotations, we trained
 an initial CNN (iCNN) to automatically identify a small spher ical region of interest around each key-point. We then used the
 iCNN to make a first set of predictions of key-point locations
 across image volumes.

We implemented the iCNN in a custom architecture that we call the 3D Compact Network (3DCN) (Figs. 2 B, S1). We designed the network to associate information over large distances when predicting neuron segmentations and key-point locations. We downsample after every two convolutions to associate information over large distances with a limited kernel size. This is similar to the downward branch of the U-Net [39]. To avoid increasing the raw size of the convolutional kernel and the corresponding increase in the number of fitting parameters for a bigger receptive field, we employ Atrous Spatial Pyramidal Pooling (ASPP) introduced in ref. [41].

The output of the iCNN should be a 3D image containing all segmented and tracked neurons with the same spatial resolution as the original image. A standard way to convert a low-resolution image in the latent space representation of a CNN 301

back to its original resolution is to apply 'upconvolutional' layers with trainable weights. However, we found that simple tricubic interpolation of the latent space representation to estimate a
high-resolution output was faster and as accurate as using upconvolutional layers. Our architecture thus resembles the FCN
presented in ref. [42] for the upsampling process.

We allowed for the possibility of un-annotated neurons in each iCNN prediction because, when tracking the brain across behavior, some neurons can drop out of the field of view or be difficult to visualize at some time points. To allow the network to be able to ignore un-annotated neurons, we chose a cross entropy loss function which is able to mask out the channels corresponding to such neurons.

We implemented the iCNN with the properties described above in the 3D Compact Network (3DCN) architecture (Figs. 2 B, S1). The 3DCN exhibited improved accuracy, stability, and speed over U-Net for our tracking task (Table 1, Fig. 3 C). The 3DCN was efficiently trained on a desktop workstation.

Selecting images for targeted augmentation. We wanted to
 minimize the number of manually annotated GT volumes that
 the pipeline would need. Our strategy was to enlarge and di versify the set of GT annotations by automatically generating
 synthetic GT annotations.

For this, the algorithm selects target images, an optimally repre-326 sentative subset of all images from the recording, which we used 327 as templates for synthetic GT annotations. The images were se-328 lected to be different from the set of GT manual annotations and 329 different from each other. To accomplish this, we needed a dis-330 tance metric to estimate image similarity. We used a convolu-331 tional autoencoder [43] to create a low-dimensional latent space 332 representation of all recorded image volumes. We reduced the 333 autoencoder's latent space representation further to two dimen-334 sions using the UMAP [36] method. The distance between two 335 points in the UMAP plane is a measure of the overall similarity 336 between the corresponding brain images (Fig. 2 C). In this way, 337 the algorithm selects a set of target images that broadly sam-338 ples the points in the latent space representation (See Extended 339 Methods). 340

Creating synthetic GT annotations. One way to create addi-341 tional GT annotations for the selected target images would be 342 to use the iCNN to make coarse predictions for neurons in the 343 target images, and then perform proofreading and manual cor-344 rection. Instead, we aimed for a less laborious and fully auto-345 mated method by leveraging the GT images and their annota-346 tions. For each target image, the method selects in turn the most 347 similar GT image, and the iCNN makes coarse predictions for 348 neurons in the target image. Next, the goal is to deform the 349 GT image to resemble the target image. For key-point annota-350 tions, we fit a low-frequency deformation field that optimally 351 maps the key-points in the manually annotated GT image onto 352 the coarse predictions of key-point locations in the target image 353 (Fig. 2 D). (See 'Segmenting and tracking volumetric objects' 354 for the modifications of this step for 3D volume tracking.) We 355 implement this fitting by minimizing the mean L1 displacement 356 of key-points after deformation and by restricting the Fourier 357 modes of the deformation field to low frequencies. We used the 358

L1 loss because it is minimally sensitive to the outlying errors 359 made by the iCNN – the inaccurate coarse predictions of the 360 iCNN that assign key-points far from their actual locations. 361

Thus our deformation field, $\mathbf{D}(\vec{x})$, is

$$\mathbf{D}(\vec{x}) = A \operatorname{Re}\left[\sum_{\vec{k}_{i}}^{|\vec{k}_{i}| < k_{max}} \theta_{i} e^{-i\vec{k}_{i} \cdot \vec{x}}\right] \quad , \tag{1}$$

362

363

where $A, \vec{\theta}$ are the free fitting parameters. The loss function is

$$\mathscr{L}(\mathbf{D}) = \sum_{j=1}^{N_{\text{key}}} |\mathbf{D}(\vec{x}_j) - \vec{d}_j| + \lambda \int |\nabla \cdot \mathbf{D}| \, dV \quad , \qquad (2)$$

where N_{key} is the number of key–points and \vec{x}_j , \vec{d}_j are, respectively, the ground truth location of the *j*th key-point and the displacement vector from the latter to the corresponding 'coarsely' predicted location. The second factor is a regularization term to keep the divergence of the deformation field reasonable.

Finally, we used the fitted deformation field to generate both a 369 synthetic image and its corresponding annotation by applying 370 the field to both the GT image and its manual annotation. The 371 synthetic image that results from this deformation will resem-372 ble the target image while propagating the correct annotations 373 from the GT image. The resulting key-point locations represent 374 a reliable annotation of the synthetic image even when the orig-375 inal coarse prediction made by the iCNN is unreliable, albeit 376 with less resemblance to the target image. This synthetic image 377 and its key-point annotations can then be added to the GT im-378 ages and their manual annotations to create a larger training set 379 for a 'targeted augmentation CNN' (taCNN) that has better per-380 formance than the iCNN because it covers more postures. The 381 targeted augmentation CNN is then applied to all images in the 382 recording. 383

Evaluation of the targeted augmentation method. We evalu-384 ated our method for targeted augmentation by segmenting and 385 tracking key-points in brain-wide recordings of the male brain 386 in freely behaving C. elegans performing their mating ritual [12] 387 (Fig. 1 A). We obtained four recordings, each containing 1500-388 3000 image volumes (5-10 min per recording). In the previ-389 ous study of male mating behavior, all recordings were fully 390 manually annotated, which required 100-200 hours per record-391 ing. We asked whether our technique using the taCNN would 392 achieve comparable reliability as full manual annotation, but 393 with a smaller set of images that needed to be manually anno-394 tated to train the neural network. We developed our technique 395 using one brain-wide recording of the male brain. We evalu-396 ated the performance of our method by comparing the results 397 of taCNN predictions to the results of full manual annotation of 398 the three other 'held-out' recordings (Table 1, Fig. 3). 399

In practice, we found that we could set the number of target images from a brain-wide recording to N_{target} =80, which saturated the accuracy of the final taCNN predictions in tests of our method (Fig. S2). The utility of our technique is measured by how much it reduces the number of manually annotated GT images that are needed to achieve the same accuracy in segmenting neurons and tracking key-point locations as the iCNN. We found that the taCNN is able to identify nearly the same

