# Local Shape Descriptors for Neuron Segmentation

Arlo Sheridan<sup>1,2</sup>, Tri Nguyen<sup>3</sup>, Diptodip Deb<sup>1</sup>, Wei-Chung Allen Lee<sup>4</sup>, Stephan Saalfeld<sup>1</sup>, Srini Turaga<sup>1</sup>, Uri Manor<sup>2</sup>, Jan Funke<sup>1</sup>

1: HHMI Janelia, Ashburn, VA, USA

2: Waitt Advanced Biophotonics Center, Salk Institute for Biological Studies, La Jolla, CA, USA3: Department of Neurobiology, Harvard Medical School, Boston, MA, USA

4: F.M. Kirby Neurobiology Center, Boston Children's Hospital, Harvard Medical School, Boston, MA, USA

Abstract-We present a simple, yet effective, auxiliary learning task for the problem of neuron segmentation in electron microscopy volumes. The auxiliary task consists of the prediction of Local Shape Descriptors (LSDs), which we combine with conventional voxel-wise direct neighbor affinities for neuron boundary detection. The shape descriptors are designed to capture local statistics about the neuron to be segmented, such as diameter, elongation, and direction. On a large study comparing several existing methods across various specimen, imaging techniques, and resolutions, we find that auxiliary learning of LSDs consistently increases segmentation accuracy of affinity-based methods over a range of metrics. Furthermore, the addition of LSDs promotes affinitybased segmentation methods to be on par with the current state of the art for neuron segmentation (Flood-Filling Networks, FFN), while being two orders of magnitudes more efficient—a critical requirement for the processing of future petabyte-sized datasets. Implementations of the new auxiliary learning task, network architectures, training, prediction, and evaluation code, as well as the datasets used in this study are publicly available as a benchmark for future method contributions

#### 1 Introduction

The goal of connectomics is the reconstruction and interpretation of neural circuits at synaptic resolution. These wiring diagrams provide insight into the inner mechanisms underlying behavior and help drive future theoretical experiments (Schneider-Mizell et al., 2016; Motta et al., 2019; Bates et al., 2020; Hulse et al., 2020). Additionally, the generation of connectomes has proven to complement existing techniques such as calcium imaging and electrophysiology where the resolution is often not sufficient to parse the circuitry in detail (Schlegel et al., 2016; Turner-Evans et al., 2020).

Currently, only electron microscopy (EM) allows imaging of neural tissue at a resolution sufficient to resolve individual synapses and fine neural processes. Two popular methods for imaging these volumes are serial block-face scanning EM (SBF-SEM) and focused ion beam scanning EM (FIB-SEM). While the former technique is faster and generates high lateral resolution, it results in lower axial resolution due to section slicing. The latter method produces isotropic resolution by etching the face of the volume with a focused ion beam before imaging. However, this method is slower than serial section approaches. Briggman and Bock (2012) provide a thorough overview of these imaging approaches and others, including ssTEM and ATUM-SEM. All methods have been used to generate invaluable datasets for the

connectomics community (Lee et al., 2016; Zheng et al., 2018; Dorkenwald et al., 2019; Schneider-Mizell et al., 2020; Turner et al., 2020; Scheffer et al., 2020; Phelps et al., 2021).

Depending on the specimen and the circuit of interest, current EM acquisitions produce datasets ranging from several hundred terabytes to petabytes. For instance, the raw data of a full adult fruit fly brain (FAFB) comprises ~50 teravoxels of neuropil (Heinrich et al., 2018). Even sub-volumes taken from vertebrate brains, which do not contain brain-spanning circuits, result in massive amounts of data: The authors of Kornfeld et al. (2017) imaged a region from a zebrafinch brain containing ~ $10^6 \mu m^3$  (~663 gigavoxels) of raw data. More recently, a larger volume of mouse visual cortex has been imaged, comprising ~ $3 \times 10^6 \mu m^3$  (~6,614 gigavoxels) (Dorkenwald et al., 2019; Turner et al., 2020; Schneider-Mizell et al., 2020; Yin et al., 2020). In order to reconstruct circuits in a full mouse brain, however, the authors of Abbott et al. (2020) argue that it will require an acquisition of around one exabyte of raw data (one million terabytes).

With datasets of this magnitude, purely manual reconstruction of connectomes is infeasible. As an example, manual tracing in the Drosophila dataset, FAFB, reaches speeds of about 11 hours per cell (Zheng et al., 2018). Even the small brain of a Drosophila contains an estimated 100,000 neurons, which would require  $\sim 125$  years of manual effort to trace each neuron to completion.

Consequently, automatic methods for the reconstruction of neurons and identification of synapses have been developed. Over the past decade, methods targeting relatively small volumes have pioneered the reconstruction of neurons (Turaga et al., 2010; Lee et al., 2017), and synapses (Kreshuk et al., 2015; Buhmann et al., 2018). More recently, these efforts have been improved to tackle the challenges of large datasets for neurons (Januszewski et al., 2018; Funke et al., 2019; Dorkenwald et al., 2019; Li et al., 2019), synaptic clefts (Heinrich et al., 2018), and synaptic partners (Huang et al., 2018; Buhmann et al., 2020). With the help of an automatic neuron segmentation method, neuron tracing times decreased by a factor of 5.4 to 11.6 minutes per µm path length (Li et al., 2019), effectively trading compute time for human tracing time.

However, given the daunting sizes of current and future EM datasets, limits on available compute time become a concern. Future algorithms do not only need to be more accurate to further decrease manual tracing time, but also computationally more

efficient to be able to process large datasets in the first place. Consider the computational time required by the current state of the art, FFN: Assuming linear scalability and the availability of 1000 contemporary GPUs (or equivalent hardware), the processing of a complete mouse brain would take about 226 years. This example alone goes to show that the objective for future method development should be the minimization of the total time spent to obtain a connectome, including computation and manual tracing. Therefore, automatic methods for connectomics need to be fast, scalable (*i.e.*, trivially parallelizable), and accurate.

# 1.1 Neuron Segmentation Methods

Neuron reconstruction is an instance segmentation problem. Unlike semantic segmentation, in which the goal is to assign every voxel to a specific class, instance segmentation assigns all voxels belonging to the same object a unique label. Since those labels can not be predicted directly, alternative local representations are sought, which permit extraction of globally unique labels in a subsequent processing step.

The most straight forward local representation is to label pixels as either foreground or background, and then perform a connected component analysis limited to foreground pixels to extract unique objects. However, in the case of 3D neuron segmentation this approach often fails to distinguish voxels in finer neurites, where the axial resolution of the data is lower (Ciresan et al., 2012). To deal with those situations, several methods have centered around the prediction of affinities (*i.e.*, the labeling of edges between neighboring voxels as "connected" or "cut"), rather than labeling the voxels themselves (Turaga et al., 2010; Funke et al., 2019).

Affinities effectively increase the resolution of the prediction, but otherwise inherit the advantages and disadvantages of voxel-wise boundary labeling: Both can be computed locally during training and inference, which allows for trivial parallelization. However, segmentations extracted from those predictions are sensitive to small errors: A few incorrectly assigned voxels (or edges between voxels) can label a boundary as foreground, resulting in two segments becoming falsely merged during post-processing.

Those so-called merge errors are the notorious failure modes of neuron segmentation methods. Merge errors have generally been considered worse than split errors, since they have the potential to propagate throughout a dataset. Even if just two neurons are merged together, the resulting segmentation can be difficult to resolve for proofreaders, if the neurons in question have several contact sites. This has been a particular concern for the first generation of proof-reading tools, which did not provide algorithmic help to split wrongly merged objects. In those situations, the origin point of a false merge would first need to be found before the objects can be separated; a task akin to searching for a needle in a haystack.

To avoid those merge errors, Turaga et al. (2009) introduced Malis, a loss function that penalizes topological errors by minimizing the Rand Index. The Rand Index naturally favors split errors over merge errors and thus helped to bias boundary predictions to split instead of merge in ambiguous situations. Funke et al. (2019) expanded this method by constraining the loss to a positive and negative pass, and by providing a maximal spanning tree formulation of the loss, which allows for a quasi-linear and exact computation of the loss during training.

More recent methods do not explicitly focus on merge errors, which is a possible consequence of improved proofreading tools

that allow users to separate objects with just a few interactions. Lee et al. (2017) found that using an increased affinity neighborhood acts as an auxiliary learning objective to improve direct neighbor affinities. This work demonstrates that auxiliary learning helps to make better use of local context in the receptive field of the neural network. The nature of this auxiliary learning approach is similar to the LSDs proposed here.

All affinity-based methods have in common that they need a subsequent agglomeration step to produce a final segmentation. Methods such as watershed variants (Wolf et al., 2018) and constrained agglomeration (Beier et al., 2017) successfully demonstrated an increase in robustness of the resulting segmentation to small errors in the predicted affinities.

Not all neuron segmentation methods are based on boundary predictions. The most notable exception are Flood-Filling Networks (FFN), the current state of the art in terms of segmentation quality (Januszewski et al., 2018). FFN eliminated the need for a multi-step segmentation process by using a recurrent convolutional neural network to fill neurons iteratively in an end-to-end fashion. Given seed points within neurons, the algorithm predicts which voxels belong to the same object as the seeds. This approach has been proven to be successful on very large volumes, although it is computationally more expensive than its affinity-based counterparts.

More recently, another promising alternative to boundary prediction has been proposed by Lee et al. (2019), which uses metric learning to produce dense voxel embeddings. The embeddings of voxels that belong to the same object are encouraged to be close in embedding space, while the embeddings of voxels of different objects are pushed away from each other. Clustering of the embeddings then reveals a segmentation. Since object similarity or dissimilarity can only be discerned locally, the method is applied in a block-wise fashion and the segmentations of neighboring blocks are stitched together to process a whole volume.