	Masked				Not Masked			
	Strict		Close		Strict		Close	
	iCNN	taCNN	iCNN	taCNN	iCNN	taCNN	iCNN	taCNN
Found Fraction (FF)	68.0 ± 18.4	82.2 ± 10.2	76.2 ± 18.2	91.4 ± 8.4	71.5 ± 16.9	81.6 ± 8.3	80.3 ± 17.1	91.2 ± 5.2
Variance of FF	6.8 ± 7.0	$3.9\!\pm\!4.7$	7.8 ± 8.3	4.4 ± 5.3	7.8 ± 8.3	1.8 ± 1.0	8.9 ± 9.9	1.9 ± 1.1
Mistake Fraction	13.6 ± 6.1	12.3 ± 5.7	5.4 ± 3.3	$3.1\!\pm\!2.5$	14.6 ± 6.6	12.9 ± 5.4	$5.8\!\pm\!3.4$	3.3 ± 1.9
Miss Fraction	18.4 ± 18.4	5.5 ± 7.9	18.4 ± 18.4	5.5 ± 7.9	13.9 ± 17.9	5.5 ± 4.2	13.9 ± 17.9	5.5 ± 4.2
Extra Fraction	7.4 ± 3.5	9.9 ± 5.0	7.4 ± 3.5	9.9 ± 5.0	10.3 ± 5.3	10.2 ± 4.8	10.3 ± 5.3	10.2 ± 4.8

	Method						
	Sti	rict	Close				
	CPD	3DCN	CPD	3DCN			
Found Fraction (FF)	25.3 ± 12.4	82.2 ± 10.2	28.0 ± 12.9	91.4 ± 8.4			
Variance of FF	1.7 ± 0.5	$3.9\!\pm\!4.7$	1.8 ± 0.5	4.4 ± 5.3			
Mistake Fraction	34.5 ± 4.0	12.3 ± 5.7	31.8 ± 4.6	$3.1\!\pm\!2.5$			
Miss Fraction	40.2 ± 8.9	$5.5\!\pm\!7.9$	40.2 ± 8.9	5.5 ± 7.9			
Extra Fraction	3.7 ± 1.8	9.9 ± 5.0	3.7 ± 1.8	9.9 ± 5.0			

	Architecture						
	Str	rict	Close				
	U-Net	3DCN	U-Net	3DCN			
Found Fraction (FF)	48.7 ± 8.6	76.2 ± 2.0	56.4 ± 9.8	87.2 ± 1.9			
Variance of FF	8.6	2.0	9.8	1.9			
Mistake Fraction	38.6 ± 6.4	14.7 ± 0.9	30.9 ± 7.6	3.7 ± 0.8			
Miss Fraction	12.7 ± 4.5	9.1 ± 2.0	12.7 ± 4.5	$9.1\!\pm\!2.0$			
Extra Fraction	19.6 ± 1.1	16.1 ± 1.0	19.6 ± 1.1	16.1 ± 1.0			
Gradient Step Time	0.486 ± 0.001	0.237 ± 0.004	0.486 ± 0.001	0.237 ± 0.004			

Table 1. Benchmarks. All fractions are in percent (%). Tests were run on data that was hidden during method development. All runs were trained until saturation of the eIOU (see extended methods). All methods assume a training set size of 10. Each value represents the mean ± standard deviation. The standard deviation around the mean describes the variability when we sample from the GT training sets of the three different recordings. The value for "Variance of FF" represents the variation when we sample from different GT training sets for the same recording. Top table: Comparison of the iCNN to the taCNN. In the "Masked" case, a second worm that was in the FOV was masked out. The "Strict" case requires the key-point to be the closest point found from the actual neuron. "Close" condition is a simple distance threshold of four pixels between the actual and the predicted key-point. Middle table: We compare our method to the CPD method (Supplemental note 1). Bottom table: We compare our network architecture to U-Net [39]. We found that the 3DCN is more than twice as fast as well as more accurate and stable.

fraction of neurons across the three held-out brain recordings 408 when trained with only 5-25 manually annotated GT images as 409 the iCNN that is trained with 100 manually annotated GT im-410 ages (Fig. 3 A). This did not change when we varied the pixel-411 distance for the threshold determining whether a key-point was 412 called correctly; the taCNN always substantially outperformed 413 the iCNN (Fig. S3). Furthermore, when neurons where found, 414 they were identified more accurately with targeted augmenta-415 tion (Fig. 3 B). 416

In image volumes containing the male brain during mating be-417 havior, the hermaphrodite often entered the field of view, a de-418 coy that challenged both the iCNN and taCNN. The taCNN out-419 performed the iCNN in segmenting and tracking neurons in the 420 male brain whether or not we erased the decoy hermaphrodite 421 from image volumes. If we increased the threshold for a cor-422 rectly identified neuron - by requiring not only that it has to 423 be within a given pixel-distance of the actual neuron but also 424 that it is the closest predicted neuron - the taCNN still always 425 outperformed the iCNN (Fig. S4). 426

We compared segmentation and tracking by taCNN with another commonly applied method, Coherent Point Drift (CPD) point-set registration [44] (Supplemental note 1). Point-set registration requires fully segmented neurons in all image volumes, 430 whereas the taCNN performs both segmentation and tracking. 431 CPD tracks neurons by finding the optimal correspondence be-432 tween neuron locations between two images. We applied CPD 433 to our brain-wide imaging datasets of male mating behavior. We 434 created reference sets of different sizes so that CPD could match 435 neurons between a new image and the closest image in the ref-436 erence set. The larger the reference set, the better CPD should 437 work, analogous to increasing the size of the training set for the 438 taCNN. We found that CPD accurately tracked fewer than 20% 439 of neurons for similar reference set sizes where taCNN correctly 440 tracked >80% of neurons (Fig. S5). 441

Segmenting and tracking volumetric objects

We asked whether the taCNN could be used to segment and 443 track the 3D shapes of neurons, not just key-point locations. In 444 C. elegans and other animals, calcium dynamics is often differ-445 ent in different parts of a cell in functionally important ways 446 (e.g., soma vs. neurites). [26] These dynamics are missed when 447 recording calcium dynamics with nuclear-localized probes, and 448 require reconstruction of the spatial distribution of calcium dy-449 namics in different neuronal compartments. 450

442

We developed a transgenic strain to measure calcium dynam-



Figure 3. Evaluation of the performance of our pipeline for key-point tracking applied to three different recordings of 5-10 min duration (1500-3000 image volumes), each indicated in a different color (gold, orange, purple). A: A neuron was considered found if the predicted key-point was within 4 pixels of the ground-truth key-point. Solid lines indicate the performance of the augmented CNN, dashed lines indicate the performance of the initial CNN. B: Mean distance from the ground-truth for found neurons. Each recording is represented by the same color as in panel A. C: Comparison of 3DCN and U-Net. The 'closest' condition requires that the predicted key-point marking a neuron is the closest predicted key-point to the actual key-point.



Figure 4. Evaluation of the performance of our pipeline for 3D volume tracking, applied to three different recordings of 10 min duration (1715 image volumes), each represented by a different color (blue, green, yellow). A: For 3D volumes, we considered objects as found if the IOU between the predicted and the manual annotation exceeded 0.33, i.e., the intersection of two 3D objects was greater than the volume of the ground-truth or the predicted 3D volumes. Solid lines indicate the performance of the augmented CNN, dashed lines indicate the performance of the initial CNN. B: Intersection-over-Union (IoU) for 3D volumes. Each recording is represented by the same color as in panel A.

ics throughout cytosolic compartments in interneurons of the 452 hermaphrodite nerve ring that are responsible for chemosensory 453 processing, including the AIY interneuron and RIA interneuron 454 that primarily exhibit calcium dynamics in their neurites and not 455 their cell bodies. This strain (SJR15) expressed red fluorescent 456 proteins (wrmScarlet and mNeptune) and GCaMP6s in the cy-457 tosols of the AIA, AIY, AIZ, and RIA interneurons. This strain 458 also expressed red fluorescent proteins in the nuclei of the AIB, 459 RIB, and RIM interneurons, and GCaMP6s in the nuclei of all 460 neurons in the nervous system. 461

We modified some steps of the taCNN method to segment and track the shapes of neurons and nerve fibers in image volumes. First, we needed to generate GT annotations of neuronal shapes and structures. We began by using adaptive thresholding to identify contiguous fluorescently-labeled objects within each image volume. All objects are associated with a number of quantitative measures – e.g., overall size, aspect ratio, brightness – that can be used as identifiers. We quantified a large set 469 of these geometric features for all objects across all image vol-470 umes. We performed k-means clustering on the set of geomet-471 ric features, resulting in the automated clustering of the same 472 objects found in different image volumes. In effect, we used 473 k-means clustering as an elementary method for tracking ob-474 jects that had reasonable accuracy (approx. 60%, see Methods). 475 Finally, the volumetric structures that were automatically seg-476 mented by adaptive thresholding and tracked by k-means clus-477 tering were then manually proofread and corrected, a step that 478 was much faster than their full manual annotation. Thus, for 479 volumetric object tracking, we used machine learning already 480 at the manual annotation step. 481

As for key-point tracking with the taCNN, we augmented the training dataset for neuronal shapes by adding synthetic GT annotations for a set of target images that were created by deforming similar manually annotated GT images. Instead of fitting a



Figure 5. Calcium activity in the second-layer interneurons RIA, RIB, and RIM. "Neurite" indicates a segment of the neurite of RIA. A: Maximum intensity projection of a 3D image of a worm expressing nuclear-localized GCaMP6s pan-neuronally, nuclear-localized mNeptune in RIB, nuclear-localized wrmScarlet in RIM, and cytosolic GCaMP6s and mNeptune in RIA. B-D: Traces of GCaMP fluorescence divided by red fluorescence (R(t)), from which the baseline R_0 was subtracted and the difference was normalized by R_0 . Green bars indicate 2-nonanone pulses. Red bars indicated IAA pulses. B: Entrained calcium activity in all three neurons. C: Entrainment begins after the third odor pulses, indicated by the arrow. D: Arrows indicate when entrained activity stops and when activity restarts later.

low-frequency deformation field, we obtained better results by
applying a non-linear optical flow transformation [45, 46] from
each GT image to each target image.

Performance evaluations of the taCNN applied to volumetric 489 tracking of neuronal shapes and structures agree with evalua-490 tions of key-point segmentation and tracking (Fig. 4). Using the 491 taCNN led to substantial improvements in comparison to the 492 iCNN in the fraction of neurons found (Fig. 4 A). We further 493 compared the accuracy in the segmentation of individual neu-494 ronal objects by computing the intersection-over-union (IoU) 495 of all objects found. By this measure, the iCNN and taCNN 496 performed similarly, with the taCNN slightly outperforming 497 the iCNN for small training sets and vice versa for medium 498 and larger training sets (Fig. 4 B, note y-axis units). As with 499 key-point tracking, using the taCNN substantially reduced the 500 amount of manual annotation and proofreading. With targeted 501 augmentation, small training sets ($N_{\rm GT} = 5$) produced results 502 similar to training sets that were 3-5 times larger and that were 503 not enhanced by targeted augmentation. 504

505 Coupling of sensory information to interneuron activity

To apply our method to measurements that could not be eas-506 ily analyzed with previous methods, we recorded calcium ac-507 tivity in the second-layer interneurons RIA, RIB, and RIM. 508 Worms expressed nuclear localized GCaMP6s pan-neuronally 509 (Fig. 5 A). Additionally, RIB expressed nuclear-localized 510 mNeptune, RIM expressed nuclear-localized wrmScarlet, and 511 RIA expressed cytosolic GCaMP6s and mNeptune. The differ-512 ent choices of red fluorescent protein allowed us to distinguish 513 the nuclei of RIB and RIM neurons. Worms were placed in mi-514 crofluidic chips encompassing a structured arena, adapted from 515 ref. [47], and were pulsed with IAA or 2-nonanone medium 516 pulses of 20 sec duration and 1 min period. The worms were im-517 aged as described above, and activity was extracted after analy-518 sis with our pipeline. 519

Periodic stimuli are powerful tools for elucidating circuit be-520 havior. [33, 34] We observed a number of rich activity patterns 521 in freely behaving C. elegans. Neuronal activity in the second-522 layer interneurons, which are thought to be closely linked to 523 locomotion, could be entrained by the odor pulses (Fig. 5 B). 524 However, this varied not only from worm to worm but also 525 from time to time for the same worm. Some animals showed 526 full entrainment to the external odor pulses while others exhib-527 ited none. Interestingly, worms could switch entrained activity 528 on or off (Fig. 5 C, D). These observations suggest that long 529 recordings of single animal activity in multiple neurons con-530 tinue to reveal novel phenomena, highlighting the importance 531 of efficient 3D image analysis techniques. 532

Graphical user interface (GUI)

We created a python-based GUI for viewing and annotating 4D 534 recordings, launching the steps of the pipeline, including tar-535 geted augmentation and the neural network training, and for 536 viewing and proofreading the results of the predictions of the 537 538 different neural networks, iCNN and taCNN (Fig. 6). The user can leaf through z-stacks, view z-projections, and annotate by 539 placing or moving key-points as well as drawing 3D masks with 540 a cubic pencil or by local thresholding ('3D bucket fill'). A 541 'neuron bar' and a dashboard are designed to make it easier

for the user to spot incomplete annotations and navigate long recordings. 543



Figure 6. Targettrack GUI for viewing and annotating 4D recordings as well as applying the pipeline.