#### 1.2 Contributions

We introduce LSDs as an auxiliary learning task for affinity predictions and demonstrate that segmentation results are competitive with the current state of the art, albeit two orders of magnitude more efficient to compute. LSDs are 10-dimensional vectors, computed for each voxel, which encode local object properties. We engineered LSDs to describe features that could be leveraged to improve boundary detection. Specifically, they consist of three parts: the local size of the object (1D), the offset to the local center of mass (3D), and the local directionality (6D) (described in detail in Section 2.1).

We conducted a large comparative study of recent neuron segmentation algorithms. Specifically, we evaluated four of the aforementioned methods against three LSD-based methods on the following datasets:

- ZEBRAFINCH: A region consisting of  ${\sim}10^6 \mu m^3$  ( ${\sim}663$  gigavoxels) of songbird neural tissue, imaged using serial blockface EM at 9x9x20 nanometer (xyz) resolution (Januszewski et al., 2018). 0.02% of the full dataset was used to train networks (dense segmentations), 12 manually traced neuron skeletons (13.5 mm) were used for validation and 50 skeletons (97 mm) were used for testing.
- **HEMI-BRAIN:** Three volumes containing  $\sim 1650 \mu m^3$ ,  $\sim 4750 \mu m^3$ , and  $\sim 10360 \mu m^3$  ( $\sim 33$  total gigavoxels) of raw

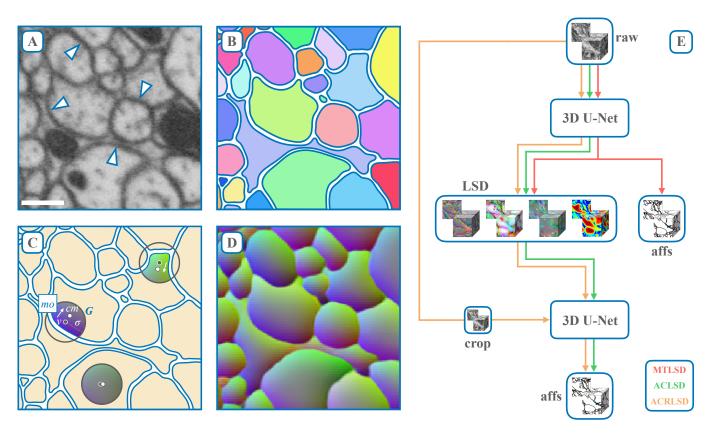


Figure 1: Local shape descriptor and network architecture overview. **A.** EM data imaged with FIB-SEM at 8 nm isotropic resolution (Hemi-brain dataset). Arrows point to example individual neuron plasma membranes. Dark blobs are mitochondria. Scale bar = 300 nm. **B.** Label colors correspond to unique neurons. **C.** LSD mean offset schematic. A Gaussian (G) with fixed sigma ( $\sigma$ ) is centered at voxel (v). The Gaussian is then intersected with the underlying label (colored region) and the center of mass of the intersection (cm) is computed. The mean offset (mo) between the given voxel and center of mass is calulated (among several other statistics, see Fig. 2 for a visualization), resulting in the first three components of the LSD for voxel (v). **D.** Predicted mean offset component of LSDs (Lsp[0:3]) for all voxels. A smooth gradient is maintained within objects while sharp contrasts are observed across boundaries. 3D vectors are RGB color encoded. **E.** Network architectures used. The 10D LSD embedding is used as an auxiliary learning task for improving affinities. In a multi-task approach (MTLsp), LSDs and affinities are directly learnt. In an auto-context approach, the predicted LSDs are used as input to a second network in order to generate affinities both without raw data (AcLsp), and with raw data (AcRLsp). For a detailed network architecture visualization, see Fig. 16.

data, cropped from the  $\sim$ 26 teravoxel dataset generated by Scheffer et al. (2020). This volume was taken from the central brain of a *Drosophila* and imaged with FIB-SEM at 8 nanometer isotropic resolution. 0.002% of the data was used for training (dense segmentation), and 0.06% for testing using a whitelist of proofread neurons.

• Fib-25: A ~1.8 × 10<sup>5</sup> µm<sup>3</sup> (~346 gigavoxels) volume from the *Drosophila* visual system was imaged with FIB-SEM at 8 nm isotropic resolution (Takemura et al., 2015). 0.09% of the data was used for training (dense segmentation). Testing was restricted to a 13.8 gigavoxel region using a whitelist of proofread neurons.

For each dataset, we compare LSD-based methods against three previous affinity-based methods: (1) direct neighbor and (2) long-range affinities with mean squared error (MSE) loss, and (3) direct neighbor affinities with Malis loss. Each affinity-based method (including our LSD methods) was trained in the same way and uses the same network architecture (where possible). We used the same segmentation extraction method (from Funke et al. (2019)) to convert the predicted affinities into segmentations.

We further include a comparison against FFN segmentations, which were made available to us by the authors of Januszewski et al. (2018).

We make the training scripts and datasets used in this study publicly available in a central repository<sup>1</sup>, in the hope that non-affinity-based methods that we did not cover in this study (like the recent deep metric learning proposed in Lee et al. (2019)) can be compared in a similar manner.

We compare the aforementioned methods against three different architectures that use LSDs as an auxiliary loss: a simple multitask approach (MTLSD) and two auto-context approaches (AcLsD and AcRLSD). We summarize those methods briefly in the following, for a detailed description see Section 2.2.

MtLsb uses a similar strategy to the long-range affinity neighborhood proposed by Lee et al. (2017). The network is taught to simultaneously predict LSDs and affinities (Fig. 1.E).

AcLsD and AcRLsD use an auto-context learning strategy, as proposed by Tu and Bai (2010). This strategy attempts to refine the quality of a prediction by using a cascade of predictors. For

<sup>1</sup>https://github.com/funkelab/lsd

(4)

voxel classification, for example, the first pass of an auto-context classifier predicts voxel labels from raw data. The second pass then uses those predictions from the first pass as input<sup>2</sup>. We loosely adapted this idea when designing our auto-context networks. We first taught a network to predict LSDs from raw EM data. The predicted LSDs were then passed into a second network in order to learn affinities (Fig. 1.E).

We generally observe an increase in affinity prediction accuracy when training to predict LSDs in an auxiliary task. This increase is most noticeable when using an auto-context setup.

#### 2 METHODS

#### 2.1 Local Shape Descriptors

The intuition behind Local Shape Descriptors (LSDs)<sup>3</sup> is to provide an auxiliary learning task that improves boundary prediction by learning statistics describing the local shape of the object close to the boundary. A similar technique was already shown to yield superior results over boundary prediction alone (Bai and Urtasun, 2017). Here, we extend on this idea by predicting for every voxel not just affinities values to neighboring voxels, but also statistics extracted from the object under the voxel aggregated over a local window, specifically: (1) the volume, (2) the voxel-relative center of mass, and (3) pairwise coordinate correlations. See Fig. 1 for a visualization.

More formally, let  $v \in \Omega \subset \mathbb{N}^3$  be the set of voxels in a volume and  $y: \Omega \mapsto \{0,\ldots,l\}$  a ground-truth segmentation. A segmentation induces ground-truth affinity values  $\operatorname{aff}_N^y$ , defined on a voxel-centered neighborhood  $N \subset \mathbb{Z}^3$ , *i.e.*:

$$\operatorname{aff}_N^y: \Omega \mapsto \{0,1\}^{|N|} \quad \operatorname{aff}_N^y(v) = \left(\delta_{y(v) = y(v+n) \neq 0} \mid n \in N\right) \quad (1)$$

where  $\delta$  is the Kronecker function. Our primary learning objective is to infer affinities from raw data  $x : \Omega \mapsto \mathbb{R}$ , *i.e.*, we are interested in learning a function

$$\operatorname{aff}_{N}^{\mathbf{x}}: \Omega \mapsto [0, 1]^{|N|} \tag{2}$$

such that  $\operatorname{aff}_{N}^{x}(v) \approx \operatorname{aff}_{N}^{y}(v)$ .

Similarly to the affinities, we introduce a function to describe the local shape of a segment  $i \in \{1,\ldots,l\}$  under a given voxel v. Let  $b^i: \Omega \mapsto \{0,1\}$  with  $b^i(v) = \delta_{y(v)=i}$  be the binary mask for segment i and  $w: \mathbb{Z}^3 \mapsto \mathbb{R}$  a kernel acting as a local window (e.g., a binary representation of a sphere centered at the origin,  $w(z) = \delta_{|z| < \sigma}$ ). The aggregation of this mask over the window yields the local size  $s^i(v)$  of segment i at position v. Formally, this operation is equal to a convolution of the binary mask with the local window:

$$s^{i}(v) = \sum_{v' \in \Omega} b^{i}(v') w(v - v') = (b^{i} * w)(v).$$
 (3)

To capture local statistics about the segment at location v, we further introduce the following voxel-wise functions m and c. Those functions aggregate the pixel coordinates v over the local window w to compute the local center of mass  $m^i(v)$  and the local covariance of voxel coordinates  $c^i(v)$  for a given segment i:

$$\mathbf{m}_{k}^{i}(v) = \frac{\left(v_{k} \, \mathbf{b}^{i} * \mathbf{w}\right)(v)}{\left(\mathbf{b}^{i} * \mathbf{w}\right)(v)} \qquad k \in \{x, y, z\}$$

$$c_{kl}^{i}(v) = \frac{\left(v_{k}v_{l} b^{i} * w\right)(v)}{\left(b^{i} * w\right)(v)} - m_{k}^{i}(v) m_{l}^{i}(v) \quad k, l \in \{x, y, z\}$$

A derivation of those equations can be found in Supplemental Section A.