545

Discussion

Many laboratories now perform whole-brain or multi-neuron 546 imaging with single-cell or subcellular resolution in different 547 animals [5, 12, 37, 38]. The current bottleneck is data anal-548 ysis: converting large-scale recordings of image volumes over time into the segmentation and tracking of individual neuronal 550 activities throughout the brain. Manual annotation is the most 551 reliable way of analyzing brain-wide recordings. However, 552 manual annotation becomes unacceptably labor-intensive when 553 recordings become numerous, long, or encompass many neu-554 rons. Manually annotating even a few minutes of recording can 555 take hundreds of hours [12], slowing progress. 556

Machine learning is ideally suited to the pattern recognition 557 task of segmenting and tracking neurons. However, any method 558 must deal with substantial image-to-image variability due to op-559 tical changes, biological differences, and movement and defor-560 mation during animal behavior in multi-neuron recordings. To 561 be accurate, neural networks must be trained with representative 562 images and annotations that span the range of image variabil-563 ity [21]. As the diversity of images and the number of neurons 564 increase, the amount of training data that is required also in-565 creases, which represents a significant burden if training data is 566 produced by manual annotation. After automated segmentation 567 and tracking using a neural network, manual proofreading also 568 becomes a burden if the network has high error rates. Thus, 569 traditional machine learning techniques involve a trade-off be-570 tween the amount of manually annotated data that is used to 571 train a neural network and the amount of manual proofreading 572 needed to correct errors. 573

We present innovations that allow a CNN to both minimize the required amounts of manually annotated training data and of proofreading. We optimized the neural network architecture to reliably segment and track neurons within a rapidly moving and deforming *C. elegans* brain. Our method generates part of its own training data based on a small number of manually annotated images using targeted augmentation. By estimating the deformation of a brain volume in a target image based on a manually annotated image, we automatically create new synthetic
images with reliable annotations. When a CNN is trained with
a small number of manually annotated images along with diverse synthetic images and annotations, its reliability increases
substantially, reducing the amount of proofreading needed.

The automatic generation of synthetic training images has pre-587 588 viously been explored in the medical domain. For example, [48] focuses on segmenting cardiac, prostate, and pancreas images, 589 and generates synthetic training examples using GANs. In con-590 trast, [49] focuses on the segmentation and registration of brain 591 and knee MRIs. The method uses a neural network to learn 592 displacement fields between images. The closest approach to 593 ours may be that of ref. [50]. This method segments brain MRIs 594 with the help of synthetic training data generated by fitting free-595 form deformations between existing labeled and unlabeled im-596 ages. Overall, the task in the present work, i.e., tracking freely 597 moving worms, is quite different from segmenting MRI images, 598 among other reasons, because the 3D volumes within a record-599 ing are much more heterogeneous. This makes segmentation 600 and tracking of our recordings substantially more challenging. 601

Because of experiment-to-experiment variability, any image 602 analysis method will be more reliable when there is ground-603 truth training data specific to each experiment. Our pipeline 604 starts with a small number of manually annotated GT images for 605 each experiment and requires no extraneous information. The 606 pipeline effectively learns the brain-wide deformations that oc-607 cur in one individual experiment. It then automatically expands 608 the originally annotated GT images into a larger training dataset 609 using the learned deformations. Our method is ideally suited 610 for use by different laboratories and changing experiments, as it 611 flexibly adapts to the specific imaging conditions of each indi-612 vidual experiment. 613

Another noteworthy aspect of our targeted augmentation 614 pipeline is that it is germane to 3D images. While 3D images 615 may be perceived as merely more difficult to analyze than 2D 616 images because of their larger sizes, they also afford additional 617 opportunities for image analysis. In three dimensions, the worm 618 brain can in principle be mapped by deformation from one im-619 age to any other. This is not true for 2D projections of the brain, 620 where, for example, two crossing lines cannot be uncrossed by 621 deformation. Thus, there may be additional unexplored oppor-622 tunities for simplifying image processing tasks in 3D. 623

We have developed and applied our approach to particularly 624 challenging problems in C. elegans brain-wide imaging. For 625 example, during its unrestrained mating behavior, the poste-626 rior brain of the male nematode exhibits rapid and dramatic 627 movements and deformations as the male interacts with a 628 hermaphrodite, itself a visual object that distracts and chal-629 lenges the performance of the neural network that is focused 630 on the male. In practice, we have required 200 hours to fully 631 manually segment and track 76 neurons in the male tail in just 632 one 10 minute recording. Using our taCNN pipeline from end-633 to-end on the same dataset, we reduced the amount of man-634 ual effort to 65 hours - 5 hours to generate a small but ade-635 quate amount of manually annotated GT images and 60 hours 636 to comprehensively proofread all automatically segmented and 637 tracked images. The latter is generally difficult to reduce as 638 proofreading is still necessary even for simpler image analysis

problems with even lower error rates [51]. Thus, the speed-up of our method expands what is feasible for brain-wide imaging of small-animal models in neuroscience. 640

We expect the effectiveness of our pipeline to improve with 643 advances in imaging. Segmenting and tracking neurons will 644 become easier with better image quality and higher spatial 645 and temporal resolution, which will be possible with improve-646 ments in microscopy and fluorophores. Multi-color imaging ap-647 proaches will allow the taCNN to incorporate more information 648 that will facilitate its reliability. In this work, we have not used 649 any information beyond a single fluorescent channel to keep the 650 method as general as possible. One strength of the CNN frame-651 work is that it is straightforward to add additional types of in-652 formation such as additional image channels. 653

Methods

Initial coarse alignment of whole-brain images

To facilitate later steps in the segmentation and tracking
pipeline, the algorithm performs an initial coarse alignment of
all whole-brain images in each recording. There exist many
techniques for the coarse alignment of images, two of which we
present here.650
657658659659

654

655

Whole-brain recordings of neurons segmented and tracked as key-points 662

When tracking many neurons inside freely moving *C. elegans*, 663 our tests were applied to a set of manually annotated recordings 664 of the posterior nervous system of the male during mating be-665 havior [12]. We used these manual annotations to train a 2D 666 U-Net to solve the problem of coarse alignment of whole-brain 667 images. To perform coarse alignment, an effective algorithm 668 must (a) determine which pixels in a 3D image correspond to 669 the brain and (b) determine the brain's orientation. We created 670 a training dataset for the 2D U-Net by converting the compre-671 hensive manual annotations of segmented and tracked neurons 672 from previous work [12] into a simple map of neuron locations 673 within the brain distinguished by their coordinates along the 674 anterior-posterior axis. After training, the 2D U-Net was able 675 to identify neurons and estimate their coordinates when given a 676 new 3D brain image. Next, the algorithm computes the gradient 677 of these estimated coordinates, which represents the orientation 678 of the worm brain in each new image. Finally, this computed 679 orientation is used to perform an affine alignment of each im-680 age. We found that the 2D U-Net network, when trained with 681 recordings of 1-3 animals, was effective for identifying neurons 682 in images of other animals and thus is useful as a general tool 683 for coarse alignment. 684

Recordings of neurons segmented and tracked as 3D volumes

We used a different coarse alignment procedure to orient neu-686 rons represented as 3D shapes. The algorithm identifies a few 687 landmark neurons, that is, particularly bright neurons that are 688 visible in all 3D brain images and that can be automatically 689 detected using a high threshold on the fluorescence intensity. 690 We then compute the coarse alignment by performing the non-691 rigid Jian-Vemuri [52] point cloud registration algorithm. We 692 approximate the point-wise registration of landmarks by rota-693 tions and translations of the images in the x-y plane. 694

605 Ground-truth image selection

To train the 3DCN to perform segmentation and tracking of neu-696 rons, we need a diverse set of manually annotated ground truth 697 (GT) images. The GT images need to be diverse to account 698 for the different postures of moving animals across a particular 699 recording. We either select images at regular intervals through-700 out the recording or select an easy-to-annotate image sequence 701 where the animal moves substantially. Either method was suf-702 ficient, and thus we did not develop a computational method to 703 select GT images. In practice, the user may benefit from flex-704 ibility when choosing images for GT manual annotation. An 705 algorithm that prescribes the 3D images to annotate may not, 706 for example, be able to account for the varying subjective diffi-707 culty of annotating particular 3D images. 708

709 Ground-truth annotations

Once GT images are selected, they need to be accurately an-710 711 notated to serve as training datasets for the 3DCN. For freely moving animals that are to be tracked by key-points, we did not 712 automate the annotation of GT images. For example, non-rigid 713 point-set registration methods such as CPD [44] (Supplemen-714 tal note 1) were not sufficiently accurate to speed up manual 715 annotation. Thus, all GT images for segmenting and tracking 716 neurons in freely moving animals are obtained by manual an-717 notation. For semi-immobilized worms, however, CPD pointset 718 registration can generate rough annotations that can be proof-719 read in less time than full manual annotation. 720

For segmenting and tracking the 3D shapes of neurons, we de-721 veloped a semi-automated method for the annotation of GT im-722 ages. First, we segment the images, that is, identify the 3D 723 structures representing the nuclei, soma, and neurites of fluo-724 rescent neurons. In each recording, we enhance each 3D im-725 age by applying a Difference-of-Gaussians filter. We apply a 726 threshold to the enhanced image to keep a percentage of the 727 brightest pixels. We compute the Euclidean distance transform 728 of the thresholded image, which is then smoothed with another 729 Gaussian kernel. The local maxima of the smoothed distance 730 transform are used as seeds for a watershed algorithm, which 731 finds the connected volumes around each local maximum. We 732 applied a user-defined threshold to discard small local maxima 733 that are too close to bigger maxima. We merge volumes that 734 were overlapping or adjoining. We adjust the algorithm to re-735 move segmentation errors (by merging pieces of the same object 736 that were erroneously split into different volumes) without cre-737 ating merge errors (avoiding the erroneous merging of different 738 adjoining objects). To do this, we only merge when the con-739 tact areas are large, that is, when the overlapping/neighbouring 740 surface divided by the smaller volume (to the power of 2/3) is 741 greater than a user-defined threshold. Volumes that are too small 742 are excluded. 743

In addition to segmentation, we created complementary tools
for the manual creation and deletion of 3D annotations in the
GUI. This tool allows 3D volume masks to be drawn with a
cubic 'pencil' or created by local thresholding.

Once all 3D individual objects are segmented and identified by
 adaptive segmentation and manual annotation, each is represented by a vector of quantitative features. These features are
 volume, total fluorescence intensity, maximum intensity, vari-

ance of intensity, ratio of diameter to volume, and eigenvalues 752 of the moments of inertia matrix. We then apply K-means clus-753 tering to segregate and locate all 3D objects in the feature space. 754 When applied to different time points and different 3D images, 755 K-means clustering should assign the same objects to the same 756 locations in the feature space. We found that the accuracy of 757 this method was about 60% (the average number of correctly 758 tracked objects in 12 frames in one recording). Proofreading 759 and correcting the results of this elementary tracking yielded 760 the GT annotations. 761

762

775

776

777

778

779

780

787

788

789

790

799

Point-to-mask conversion

For key-point tracking, the method should predict a key-point 763 corresponding to the location of each tracked neuron. During 764 the training of the neural network, however, it did not suffice 765 to supply a single pixel to be predicted for each neuron. In-766 stead, for each annotation, we generate a mask in which all 767 pixels within a radius of 4 pixels from the ground-truth key-768 point are labeled as the neuron which is to be predicted. So, 769 the neural network is trained to predict a 4-pixel ball of pixels 770 to identify a neuron. When the neural network is then applied 771 to new images in the recording, we straightforwardly reduce 772 the set of predicted label pixels to a single key-point pixel (see 773 'Post-processing' below). 774

initial CNN

Once we have an initial set of GT images, we use them to train an initial CNN. The architecture of the CNN, which we call 3DCN, is illustrated in Figs. 2 C, S1. Several features make our CNN more accurate, as well as faster to train and apply than the popular U-Net [39].

First, we accounted for anisotropic resolution. In most 3D 781 light microscopy – whether confocal, two-photon, or lightsheet microscopy – the resolution in the xy directions is higher 783 than the resolution in the z-direction. We thus applied 2x2x1 784 down-sampling to compensate for the difference in xy- and zresolution. 786

We found that the 3DCN did not need to train the upsampling layer to generate the final predictions. Instead, we found that a simple tricubic interpolation was an effective and computationally efficient way to extract predictions.

We found that the trained network needed to account for long-791 range correlations to accurately identify neurons. Deformations 792 in distant parts of the animal contain information about the ani-793 mal's overall posture. To capture long-range information with-794 out large kernel sizes, we employed atrous convolutions in the 795 ASPP module [41]. In brief, this method enlarges the field of 796 view of the kernel by skipping over features that are adjacent to 797 features already captured. 798

Target set for augmentation

For targeted augmentation, the algorithm selects a set of diverse 800 3D images from the recording as templates for deforming the 801 GT images and annotations. This step requires image volumes 802 to be compared and their similarities to be quantified. To com-803 pute the similarity between any two images, all 3D images from 804 each recording are used to train a convolutional autoencoder. 805 Thus, the autoencoder maps each 3D image to a compressed 806 representation in the network's latent space. In principle, the 807

relevant information from the 3D images, e.g., noise, intensity 808 changes, worm body deformations, other objects in the field of 809 view, and the field of view, are captured in the latent space. We 810 used the L2 loss function so that the network focused on the 811 bright part of the image rather than trying to summarize the 812 background noise distribution, which is a majority of the pix-813 els. Because most deformations are in the x-y plane, it sufficed 814 to perform a maximum intensity projection in the z-direction 815 first, and apply a 2D autoencoder. After normalizing the latent 816 vectors from the autoencoder to zero-mean and unity-standard 817 deviation, the representation of each image was mapped onto a 818 plane using UMAP [36]. 819

Targeted augmentation 820

Targeted augmentation succeeds when the coarse annotations 821 by the initial CNN, while not accurate enough to be satisfac-822 tory as the final results of the pipeline, suffice to deform the GT 823 annotations to match the target images. We use the coarse anno-824 tations to match the nearest GT image to each target image by 825 computing the most effective and smoothed deformation field. 826 For key-points, we compute the deformation field with a wave-827 length cut-off in Fourier space by fitting the vectors pointing 828 from the neurons in the GT annotation to the neurons predicted 829 in the target image by the initial CNN. 830