To obtain a dense volume of shape descriptors, we compute the above statistics for each voxel with respect to the segment this voxel belongs to. Formally, we evaluate for each voxel v:

$$\tilde{\mathbf{s}}(v) = \mathbf{s}^{\mathbf{y}(v)}(v) \tag{5}$$

$$\tilde{\mathbf{m}}(v) = \left( \mathbf{m}_{x}^{y(v)}(v), \mathbf{m}_{y}^{y(v)}(v), \mathbf{m}_{z}^{y(v)}(v) \right)$$
(6)

$$\tilde{c}(v) = \left(c_{xx}^{y(v)}(v), c_{yy}^{y(v)}(v), \dots, c_{xz}^{y(v)}(v), c_{yz}^{y(v)}(v)\right) \tag{7}$$

The final local shape descriptor  $lsd^y : \Omega \mapsto \mathbb{R}^{10}$  for a voxel v is a concatenation of the size, center offset, and coordinate covariance, *i.e.*,

$$\operatorname{lsd}^{y}(v) = (\underbrace{\tilde{s}(v),}_{\text{size}} \underbrace{\tilde{m}(v) - v,}_{\text{center offset covariance}} \underbrace{\tilde{c}(v)}_{}). \tag{8}$$

For a visual representation of the components of the LSDs, see Fig. 2. We use  $lsd^y(\nu)$  to formulate an auxiliary learning task that complements the prediction of affinities. For that, we use the same neural network to simultaneously learn the functions  $aff^x: \Omega \mapsto [0,1]^{|N|}$  and  $lsd^x: \Omega \mapsto \mathbb{R}^{10}$  directly from raw data x, sharing all but the last convolutional layer of the network.

#### 2.2 Network Architectures

We implement the LSDs using three network architectures. The first is a multitask approach, MTLsD, in which the LSDs are output from a 3D U-NET (Çiçek et al., 2016), along with nearest neighbor affinities in a single pass. The other two methods, AcLsD and AcRLsD are both auto-context setups in which the LSDs from one U-NET are fed into a second U-NET in order to produce the affinities. The former relies solely on the LSDs while the latter also sees the raw data in the second pass. For a visualization of network architectures, see Fig. 1.E and Supplemental Fig. 16. All training was done using Gunpowder<sup>4</sup> and Tensorflow<sup>5</sup>, using the same 3D U-NET architecture (Funke et al., 2019).

## 2.3 Post-Processing

Prediction and post-processing (*i.e.*, watershed and agglomeration) was done in a block-wise fashion using the method described in Funke et al. (2019).

# 3 RESULTS

In this section, we present experimental results of the proposed LSDs for neuron segmentation. We compare the accuracy of the segmentations against several alternative methods for affinity prediction and FFN on three large and diverse datasets. Furthermore, we compare the computational efficiency of different methods and analyze the relationship between different error metrics for neuron segmentations.

<sup>&</sup>lt;sup>2</sup>See https://www.ilastik.org/documentation/autocontext/autocontext for a popular example of this strategy.

<sup>&</sup>lt;sup>3</sup>Distinct from shape descriptors in Maitin-Shepard et al. (2016).

<sup>4</sup>http://funkey.science/gunpowder

<sup>5</sup>https://www.tensorflow.org/

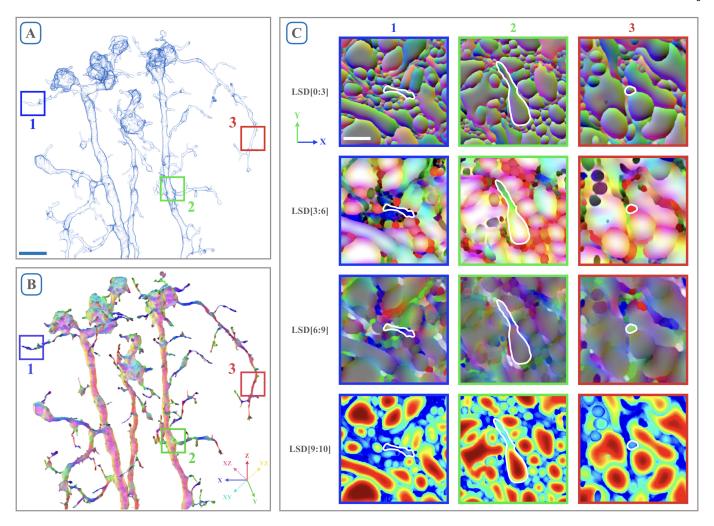


Figure 2: Visualization of LSD components. **A.** Surface mesh of a segmented neuron from FIB-SEM data. Scale bar = 1  $\mu$ m. **B.** RGB mapping of LSD direction vectors (LSD[3:6]). Neural processes are colored with respect to the directions they travel. Intermediate directions are mapped accordingly (see cartesian coordinate inset). **C.** LSD predictions in corresponding 2D slices to the three boxes shown in A and B; neuron highlighted in white. Columns signify neuron orientation (blue = lateral movement, green = vertical movement, red = through-plane movement). Rows correspond to components of the LSDs. First row = mean offset (as seen in Fig. 1). Second and third rows = direction vectors. Second row shows mapping seen in B. Last row = size (number of voxels inside intersected Gaussian). Scale bar = 250 nm.

#### 3.1 Methods

For each dataset we investigated seven methods:

- **Direct neighbor affinities (BASELINE):** Baseline network with a single voxel affinity neighborhood and MSE loss as proposed by Turaga et al. (2010). We trained a 3D U-NET to predict affinities.
- Long range affinities (LR): Same approach as the BASELINE network, but uses an extended affinity neighborhood with three extra neighbors per direction, as proposed by Lee et al. (2017). The extended neighborhood functions as an auxiliary learning task that was shown to improve the direct neighbor affinities.
- MALIS loss (MALIS): Same approach as the Baseline network, but using the loss described in Funke et al. (2019) instead of plain MSE.
- Flood filling networks (FFN): A single segmentation per investigated dataset, provided by the authors of Januszewski

et al. (2018).

- Multi-task LSDs (MTLsD): A network to predict both LSDs and direct neighbor affinities in a single pass, as seen in Fig. 1.E. Similar to LR, the LSDs act as an auxiliary learning task for the direct neighbor affinities.
- Auto-context LSDs (AcLsd): An auto-context setup, where LSDs were predicted from one network and then used as input to a second network in which affinities were predicted.
- Auto-context LSDs with raw (AcRLsD): Same approach as AcLsD, but the second network also receives the raw data as input in addition to the LSDs generated by the first network.

All network architectures are described in detail in Supplemental Section C.1.2 and Supplemental Section D.1.2 for the Zebrafinch and FIB-SEM volumes, respectively.

Each affinity-based network was trained with the same pipeline (i.e. data augmentations and optimizer) and different hyper-parameters with respect to whether the data was anisotropic

(ZEBRAFINCH, Supplemental Section C.1.3) or isotropic (FIB-SEM volumes, Supplemental Section D.1.3).

#### 3.2 Metrics for Neuron Segmentation

Since proofreading of segmentation errors is currently the main bottleneck in obtaining a connectome (Scheffer et al., 2020), metrics to assess neuron segmentation quality should ideally reflect the time needed for proofreading. This requirement is not easily met, since it depends on the tools and strategies used in a proofreading workflow. Currently used metrics aim to correlate scores with the time needed to correct errors based on assumptions about the gravity of certain types of errors. A common assumption has been that false merges take significantly more time to correct than false splits, although next generation proofreading tools challenge this conception (Plaza and Funke, 2018).

In this study, we report neuron segmentation quality with two established metrics: Variation of Information (VoI) and Expected Run-Length (ERL). In addition to those metrics, we propose a new metric, which we refer to as the Min-Cut Metric (MCM), designed to measure the number of graph edit operations needed to perform in a hypothetical proofreading tool.

- Variation of Information (VoI): A metric to compare clusterings (Meilă, 2007), which became an established metric to assess neuron segmentation accuracy. VoI measures the disagreement between two segmentations in terms of the average number of bits needed to guess the segment ID of a randomly chosen voxel in one segmentation, given only its label in the other segmentation. This measurement is performed in both directions, giving rise to the two additive components of VoI, a measure for split and merge errors. Lower values are better, with equivalent segmentations (up to label permutations) having a value of zero.
- Expected Run-Length (ERL): Following the assumption that false merges are in practice harder to correct than false splits, Januszewski et al. (2018) proposed to measure accuracy in terms of the Expected Run Length (ERL), which measures the expected length of an error-free path along neurons in a volume. Notably, all paths contained in falsely merged segments are considered erroneous, thus ERL emphasizes merge errors disproportionally. An appealing aspect of ERL is that it relates segmentation errors to cable length, a commonly used and interpretable feature of neurons.
- Min-Cut Metric (MCM): A metric that assumes that a user can directly interact with agglomerated fragments from an oversegmentation. In particular, we assume that users can split segments by means of a min-cut through the fragment graph between two selected fragments, where edge costs correspond to the merge scores used during agglomeration. In this context, the MCM reports the number of split and merge operations needed to be performed by a human annotator to obtain the desired segmentation. The details of this metric are described in Supplemental Section B.

## 3.3 Datasets

The following describes the datasets, regions of interest (RoIs), and ground-truth used in this study. Dataset overviews can be seen in Fig. 3, Fig. 4, Fig. 5 top panels, and are described in Table 1.

#### 3.3.1 Zebrafinch

The largest dataset used in this study was the songbird dataset also used in Januszewski et al. (2018). This volume consists of a  $\sim 10^6 \mu m^3$  region of a zebrafinch brain, imaged with serial blockface EM at a resolution of 9x9x20 nm (x,y,z). For our experiments, we used a slightly smaller region completely contained inside the raw data with edge lengths of 87.3, 83.7, and 106  $\mu$ m (x,y,z). We refer to this region as the Benchmark Roi.

For each affinity-based network described in Section 3.1, we used 33 volumes containing a total of  $\sim 200 \mu m^3$  ( $\sim 6 \mu m^3$  average per volume) of labeled data<sup>6</sup> for training. We then ran prediction on the Benchmark Roi, using a block-wise processing scheme.

Using the resulting affinities, we generated two sets of supervoxels: one without any masking and one constrained to neuropil using a mask<sup>6</sup>. Additionally, we filtered supervoxels in regions in which the average affinites were lower than a predefined value (e.g., glia). Supervoxels were agglomerated using one of two merge functions, described in Funke et al. (2019), to produce the region adjacency graphs used for evaluation.

We then produced segmentations for RoIs of varying size centered in the Benchmark Roi, in order to assess how segmentation measures scale with volume size. In total, we cropped 10 cubic RoIs ranging from ~11µm to ~76µm edge lengths, in addition to the whole Benchmark Roi. We will refer to the respective RoIs by their edge lengths. For each affinity-based network, in each RoI, we created segmentations for a range of agglomeration thresholds (resulting in a sequence of segmentations ranging from over-to undersegmentation). Additionally, we cropped the provided FFN segmentation accordingly and relabeled connected components.

We used a set of 50 manually ground-truthed skeletons<sup>6</sup>, comprising 97 mm, for evaluation. For each network we assessed VoI and ERL on each RoI. For affinity-based methods we also computed the MCM on the 11, 18, and 25µm RoIs. Additionally, we used 12 validation skeletons consisting of 13.5 mm to determine the optimal thresholds for each network on the Benchmark Roi. Further details can be found in Supplemental Section C.

#### 3.3.2 Hemi-Brain

The so-called Hemi-brain is a FIB-SEM volume of the *Drosophila melanogaster* central brain, imaged at 8nm isotropic resolution (Scheffer et al., 2020), comprising a total of 26 teravoxels of image data. We evaluate all investigated methods on regions restricted to the Ellipsoid Body, a neuropil implicated in spatial navigation, (Turner-Evans and Jayaraman, 2016), which contained ample ground-truth data for evaluation.

We used eight volumes of densely annotated ground-truth volumes containing  $\sim 450 \mu m^3$  of labeled data for training. Three RoIs with  $\sim 12 \mu m$ ,  $\sim 17 \mu m$ , and  $\sim 22 \mu m$  edge lengths were cropped from the larger volume, and prediction was done directly on each RoI.

Supervoxels were limited to the Ellipsoid Body using a mask<sup>7</sup> and then agglomerated using the same two merge functions as in the Zebrafinch dataset.

We produced segmentations for each network over a range of thresholds on the RoIs, and consolidated a single FFN segmentation<sup>6</sup>. A densely labeled ground-truth volume was cropped and filtered using a list of neuron IDs<sup>7</sup> deemed to be

<sup>&</sup>lt;sup>6</sup>Kindly provided by the authors of Januszewski et al. (2018)

<sup>&</sup>lt;sup>7</sup>Kindly provided by the Janelia FlyEM project team (https://janelia.org/project-team/flyem)

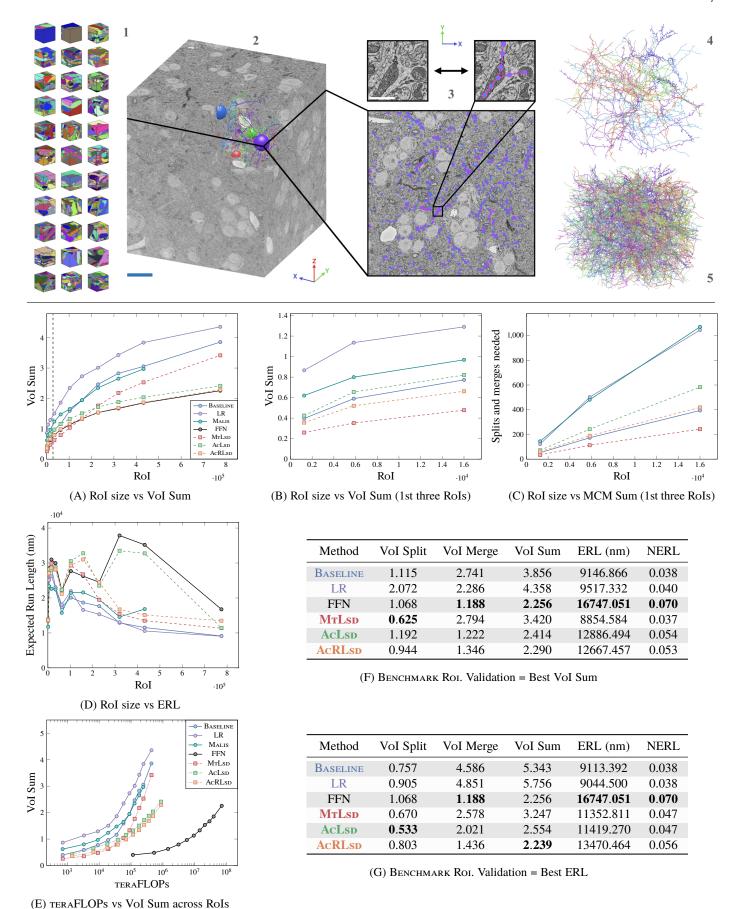


Figure 3: Zebrafinch dataset. Top panel shows data used for training and testing. 1. 33 gound truth volumes were used for training. 2. Full raw dataset, scale bar =  $15 \,\mu\text{m}$ . 3. Single section shows ground-truth skeletons, zoom in scale bar =  $500 \,\text{nm}$ . 4. Validation skeletons (n=12). 5. Testing skeletons (n=50). Main results shown below top panel. Points in Zebrafinch plots correspond to optimal thresholds from validation set. Each point represents an RoI. For VoI and MCM, lower scores are better, for ERL higher scores are better. A. VoI Sum vs RoI size ( $\mu$ m<sup>3</sup>). B,C. VoI Sum and MCM Sum vs RoI size (first three RoIs), respectively. Dashed line in A corresponds to RoIs shown in B,C. D. ERL (nanometers) vs RoI size. E. TeraFLOPs vs VoI Sum across RoIs (same as in panels A and D). Tables (F,G) show best network scores in bold, with respect to which metric was optimized during validation.

Dataset	Imaging Method	Tissue	Resolution (xyz)	Training Data	Testing Data
Zebrafinch	SBFSEM	songbird	9x9x20 nm	33 dense volumes ( $\sim 200  \mu m^3$ )	50 skeletons (97 mm)
Hemi-brain	FIBSEM	Drosophila	8x8x8 nm	8 dense volumes ( $\sim 450  \mu m^3$ )	3 whitelisted volumes ( $\sim 1.7 \times 10^4 \ \mu m^3$ )
Fib-25	FIBSEM	Drosophila	8x8x8 nm	4 dense volumes ( $\sim 160  \mu m^3$ )	1 whitelisted volume ( $\sim 7.7 \times 10^3 \ \mu m^3$ )

Table 1: Overview of datasets used in study

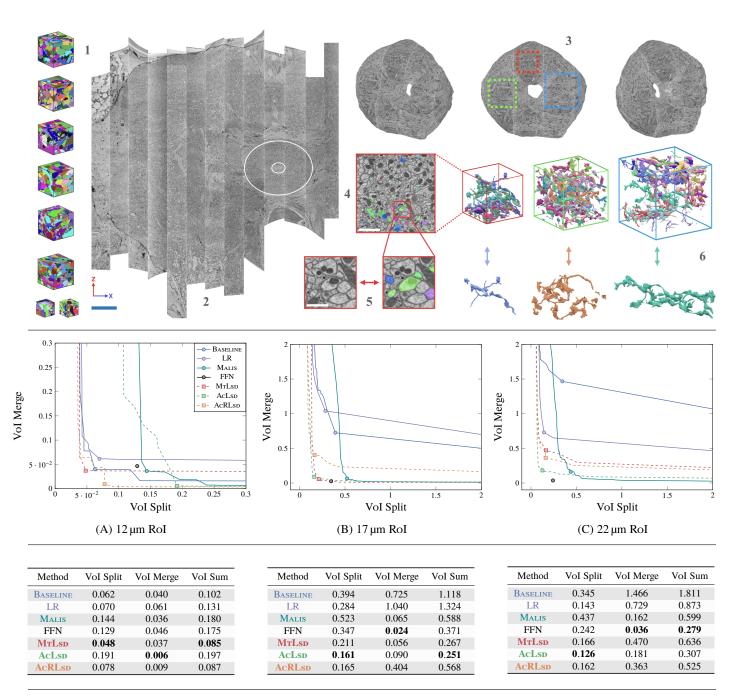


Figure 4: Hemi-brain dataset. Top panel shows data used for training and testing. 1. 8 ground-truth volumes were used for training. 2. Full Hemi-brain volume, scale bar =  $30 \,\mu m$ . Experiments were restricted to Ellipsoid Body (circled region). 3. Volumes used for testing. 4. Example sparse ground-truth testing data, scale bar =  $2.5 \,\mu m$ . 5. Zoom-in, scale bar =  $800 \,nm$ . 6. Example 3D renderings of selected neurons. Main results below top panel. Plot curves show results over range of thresholds for each RoI ( $A = 12 \,\mu m$  RoI,  $B = 17 \,\mu m$  RoI,  $C = 22 \,\mu m$  RoI). Points correspond to optimal thresholds on testing set, no validation set was available. Lower scores are better. Tables show best network scores in bold for each RoI.