For 3D objects, we compute the deformation field in multiple 831 steps. First, we perform a better coarse alignment of the GT 832 image and the target image. To do this, all objects in the GT 833 image and all predicted objects in the target image are approx-834 imated as clouds of points. After this, all objects are matched 835 with the Jian-Vemuri [52] point cloud registration method. The 836 Jian-Vemuri method [52] ignores the identities of the neurons 837 and matches the constellation of point clouds. It generates vec-838 tors representing the match of each point in the point cloud in 839 the GT image to a corresponding point in the point cloud in the 840 target image. Our algorithm then approximates the transforma-841 tion represented by these vectors as an overall translation and 842 rotation of the whole GT image. Subsequently, optical flow [45, 843 46] finely and non-rigidly registers the rotated and translated GT 844 with the target image. The deformation generates a new GT im-845 age and annotation which are often close to a valid annotation 846 of the target image. The deformation field, which represents the 847 translation-rotation and the optical flow registration, can then 848 also be applied to the annotations in the GT image. To prevent 849 objects in the annotation from being torn due to the optical flow 850 step, we post-process them by computing the nearest α -shape 851 (based on a radius of 5 pixels) that corresponds to each individ-852 ual object in each 2D mask slice [53]. 853

To explain the improvement of the annotations by the aug-854 mented CNN compared to the initial CNN, we speculate that 855 even when the deformations are small, the deformed images 856 force the neural network to have a more consistent represen-857 tation of the input images. If the results of the deformation 858 are imperfect, the deformed image and annotations are not ex-859 860 pected, in principle, to harm the neural network performance. When this happens, the augmentation just fails to increase the 861 diversity of realistic 3D brain postures. The validity of this idea 862 can be quantitatively assessed based on the data presented in 863 Figs. 3, 4.

Post-processing (not shown in Fig. 2)

For key-point tracking, the neural network predicts a set of 866 points as labeling a neuron, not a single key-point pixel. This 867 is because the training is performed with a 4-pixel ball around 868 each key-point annotation (see 'Point-to-mask conversion'). 869 Consequently, to generate key-points from the predicted labels, 870 we take the largest connected components of predicted pixels 871 for each neuron, and calculate its 'center of mass' where the 872 weight of each pixel is given by the fluorescence intensity in 873 the recorded image. The center of mass is then assigned as 874 the predicted key-point. This step also addresses the problem 875 that sometimes the neural network predicts disconnected pixels 876 to label an individual neuron; by taking the largest connected 877 component, stray pixel labels are ignored. 878

865

879

880

881

882

885

886

889

890

891

892

894

902

917

For 3D volumes, the neural network sometimes mislabels individual pixels that are part of one neuron as belonging to another neuron. Thus, we check all pixels that are part of the same connected component, and if there are pixels of one neuron touching or inside another neuron, they are merged with the larger 883 object.

Strains

For imaging interneurons, we used ADS1001 [Prgf-1:NLS-GCaMP6s] from ref. [12], which expresses nuclear-localized 887 GCaMP6s pan-neuronally. For recordings of secondlayer interneurons, sixIs9 was generated by integrating sjxEx9[Pglr-3::mNeptune::GCaMP6s; Psto-3::NLS-mNeptune; Pcex-1::NLS-wrmScarlet; Punc-122::dsRed; lin-15] into lin-15(n765) mutants. The integrant was outcrossed with N2 three times and crossed into SJR1 to make SJR16 (used for recordings 893 in Fig. 5 A-D).

Finally, for recording first- and second-layer interneurons si-895 multaneously, *sjxIs8* was generated by integrating *sjxEx8*[Pnpr-896 9::NLS-wrmScarlet; *Pttx-3::mNeptune::GCaMP6s*; Pceh-897 16::mNeptune::GCaMP6s; Pgcy-28d::mCherry::GCaMP6s; 898 lin-15] into lin-15(n765) mutants. The integrant was outcrossed 899 three times with N2 and crossed into ADS1001 to generate 900 SJR15 (used for recordings in Fig. 4). 901

Cultivation and microscopy

The animals were grown in a 20°C incubator on NGM plates 903 seeded with OP50 bacteria. At the stage of young adult, 904 they were transferred into a microfluidic polydimethylsiloxane 905 (PDMS) arena for recording. The microfluidic chip was cus-906 tomized based on the design presented in ref. [47]. 907

The recording was performed using a spinning disc confocal 908 microscope (Nikon Eclipse Ti2 and Yokogawa CSUX1FW) and 909 two Andor Zyla 4.2MP Plus cameras, one recording the green 910 and the other recording the red channel. High-resolution images 911 were collected through a 40× Nikon Plan Fluor Oil DIC N.A. 912 1.30 objective. Green (GCaMP6s) and red (mCherry, mNep-913 tune, wrmScarlet) channel 3D volumetric stacks were obtained 914 with an exposure time of 10 ms at approximately 3 volumes per 915 second. 916

Extraction of calcium activity

To extract the calcium activity of each neuron, we identified the 918 30% brightest red (mCherry, mNeptune, wrmScarlet) channel 919 pixels in each neuronal volume. We then computed for these 920

pixels the ratio (R(t)) of mean green (GCaMP6s) intensity to 921 mean red intensity. Furthermore, to exclude the effects of out-922 liers resulting from the poor annotation of frames, we dropped 923 the second lowest and second highest percentile neuronal activ-924 ities and smoothed the remaining recording tracks using a 1-D 925 Gaussian filter with a standard deviation of 3. 926

Finally, the activity of each neuron at time t was computed using:

$$\Delta R(t)/R_0 = (R(t) - R_0)/R_0, \qquad (3)$$

 R_0 is the lowest 1st percentile of R(t) in the recording. 927

CODE AND DATA AVAILABILITY 928

4D datasets of pan-neuronal nuclear marked The code and sample 929 multi-neuron marked are available worms and cytosolic worms 930 at https://github.com/lpbsscientist/targettrack . 931

ACKNOWLEDGEMENTS 932

AD, KK, MBK, SJR were supported by the École polytechnique fédérale de Lausanne 933 (EPFL), the Helmut-Horten Foundation, the Swiss Data Science Center grant C20-12, 934 and an EPFL Interdisciplinary Seed Fund. We thank Nicholas Greensmith for help 935 developing the GUI; Guillaume Obozinski for feedback; Albert Lin for help constructing 936 strains; Matthieu Schmidt and Alice Gross for collecting and analyzing data. 937

AUTHOR CONTRIBUTIONS 938

ADTS, CFP, KK, SJR conceived the project; CFP, KK, MBK, VS collected the data; CFP 939 prepared the data, developed the neural network, suggested targeted augmentation, 940 and implemented the method; MBK and CLJ adapted the method for 3D volumes; CFP 941 MBK ran the evaluations; CFP, MBK, AD developed the GUI; CFP, ADTS, SJR, MBK, 942 943 CLJ wrote the manuscript; ADTS, SJR initiated and supervised the project.