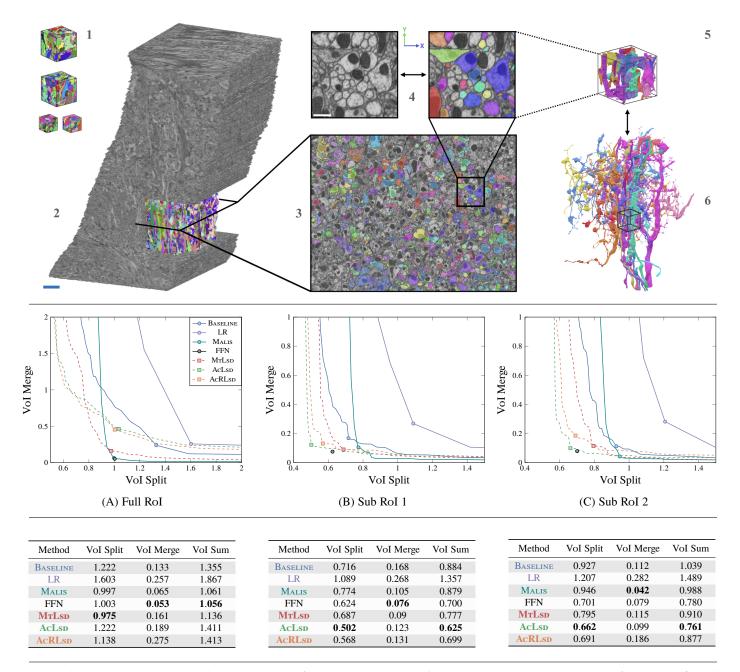


Figure 5: Fib-25 dataset. Top panel shows data used for training and testing. 1. 4 ground-truth volumes were used for training. 2. Full volume with cutout showing testing region, scale bar =  $5 \,\mu m$ . 3. Cross section with sparsely labeled testing ground-truth. 4. Zoom-in, scale bar =  $750 \, nm$ . 5. Sub-volume corresponding to zoomed-in plane. 6. Subset of full RoI testing neurons. Small volume shown for context. Main results below top panel. A. Full testing RoI. B,C. Two sub RoIs contained within full RoI. Plot curves show results over range of thresholds. Points correspond to optimal thresholds on testing set, no validation set was available. Lower scores are better. Tables show best network scores in bold, for each RoI.

completely traced by expert proofreaders. Since the ground-truth is comprised of voxel data rather than skeletons, we report only VoI. Further details can be found in Supplemental Section D.

#### 3.3.3 FIB-25

The Fib-25 dataset, produced by Takemura et al. (2015), is another Fib-SEM volume imaged at 8nm resolution, containing  $\sim\!1.8\times10^5\mu\text{m}^3$  of raw data taken from the *Drosophila* visual system. Four volumes with  $\sim\!160\mu\text{m}^3$  of labeled data  $^7$  were used for training.

We predicted on the full raw data and then restricted supervoxels to an irregularly shaped neuropil mask  $^6$ . Agglomeration was done in the same fashion as the aforementioned volumes. We created segmentations in two ways: for the first method, following the procedure described in Januszewski et al. (2018), we segmented neurons within the entire neuropil mask. For the second method, we limited segmentation to two sub-RoIs contained within both the neuropil mask and testing region (sections 5074–7950). The sub-RoIs have a size of  $\sim\!2.2\times10^3\mu\text{m}^3$  and  $\sim\!2.6\times10^3\mu\text{m}^3$ , respectively. For FFN, we cropped the provided segmentation  $^6$  and relabeled connected components.

Evaluation was limited to a list of proofread, voxel ground-truth labels<sup>7</sup> contained inside the testing region. Depending on the segmentation method, connected components on both ground-truth and segmentation volumes were either cropped and relabelled or untouched. Since the ground-truth is comprised of voxel data rather than skeletons, we report only VoI. Further details can be found in Supplemental Section D.

# 3.4 Neuron Segmentation Accuracy

## 3.4.1 Zebrafinch

We find that LSDs are useful for improving the accuracy of direct neighbor affinities, and subsequently the resulting segmentations. Specifically, LSD-based methods consistently outperform other affinity-based methods over a range of RoIs, whether used in a multitask (MTLSD) or auto-context (AcLsD and AcRLSD) architecture (Fig. 3.A). In terms of segmentation accuracy according to VoI, the best auto-context network (AcRLSD) performs on par with FFN (Fig. 3.A).

Interestingly, we find that the ranking of methods depends on the size of the evaluation RoI. Even for monotonic metrics like VoI, we see that performance on the smallest RoIs (up to  $54\,\mu m$ ) does not extrapolate to the performance on larger datasets.

We also investigated how ERL varies over different RoI sizes. To this end, we cropped the skeleton ground-truth to the respective RoIs and relabeled connected components (as we did for the VoI evaluation). However, the resulting fragmentation of skeletons heavily impacts ERL scores: ERL can not exceed the average length of skeletons, and thus the addition of shorter skeletons fragments can result in a decrease of ERL, even in the absence of errors. We see this effect prominently in Fig. 3.D: ERL measures do not progress monotonically over RoI sizes and absolute values are likely not comparable across different dataset sizes. In addition, the ranking of methods for a given RoI size varies significantly over different RoI sizes.

The high variability between metrics and RoI sizes prompted us to develop a metric that aims to measure proofreading effort. We developed MCM to count the number of interactions needed to split and merge neurons in order to correctly segment the ground-truth skeletons, assuming that a min-cut-based split tool

is available. Due to the computational cost associated with MCM (stemming from repeated min-cuts in large fragment graphs), we limited its computation to the three smallest investigated RoIs in this dataset. As expected, we observe a linear increase in MCM with RoI size across different methods (Fig. 3.C). Furthermore, we see that MCM and VoI mostly agree on the ranking of methods (Fig. 3.B, Supplemental Fig. 14), which suggests that VoI should be preferred to compare segmentation quality in the context of a proofreading workflow that allows annotators to split false merges using a min-cut on the fragment graph. We could not compute MCM for FFN since the MCM requires a fragment graph.

# 3.4.2 Hemi-Brain

We observe a similar variability of method rankings over RoI sizes on the Hemi-brain dataset.

On the largest investigated RoI ( $22\mu m$  edge length), AcLso clearly performs best among all affinity-based methods (VoI sum of 0.307 vs. 0.525 for the second best, see Fig. 4.C). As such, AcLso is competitive with FFN (VoI sum 0.279), with the notable difference of performing more merge errors, but significantly less split errors than FFN.

On the 17 $\mu$ m RoI, AcLsD again performs better than all other affinity-based methods and also better than FFN (VoI sum of 0.251 vs. 0.371, see Fig. 4.B). MtLsD is on par with AcLsD on this RoI, which stands in contrast to its significantly worse performance on the 22 $\mu$ m RoI. This observation further puts into question to what extent the accuracy of a segmentation can be extrapolated from smaller to larger RoI sizes.

These concerns become even more evident when turning to the results on the smallest RoI of 12µm edge length. Here, AcLsp performs worse than all other methods, with a significant margin to the best performing method, MtLsp (VoI sum 0.197 vs 0.085). Even Baseline achieves very good results on this RoI (VoI sum 0.102), although it would be a poor choice in production given its detrimental performance on the larger RoIs. We have to conclude that the size of this RoI is likely not large enough to accurately deduce whether the differences in method performance are due to model accuracies or data biases.

A somewhat surprising result is the performance of AcRLsD on this dataset. Although architecturally very similar to AcLsD (the only difference is that AcRLsD receives raw data and LSDs in the second pass, while AcLsD receives only LSDs), AcRLsD is significantly worse on the two larger RoIs. This stands in contrast to the results we obtained on the Zebrafinch dataset. Our results do not allow us to say with confidence whether this artifact is due to overfitting to the training data (which might be more likely to happen for AcRLsD) or due to model noise introduced by the random initialization during training.

#### 3.4.3 FIB-25

We first evaluated all methods on the full testing RoI of the Fib-25 dataset (Fig. 5.A). On the full RoI, LSDs generally do not perform well. We observe that the best auto-context method (AcLsd) performs worse than the Baseline (VoI sum 1.413 vs 1.355). Interestingly, MtLsd achieves higher accuracy than both auto-context networks. FFN exceeds all other methods (VoI sum 1.056)<sup>8</sup>, and Malis is not far behind (VoI sum 1.061). Since

 $^8{\rm The}$  authors of Januszewski et al. (2018) report a VoI split of 0.8837. While we were able to replicate the reported VoI merge score of 0.053, we found VoI split to be 1.003.

those results are not consistent with the results seen on the Zebrafinch and Hemi-brain volumes, we visually inspected the segmentations. We found a high rate of false merges occurring in the periphery of the testing RoI, stemming from nuclei and boundaries of the imaged volume, which are not contained in the training data. As such, the full testing RoI of this dataset favors "conservative" methods, *i.e.*, methods that have higher split rates.

To test the plain neuropil segmentation accuracy, we further cropped two RoIs ( $\sim 2.2 \times 10^3 \mu m^3$  and  $\sim 2.6 \times 10^3 \mu m^3$ ) inside the testing region, such that they contain only dense neuropil (see Supplemental Section D.3.3 for details).

On the two sub RoIs, LSDs outperform other affinity-based methods and are comparable to FFN, (Fig. 5.B,C). Consistent with the Zebrafinch and Hemi-brain results, using an auto-context approach seems to generate the best results. On sub RoI 1, AcLso slightly exceeds the accuracy of FFN (VoI sum: 0.625 vs 0.700, respectively). We observe similar results on sub RoI 2 (VoI sum: 0.761 vs 0.780, respectively).

These results highlight the need for masking neuropil when processing large volumes, as was done in the Zebrafinch dataset. Interestingly, LR affinities perform poorly across RoIs, which might suggest that an increased affinity neighborhood is sometimes not sufficient for improving direct neighbor affinities.

# 3.5 Throughput

In addition to being accurate, it is important for neuron segmentation methods to be fast and computationally inexpensive. As described in the introduction, the acquisition size of datasets is growing rapidly and approaches should therefore aim to complement this trajectory. LSD-based methods can be parallelized in the same manner as BASELINE affinities, making them a good candidate for the processing of very large datasets or environments with limited computing resources. In our experiments, prediction and segmentation of affinity-based methods was done in a blockwise fashion, allowing parallel processing across many workers, see Fig. 7 for an overview. This allowed for efficient segmentation following prediction (Table 2.B).

When considering computational costs in terms of floating point operations (FLOPs), we find that the AcRLsD network (the computationally most expensive of all LSD architectures) is two orders of magnitude more efficient than FFN, while producing a segmentation of comparable quality (Fig. 3.E). For this comparison, we computed FLOPs of all affinity-based methods during prediction (see Supplemental Section E for details). For FFN, we used the numbers reported in Januszewski et al. (2018), limited to the forward and backward passes of the network, i.e., the equivalent of the prediction pass for affinity-based methods. We limit the computational cost analysis to GPU operations, since FLOP estimates on CPUs are unreliable and the overall throughput is dominated by GPU operations. We therefore only consider inference costs for all affinity-based networks, since agglomeration is a post-processing step done on the CPU. To keep the comparison to FFN fair, we do not count FLOPs during FFN agglomeration, although it involves a significant amount of GPU operations. Table 2.A summarizes the computational costs of all investigated methods and provides throughput numbers on the hardware used in this study. Generally, affinity-based methods are more computationally efficient than FFN by two orders of magnitude.

# 4 Discussion

The main contribution of this work is the introduction of LSDs as an auxiliary learning task for neuron segmentation. We demonstrated that, when compared to other affinity-based methods, the LSDs consistently help to improve neuron segmentations across specimen, resolution, and imaging techniques. We also found results to be competitive with the current state of the art approach, while being two orders of magnitude faster. All methods, datasets, and results are publicly available<sup>9</sup>, which we hope will be a useful starting point for further extensions and a benchmark to evaluate future approaches in a comparable manner.

On the Zebrafinch dataset, the largest dataset both in terms of image data and available ground-truth, LR and Malis did not exceed the accuracy of the Baseline network. This is surprising considering that LR also uses an auxiliary task to improve direct neighbor affinity predictions. This observation might well be due to differences in how masks are handled during training, which we will discuss in more detail in Section 4.4. We found those results generally confirmed on the Hemi-brain dataset, with the exception of Malis, which performed better on this dataset than on Zebrafinch.

Although relatively small in comparison, the Fib-25 dataset nevertheless provided insights into network performance in the periphery of dense neuropil. On the full RoI, both auto-context LSD networks performed very poorly, while Malis was on par with FFN. This is due to increased false merges in the periphery, which proliferate into the testing RoI. Networks that favor splitting over merging, such as Malis and FFN, are consequently less affected. Further analysis on two sub RoIs, both contained entirely within the neuropil mask and the testing RoI, confirmed this to be the case: here, LSD networks improve over other affinity networks and are again competitive with FFN. These results highlight the importance of masking when processing large volumes.

# 4.1 Metric Evaluation

An important but challenging task is finding a robust metric for assessing the quality of a neuron segmentation. Ideally, such a metric reflects the amount of time needed to proofread a segmentation. Here, we presented results in terms of VoI and ERL, two commonly used metrics for this task.

VoI directly reports the amount of split and merge errors. Being a voxel-wise metric, however, VoI can be sensitive to slight, but systematic, shifts in boundaries. At the same time, small topological changes might go unnoticed, which is especially problematic in fine neurites in the vicinity of synapses (Funke et al., 2017).

ERL reports the expected error-free path-length of a reconstruction with respect to skeleton ground-truth. Similar to VoI, ERL is not sensitive to small topological changes close to terminals. Furthermore, ERL disproportionately punishes merge errors and subsequently favors split-preferring methods (Plaza and Funke, 2018). Additionally, we found that ERL increases non-monotonically with varying volume sizes (Fig. 3.D), which is due to the fragmentation of skeletons in volumes that are not large enough to contain entire neurons.

Consequently, neither method directly reflects the labor required for proofreading a segmentation, which is arguably the relevant quantity to optimize (Plaza, 2016; Funke et al., 2017).

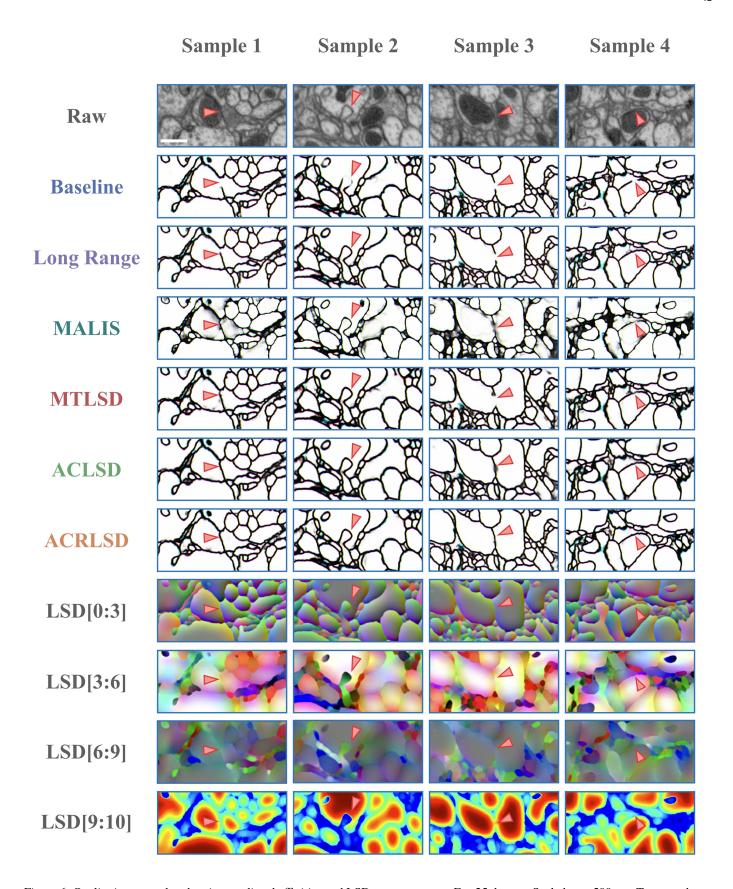


Figure 6: Qualitative examples showing predicted affinities and LSD components on Fib-25 dataset. Scale bar = 500 nm. Top row shows raw data. Arrows correspond to example plasma membranes. A baseline architecture (row two) fails to distinguish these boundaries, which could lead to possible errors in a resulting segmentation. The LSDs (bottom 4 rows) help to improve the affinites generated by the Mtlsb, Aclsb and AcRlsb networks.

#### (A) Prediction

Method	Workers	Total Process Time (seconds)	Throughput ( $\mu m^3$ / GPU seconds)	teraFLOPs
LSDs	100 (2080Ti)	7,090	1.093	437,000
BASELINE	60 (2080Ti)	8,449	1.528	437,000
LR	60 (2080Ti)	10,596	1.218	440,000
Malis	60 (2080Ti)	8,522	1.515	437,000
FFN	1,000 (P100)	11,054	0.07	70,500,000
MtLsd	60 (2080Ti)	9,193	1.404	440,000
AcLsd	15 (V100)	43,972	1.174	874,000
AcRLsd	24 (V100)	37,837	0.854	874,000

## (B) Segmentation (MTLsD example)

Method	Workers	Total Process Time (seconds)	Throughput ( $\mu m^3$ / CPU seconds)
Watershed	100 CPUs	7,780	0.996
Agglomeration	100 CPUs	19,859	0.390

Table 2: Computational costs on Zebrafinch Benchmark Roi.

This quantity depends on the available tools for proofreading, and in particular on the amount of interactions needed to fix errors of different kinds: False splits might be hard to find, but do require only one interaction to merge. False merges, on the other hand, might be easy to spot, but the number of interactions needed to fix them depends greatly on the proofreading tool. Current proofreading tools 10,11 allow annotators to correct merge errors with a few interactions. We therefore introduced the MCM, a metric which uses graph cuts to emulate the amount of interactions required to correct false merges in a segmentation. We observed a linear growth of MCM with volume size (Fig. 3.C), which is a necessary condition for neuron segmentation metrics that measure the amount of proofreading effort needed (assuming an equal distribution of errors).

Unfortunately, MCM is computationally quite expensive. The sequence of graph-cuts needed for the evaluation of merge errors quickly becomes infeasible on large volumes. However, MCM shows general ranking agreement with VoI, evaluated on 11 randomly sampled sub-RoIs in ZEBRAFINCH (Fig. 10) and across different thresholds (Supplemental Fig. 14). These findings suggest that VoI can serve as a reasonable proxy to rank methods based on their expected proofreading time.

Additionally, we find VoI to be a robust metric for the validation of method parameters: For each affinity-based network, the threshold minimizing VoI sum on the validation set is also close to the best threshold on the testing set (Supplemental Fig. 12). This property is of practical relevance, as in any real-world scenario hyperparameters have to be adjusted on a volume that is significantly smaller than the target volume. Unfortunately, ERL does not seem to exhibit this property to the same degree: the best validation thresholds gradually diverge from the best testing thresholds as scale increases (Supplemental Fig. 13), which makes it difficult to extrapolate segmentation accuracy from a validation dataset.