COMPETING INTERESTS 944

The authors declare having no competing interests 945

References 946

- [1] C. Dupre and R. Yuste, "Non-overlapping neural net-947 works in Hydra vulgaris," Current Biology, vol. 27, 8 948 2017. 949
- S. Kato et al., "Global brain dynamics embed the mo-[2] 950 tor command sequence of Caenorhabditis elegans," Cell, 951 vol. 163, no. 3, pp. 656-669, 2015. 952
- W. C. Lemon et al., "Whole-central nervous system func-[3] 953 tional imaging in larval Drosophila," Nature Communi-954 cations, vol. 6, no. May, 2015. 955
- [4] K. Mann, C. L. Gallen, and T. R. Clandinin, "Whole-956 brain calcium imaging reveals an intrinsic functional net-957 work in Drosophila," Current Biology, vol. 27, no. 15, 958 pp. 2389-2396, 2017. 959
- V. Venkatachalam et al., "Pan-neuronal imaging in roam-[5] 960 ing Caenorhabditis elegans.," Proceedings of the Na-961 tional Academy of Sciences of the United States of Amer-962 ica, vol. 113, no. 8, pp. 1082-8, 2016. 963
- T. Schrödel, R. Prevedel, K. Aumayr, M. Zimmer, and [6] 964 A. Vaziri, "Brain-wide 3D imaging of neuronal activity 965 in Caenorhabditis elegans with sculpted light," Nature 966 Methods, vol. 10, no. 10, pp. 1013-1020, 2013. 967
- J. P. Nguyen et al., "Whole-brain calcium imaging with 968 [7] cellular resolution in freely behaving Caenorhabditis el-969 egans," Proceedings of the National Academy of Sciences 970 of the United States of America, vol. 113, no. 8, pp. 1074-971 81, 2016. 972

- R. Prevedel et al., "Simultaneous whole-animal 3D imag-[8] 973 ing of neuronal activity using light-field microscopy," 974 Nature Methods, vol. 11, no. 7, pp. 727-730, 2014. 975
- S. Abrahamsson et al., "Multifocus microscopy with pre-[9] 976 cise color multi-phase diffractive optics applied in func-977 tional neuronal imaging," Biomedical Optics Express, 978 vol. 7, 3 2016. 979
- V. Voleti et al., "Real-time volumetric microscopy of in [10] vivo dynamics and large-scale samples with SCAPE 2.0," Nature Methods, vol. 16, no. 10, pp. 1054-1062, 2019.

980

981

982

990

991

992

993

994

995

996

998

- K. M. Hallinen et al., "Decoding locomotion from pop-[11] 983 ulation neural activity in moving C. elegans," bioRxiv, 984 2021. eprint: https://www.biorxiv.org/content/early/ 985 2021/01/15/445643.full.pdf. 986
- V. Susoy et al., "Natural sensory context drives di-[12] 987 verse brain-wide activity during C. elegans mating," Cell, 988 vol. 184, no. 20, pp. 5122-5137, 2021.
- J. C. Marques, M. Li, D. Schaak, D. N. Robson, and J. M. [13] Li, "Internal state dynamics shape brainwide activity and foraging behaviour," Nature, vol. 577, no. 7789, pp. 239-243, 2020.
- Y. Toyoshima et al., "Accurate automatic detection of [14] densely distributed cell nuclei in 3D space," PLoS computational biology, vol. 12, no. 6, e1004970, 2016.
- J. Ma and A. Yuille, "Nonrigid point set registration by [15] 997 preserving global and local structures," IEEE Transactions on Image Processing, vol. 25, no. 1, pp. 53-62, 999 2016. 1000
- [16] J. P. Nguyen, A. N. Linder, G. S. Plummer, J. W. Shae-1001 vitz, and A. M. Leifer, "Automatically tracking neurons 1002 in a moving and deforming brain," PLoS Computational 1003 Biology, vol. 13, no. 5, 2017. 1004
- S. Chaudhary, S. A. Lee, Y. Li, D. S. Patel, and H. Lu, [17] 1005 "Automated annotation of cell identities in dense cellular 1006 images," en, bioRxiv, p. 2020.03.10.986356, 2020. 1007
- D. Witvliet et al., "Connectomes across development re-[18] 1008 veal principles of brain maturation in C. elegans," Na-1009 ture, no. May 2020, 2021. 1010
- J. G. White, E. Southgate, J. N. Thomson, and S. Bren-[19] 1011 ner, "The structure of the nervous system of the nematode 1012 Caenorhabditis elegans," Philosophical Transactions of 1013 the Royal Society of London B, vol. 314, pp. 1–340, 1986. 1014
- T. Lagache, A. Hanson, A. Fairhall, and R. Yuste, "Ro-[20] 1015 bust single neuron tracking of calcium imaging in behav-1016 ing hydra," bioRxiv, pp. 1-30, 2020. 1017
- E. Moen, D. Bannon, T. Kudo, W. Graf, M. Covert, and [21] 1018 D. Van Valen, "Deep learning for cellular image anal-1019 ysis," Nature Methods, vol. 16, no. 12, pp. 1233-1246, 1020 2019. 1021
- [22] C. Wen et al., "3DeeCellTracker, a deep learning-based 1022 pipeline for segmenting and tracking cells in 3D time 1023 lapse images," eLife, vol. 10, no. 1, 2021. 1024