#### 4.2 Auxiliary Learning for Boundary Prediction

Auxiliary learning tasks have been shown to improve network performance across different applications. One possible explanation for why auxiliary learning is also helpful for the prediction of neuron boundaries is that the additional task incentivizes the network to consider higher-level features. Predicting LSDs is likely harder than boundaries, since additional local structure of the object has to be considered. Merely detecting an oriented, dark sheet (e.g., plasma membranes) is not sufficient; statistics of the whole neural process have to be taken into account. Those statistics rely on features that are not restricted to the boundary in question. Therefore, the network is forced to make use of more information in its receptive field than is necessary for boundary prediction alone. This, in turn, increases robustness to local ambiguities and noise for the prediction of LSDs. As a welcome side-effect, it seems that the network learns to correlate boundary prediction with LSD prediction, which explains why the boundary prediction benefits from using the LSDs as an auxiliary objective.

Surprisingly, we see that LR affinities do not perform as well across the investigated datasets. While LR affinities share some of the benefits of LSDs, they might not be as efficient as LSDs in encoding higher-level features. For example, LR affinities have blind spots (missing neighborhood steps), whereas LSDs are spatially homogeneous. Additionally, we found LR affinities to be detrimental when used with masking of glia and other structures. It is likely harder to correlate nearest neighbor affinities with a long-range neighborhood in the presence of masks.

While Lee et al. (2017) saw superhuman accuracy using an increased affinity neighborhood on the SNEMI3D challenge, the processed volume was relatively small (~110µm<sup>3</sup>). Our results suggest that it is hard to correlate accuracy on small volumes to accuracy on large volumes (Fig. 3.A)12. Additionally, we only consider an increased affinity neighborhood and not other aspects of the original LR implementation, such as residual modules in the U-NeT and inference blending, which might be essential for further performance increases. Finally, the SNEMI3D dataset has an anisotropy factor of  $\sim$ 5, whereas the data we test on here has an anisotropy factor of either ~2 (Zebrafinch) or is isotropic (Hemi-BRAIN, FIB-25).

## 4.3 Auto-Context Refinement

We found that using an auto-context approach greatly improved resulting segmentations (Fig. 3.A). We tested if this increase

<sup>10</sup>https://flywire.ai

<sup>11</sup>https://github.com/janelia-flyem/NeuTu

<sup>&</sup>lt;sup>12</sup>Also shown by the authors of Januszewski et al. (2018) in supplementary tables 4 and 5.

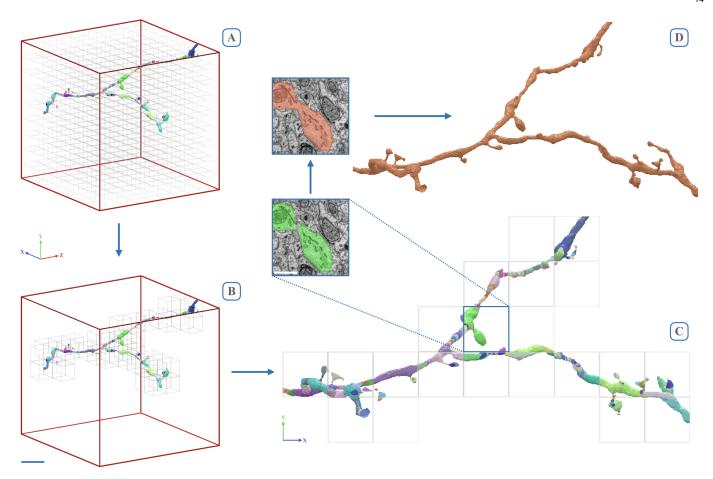


Figure 7: Overview of block-wise processing scheme. **A.** Example 32  $\mu$ m RoI showing total block grid. **B.** Required blocks to process example neuron. Scale bar =  $\sim 6\mu$ m. **C.** Corrsponding orthographic view highlights supervoxels generated during watershed. Block size = 3.6  $\mu$ m. Inset shows respective raw data inside single block (scale bar =  $\sim 1\mu$ m). **D.** Supervoxels are then agglomerated to obtain a resulting segment. **Note:** While this example shows processing of a single neuron, in reality all neurons are processed simultaneosuly.

in accuracy was consistent when using affinities as the input to the second network (*i.e.*, a Baseline auto-context approach, AcBaseline), and found that it made no significant improvements to the Baseline network (Fig. 8).

We hypothesize that predicting affinities from affinities is too similar to predicting affinities from raw EM data. Specifically, we suspect that the ACBASELINE network simply copies data in the second pass rather than learning anything new. Easy solutions, such as looking for features like oriented bars, already produce relatively accurate boundaries in the first pass. Consequently, there is little incentive for the network to change course in the second pass. Translating from LSDs to affinities, on the other hand, is a comparatively different task, which forces the network to incorporate the features from the LSDs in the second pass. The subsequent boundary predictions seem to benefit from this.

## 4.4 Masking

Binary masks are commonly used to limit neuron segmentation to dense neuropil and exclude confounding structures like glia cells. Recent approaches to processing large volumes have incorporated tissue masking at various points in the pipeline (Januszewski et al., 2018; Li et al., 2019; Dorkenwald et al., 2019; Scheffer et al., 2020), to prevent errors in areas that were underrepresented in the training data.

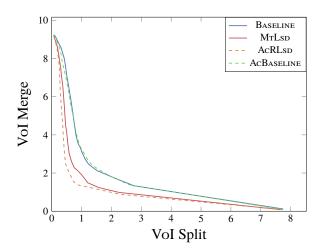


Figure 8: Zebrafinch, Benchmark Roi, VoI split vs VoI merge, auto-context comparison.

Our results confirm the importance of masking. We used a neuropil mask which excluded cell bodies, blood vessels, myelin, and out-of-sample (background) voxels (see Supplemental Fig. 20 for a visualization). Across all investigated methods, the accuracy

degraded substantially on larger RoIs when processed without masking (see Fig. 9, Supplemental Fig. 11).

Masking of irrelevant structures can also be incorporated in the training process. The Zebrafinch training volumes already had some glial processes masked out. We trained all networks to predict zero affinities in these regions (see Supplemental Fig. 18 for a visualization). We then discarded fragments with close to zero affinity values during agglomeration. Methods which succeeded in learning to mask these areas (Baseline, MtLsd, AcLsd, AcRLsd) produced better results than those that did not (Malis, LR).

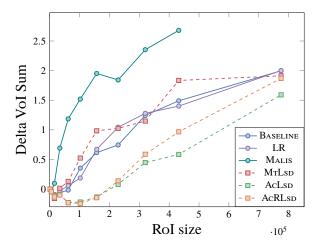


Figure 9: Zebrafinch, mask delta VoI sum vs RoI.

# 4.5 Accuracy Extrapolation

One of the challenges of deep learning is to find representative testing data and metrics to infer production performance. This is especially challenging for neuron segmentation, considering the diversity of neural ultrastructure and morphology found in EM volumes. While challenges like CREMI<sup>13</sup> and SNEMI3D<sup>14</sup> make an effort to include representative training and testing data, the implications for model performance on larger datasets are not straight forward.

Our results suggest that testing on small volumes provides limited insight into the quality of a method when applied to larger volumes. For example, the total volume of the three CREMI testing datasets ( $\sim\!1056~\mu m^3$ ) is still less than the smallest Zebrafinch ( $\sim\!1260~\mu m^3$ ) and Hemi-brain ( $\sim\!1643~\mu m^3$ ) RoIs. As shown in Fig. 3 and Fig. 4, the smallest volumes are not indicative of performance on larger volumes. In this context, it seems difficult to declare a clear "winner" when it comes to neuron segmentation accuracy. Dataset sizes and the choice of evaluation metrics greatly influence which method is considered successful.

# 5 CONCLUSIONS AND FUTURE DIRECTIONS

Although adding LSDs as an auxiliary learning task significantly increases accuracy, it is unclear whether different shape descriptors could lead to further improvements. The LSDs proposed here were subjectively engineered based on features that we expected to be important to encode object shape. Future experiments could incorporate different features or focus on learning an optimal embedding rather than a hand-designed one. In that context, we

note that it is not clear whether each component of the LSD embedding contributes equally to the improvement of affinity predictions.

Currently, we only use LSDs as an auxiliary learning task. As a result, affinities are still required to produce a segmentation. Whether this is really needed is an open question, since the predicted LSDs already identify objects reliably. An interesting future direction would be to use the predicted local shape information directly for fragment agglomeration. As an intermediate step, LSDs can serve to provide a second source of information for identifying errors in a segmentation. Once a segmentation is generated, LSDs could be calculated on the labels and then compared with the initial LSD predictions. The difference between the two would likely highlight regions containing errors (Fig. 23.C).

LSDs were designed for the goal of neuron segmentation but might also be applicable to other instance segmentation problems. As an example, LSDs have already been successfully used to generate segmentations of cell bodies and mitochondria (Fig. 23.A,B). Additionally, Gallusser et al. (2020) demonstrate promising results using LSDs for Golgi apparatus and endoplasmic reticulum segmentation. In general, we believe that objects that have a blob-like structure such as other organelles and various cell types would likely benefit from LSDs. Furthermore, the direction vectors of the LSDs provide insight into neuropil vs. tract regions of the brain (Fig. 23.D). These predictions could be leveraged in order to generate better tissue masks. While the LSDs presented here were conceived for a specific instance segmentation task, it would be interesting to see the LSDs extended and applied to other microscopy problems.