- X. Yu, M. S. Creamer, F. Randi, A. K. Sharma, S. W. Linderman, and A. M. Leifer, "Fast deep learning correspondence for neuron tracking and identification in *C. elegans* using synthetic training," *arXiv*, 2021. eprint: 2101.
 08211 (q-bio.QM).
- I. M. Gray, J. J. Hill, and C. I. Bargmann, "A circuit for navigation in *Caenorhabditis elegans*," *Proc. Natl. Acad. Sci. U. S. A.*, vol. 102, no. 9, pp. 3184–3191, 2005. eprint: https://www.pnas.org/content/102/9/3184.full.pdf.
- H. S. Kaplan, A. L. Nichols, and M. Zimmer, "Sensorimotor integration in *Caenorhabditis elegans*: A reappraisal towards dynamic and distributed computations," *Phil. Trans. R. Soc. B.*, vol. 373, no. 1758, p. 20170371, 2018. eprint: https://royalsocietypublishing.org/doi/pdf/ 10.1098/rstb.2017.0371.
- M. Hendricks, "Compartmentalized calcium dynamics in a *C. elegans* interneuron encode head movement," *Nature*, vol. 487, no. 7405, pp. 99–103, 2012.
- [27] M. H. Ouellette, M. J. Desrochers, I. Gheta, R. Ramos, and M. Hendricks, "A gate-and-switch model for head orientation behaviors in *C. elegans*," *bioRxiv*, vol. 2014, no. December, pp. 1–13, 2018.
- [28] Z. Li, J. Liu, M. Zheng, and X. S. Xu, "Encoding of both analog- and digital-like behavioral outputs by one *C. elegans* interneuron," *Cell*, vol. 159, no. 4, pp. 751–765, 2014.
- [29] H. S. Kaplan, O. Salazar Thula, N. Khoss, and M. Zimmer, "Nested neuronal dynamics orchestrate a behavioral hierarchy across timescales," *Neuron*, vol. 105, no. 3, pp. 562–576, 2020.
- [30] Y. Wang *et al.*, "Flexible motor sequence generation during stereotyped escape responses," *eLife*, vol. 9, M. Zimmer and P. Sengupta, Eds., e56942, 2020.
- [31] A. Sordillo and C. I. Bargmann, "Behavioral control by depolarized and hyperpolarized states of an integrating neuron," *eLife*, vol. 10, Y. Iino, R. L. Calabrese, and Y. Iino, Eds., e67723, 2021.
- [32] S. J. Rahi, K. Pecani, A. Ondracka, C. Oikonomou, and
 F. R. Cross, "The CDK-APC/C oscillator predominantly
 entrains periodic cell-cycle transcription," *Cell*, vol. 165,
 no. 2, pp. 475–487, 2016.
- [33] S. J. Rahi *et al.*, "Oscillatory stimuli differentiate adapting circuit topologies," *Nature Methods*, vol. 14, pp. 1010–1016, 2017.
- [34] E. V. Nikolaev, S. J. Rahi, and E. D. Sontag, "Subharmonics and chaos in simple periodically forced biomolecular models," *Biophys. J.*, vol. 114, no. 5, pp. 1232–1240, 2018.
- [35] N. Ji *et al.*, "Corollary discharge promotes a sustained motor state in a neural circuit for navigation," eng, *eLife*, vol. 10, 2021.
- [36] L. McInnes, J. Healy, and J. Melville, "Umap: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.

- [37] M. B. Bouchard *et al.*, "Swept confocally-aligned planar excitation (scape) microscopy for high-speed volumetric imaging of behaving organisms," *Nature photonics*, vol. 9, no. 2, pp. 113–119, 2015.
- [38] V. Voleti *et al.*, "Real-time volumetric microscopy of in vivo dynamics and large-scale samples with scape 2.0," *Nature methods*, vol. 16, no. 10, pp. 1054–1062, 2019. 1085
- [39] O. Ronneberger, P. Fischer, and T. Brox, *U-net: Convolutional networks for biomedical image segmentation*, 1087 2015. arXiv: 1505.04597 [cs.CV].
- [40] G. Bradski, "The OpenCV Library," Dr. Dobb's Journal 1089 of Software Tools, 2000.
- [41] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, *Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs*, 2017. arXiv: 1606.00915 [cs.CV].
- [42] J. Long, E. Shelhamer, and T. Darrell, *Fully convolutional networks for semantic segmentation*, 2015. arXiv: 1096 1411.4038 [cs.CV].
- [43] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, 1098 "Stacked convolutional auto-encoders for hierarchical 1099 feature extraction," in *International conference on artificial neural networks*, Springer, 2011, pp. 52–59. 1101
- [44] A. Myronenko and X. Song, "Point set registration: Coherent point drift," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 12, pp. 2262–2275, 2010.
- [45] L. Alvarez, J. Sánchez, and J. Weickert, "A scale-space 1105 approach to nonlocal optical flow calculations," in *Scale-Space Theories in Computer Vision*, M. Nielsen, P. Johansen, O. F. Olsen, and J. Weickert, Eds., Berlin, Heidelberg: Springer Berlin Heidelberg, 1999, pp. 235–246. 1109
- [46] C. Zach, T. Pock, and H. Bischof, "A duality based approach for realtime tv-11 optical flow," in *Pattern Recognition*, F. A. Hamprecht, C. Schnörr, and B. Jähne, Eds., Berlin, Heidelberg: Springer Berlin Heidelberg, 2007, pp. 214–223.
- [47] D. R. Albrecht and C. I. Bargmann, "High-content behavioral analysis of *Caenorhabditis elegans* in precise spatiotemporal chemical environments," *Nat. Methods*, 1117 vol. 8, no. 7, pp. 599–605, 2011.
- [48] K. Chaitanya *et al.*, "Semi-supervised task-driven data understand augmentation for medical image segmentation," *Medical Image Anal.*, vol. 68, p. 101 934, 2021.
- [49] Z. Xu and M. Niethammer, "Deepatlas: Joint semisupervised learning of image registration and segmentation," in *International Conference on Medical Image Computing and Computer Assisted Intervention*, D. Shen *et al.*, Eds., ser. Lecture Notes in Computer Science, vol. 11765, Springer, 2019, pp. 420–429.
- [50] J. Nalepa *et al.*, "Data augmentation via image registration," in *IEEE International Conference on Image Processing*, IEEE, 2019, pp. 4250–4254. 1128
- [51] N. Dietler *et al.*, "A convolutional neural network segments yeast microscopy images with high accuracy," *Nat. Commun.*, vol. 11, no. 1, p. 5723, 2020.

- 1134 [52] B. Jian and B. C. Vemuri, "Robust point set registration
- using Gaussian mixture models," *IEEE T. Pattern. Anal.*,
 vol. 33, no. 8, pp. 1633–1645, 2011.
- II37 [53] H. Edelsbrunner, D. G. Kirkpatrick, and R. Seidel, "On
 II38 the shape of a set of points in the plane," *IEEE Trans. Inf.*II39 *Theory*, vol. 29, no. 4, pp. 551–558, 1983.

Supplemental figures



Figure S1. Illustration of the neural network architecture and of the atrous convolutions.



Figure S2. Fraction of objects found for one recording (blue plot in Fig. 4 A) with $N_{GT} = 5$ (stars), $N_{GT} = 15$ (triangles), or $N_{GT} = 25$ (circles) ground-truth annotations as a function of the target set size N_{target} .



Figure S3. Fraction of key-points found depending on the pixel threshold for $N_{GT} = 5$ (left) and $N_{GT} = 15$ (right).



Figure S4. (top left) Fraction of key-points successfully found when the strict criterion is applied and the hermaphrodite is masked out (top right). Fraction of key-points successfully found when the strict criterion is applied and the hermaphrodite is not masked out (bottom left). Fraction of key-points successfully found when the strict criterion is not applied and the hermaphrodite is masked out (bottom right). Fraction of key-points successfully found when the strict criterion is not applied and the hermaphrodite is not masked out.



Figure S5. Fraction found for the CPD method.

¹¹⁴¹ Supplemental note 1: Details of the CPD method

We use CPD (Coherent Point Drift) together with a nearest neighbor matching scheme as a simple tracking method for comparison. We assume a perfectly solved segmentation problem using the ground truth pointset from manual annotation. The algorithm is as following:

- 1145 1. We begin with N annotated pointsets, i.e., points are provided with labels.
- 1146 2. For each non-annotated pointset:
- (a) The nearest annotated pointset among the N pointsets is found.
- (b) The nearest annotated pointset found is deformed into the current pointset using CPD.
- (c) The label of each point is assigned to be the label of the nearest deformed annotated pointset.