# 6 ACKNOWLEDGEMENTS

We thank Caroline Malin-Mayor, William Patton, and Julia Buhmann for code contributions; Nils Eckstein and Julia Buhmann for helpful discussions; Stuart Berg for code to help with data acquisition; Jeremy Maitin-Shepard for helpful feedback on Neuroglancer; Viren Jain, Michał Januszewski, Jörgen Kornfeld, and Steve Plaza for access to data used for training and evaluation. **Funding.** This work was supported by the Howard Hughes Medical Institute. Author Contributions. Conceptualization: Srinivas Turaga, Jan Funke. Funding acquisition: Stephan Saalfeld, Srinivas Turaga, Jan Funke. Software: Arlo Sheridan, Diptodip Deb, Tri Nguyen, Jan Funke. Data consolidation: Arlo Sheridan, Diptodip Deb, Jan Funke. Evaluation: Arlo Sheridan, Jan Funke. Data dissemination: Arlo Sheridan, Jan Funke. Visualization: Arlo Sheridan, Jan Funke. Writing (original draft): Arlo Sheridan, Jan Funke. Writing (review and editing): Arlo Sheridan, Tri Nguyen, Diptodip Deb, Wei-Chung Allen Lee, Uri Manor, Jan Funke.

# REFERENCES

Abbott, L. F., Bock, D. D., Callaway, E. M., Denk, W., Dulac, C., Fairhall, A. L., Fiete, I., Harris, K. M., Helmstaedter, M., Jain, V., Kasthuri, N., LeCun, Y., Lichtman, J. W., Littlewood, P. B., Luo, L., Maunsell, J. H. R., Reid, R. C., Rosen, B. R., Rubin, G. M., Sejnowski, T. J., Seung, H. S., Svoboda, K., Tank, D. W., Tsao, D., and Van Essen, D. C. (2020). The Mind of a Mouse. *Cell*, 182(6):1372–1376.

Bai, M. and Urtasun, R. (2017). Deep Watershed Transform for Instance Segmentation. pages 5221–5229.

Bates, A. S., Schlegel, P., Roberts, R. J. V., Drummond, N., Tamimi, I. F. M., Turnbull, R., Zhao, X., Marin, E. C., Popovici, P. D., Dhawan, S., Jamasb, A., Javier, A., Serratosa Capdevila, L., Li, F.,

<sup>13</sup>https://cremi.org

<sup>14</sup>http://brainiac2.mit.edu/SNEMI3D

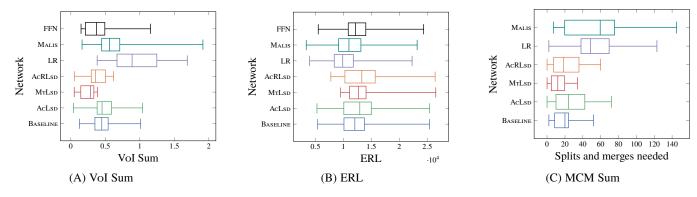


Figure 10: Metric distribution across 11 randomly sampled non-masked sub RoIs in the Zebrafinch.

Rubin, G. M., Waddell, S., Bock, D. D., Costa, M., and Jefferis, G. S. X. E. (2020). Complete Connectomic Reconstruction of Olfactory Projection Neurons in the Fly Brain. *Current Biology*, 30(16):3183–3199.e6.

Beier, T., Pape, C., Rahaman, N., Prange, T., Berg, S., Bock, D. D., Cardona, A., Knott, G. W., Plaza, S. M., Scheffer, L. K., et al. (2017). Multicut brings automated neurite segmentation closer to human performance. *Nature methods*, 14(2):101–102.

Briggman, K. L. and Bock, D. D. (2012). Volume electron microscopy for neuronal circuit reconstruction. *Current Opinion in Neurobiol*ogy, 22(1):154–161.

Buhmann, J., Krause, R., Lentini, R. C., Eckstein, N., Cook, M., Turaga, S., and Funke, J. (2018). Synaptic partner prediction from point annotations in insect brains. arXiv:1806.08205 [cs]. arXiv: 1806.08205.

Buhmann, J., Sheridan, A., Gerhard, S., Krause, R., Nguyen, T.,
Heinrich, L., Schlegel, P., Lee, W.-C. A., Wilson, R., Saalfeld,
S., Jefferis, G., Bock, D., Turaga, S., Cook, M., and Funke, J.
(2020). Automatic Detection of Synaptic Partners in a Whole-Brain Drosophila EM Dataset. bioRxiv, page 2019.12.12.874172.
Publisher: Cold Spring Harbor Laboratory Section: New Results.

Çiçek, O., Abdulkadir, A., Lienkamp, S. S., Brox, T., and Ronneberger, O. (2016). 3d U-Net: Learning Dense Volumetric Segmentation from Sparse Annotation. In *Medical Image Computing* and Computer-Assisted Intervention – MICCAI 2016, pages 424– 432.

Ciresan, D., Giusti, A., Gambardella, L. M., and Schmidhuber, J. (2012). Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, Advances in Neural Information Processing Systems 25, pages 2843–2851. Curran Associates, Inc.

Dorkenwald, S., Turner, N. L., Macrina, T., Lee, K., Lu, R., Wu, J., Bodor, A. L., Bleckert, A. A., Brittain, D., Kemnitz, N., Silversmith, W. M., Ih, D., Zung, J., Zlateski, A., Tartavull, I., Yu, S.-C., Popovych, S., Wong, W., Castro, M., Jordan, C. S., Wilson, A. M., Froudarakis, E., Buchanan, J., Takeno, M., Torres, R., Mahalingam, G., Collman, F., Schneider-Mizell, C., Bumbarger, D. J., Li, Y., Becker, L., Suckow, S., Reimer, J., Tolias, A. S., Costa, N. M. d., Reid, R. C., and Seung, H. S. (2019). Binary and analog variation of synapses between cortical pyramidal neurons. *bioRxiv*, page 2019.12.29.890319. Publisher: Cold Spring Harbor Laboratory Section: New Results.

Funke, J., Klein, J., Moreno-Noguer, F., Cardona, A., and Cook, M. (2017). TED: A Tolerant Edit Distance for segmentation evaluation. *Methods*, 115:119–127.

Funke, J., Tschopp, F., Grisaitis, W., Sheridan, A., Singh, C., Saalfeld, S., and Turaga, S. C. (2019). Large Scale Image Segmentation with Structured Loss Based Deep Learning for Connectome Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(7):1669–1680.

Gallusser, B., Vadakkan, T. J., Sahasrabudhe, M., Di Caprio, G., and

Kirchhausen, T. (2020). Deep-learning-based automatic organelle segmentation in 3d electron microscopy datasets. *in preparation*.

Heinrich, L., Funke, J., Pape, C., Nunez-Iglesias, J., and Saalfeld, S. (2018). Synaptic Cleft Segmentation in Non-Isotropic Volume Electron Microscopy of the Complete Drosophila Brain. arXiv:1805.02718 [cs]. arXiv: 1805.02718.

Hildebrand, D. G. C., Cicconet, M., Torres, R. M., Choi, W., Quan, T. M., Moon, J., Wetzel, A. W., Scott Champion, A., Graham, B. J., Randlett, O., Plummer, G. S., Portugues, R., Bianco, I. H., Saalfeld, S., Baden, A. D., Lillaney, K., Burns, R., Vogelstein, J. T., Schier, A. F., Lee, W.-C. A., Jeong, W.-K., Lichtman, J. W., and Engert, F. (2017). Whole-brain serial-section electron microscopy in larval zebrafish. *Nature*, 545(7654):345–349. Number: 7654 Publisher: Nature Publishing Group.

Huang, G. B., Scheffer, L. K., and Plaza, S. M. (2018). Fully-Automatic Synapse Prediction and Validation on a Large Data Set. Frontiers in Neural Circuits, 12. Publisher: Frontiers.

Hulse, B. K., Haberkern, H., Franconville, R., Turner-Evans, D. B.,
Takemura, S.-y., Wolff, T., Noorman, M., Dreher, M., Dan, C.,
Parekh, R., Hermundstad, A. M., Rubin, G. M., and Jayaraman,
V. (2020). A connectome of the Drosophila central complex reveals network motifs suitable for flexible navigation and context-dependent action selection. bioRxiv, page 2020.12.08.413955.
Publisher: Cold Spring Harbor Laboratory Section: New Results.

Januszewski, M., Kornfeld, J., Li, P. H., Pope, A., Blakely, T., Lindsey, L., Maitin-Shepard, J., Tyka, M., Denk, W., and Jain, V. (2018). High-precision automated reconstruction of neurons with flood-filling networks. *Nature Methods*, 15(8):605.

Kornfeld, J., Benezra, S. E., Narayanan, R. T., Svara, F., Egger, R., Oberlaender, M., Denk, W., and Long, M. A. (2017). EM connectomics reveals axonal target variation in a sequence-generating network. *eLife*, 6:e24364. Publisher: eLife Sciences Publications, Ltd.

Kreshuk, A., Funke, J., Cardona, A., and Hamprecht, F. A. (2015). Who Is Talking to Whom: Synaptic Partner Detection in Anisotropic Volumes of Insect Brain. In Navab, N., Hornegger, J., Wells, W. M., and Frangi, A., editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Lecture Notes in Computer Science, pages 661–668, Cham. Springer International Publishing.

Lee, K., Lu, R., Luther, K., and Seung, H. S. (2019). Learning Dense Voxel Embeddings for 3D Neuron Reconstruction. arXiv:1909.09872 [cs]. arXiv: 1909.09872.

Lee, K., Zung, J., Li, P., Jain, V., and Seung, H. S. (2017). Superhuman Accuracy on the SNEMI3d Connectomics Challenge. *arXiv:1706.00120 [cs]*.

Lee, W.-C. A., Bonin, V., Reed, M., Graham, B. J., Hood, G., Glattfelder, K., and Reid, R. C. (2016). Anatomy and function of an excitatory network in the visual cortex. *Nature*, 532(7599):370–374. Number: 7599 Publisher: Nature Publishing Group.

Li, P. H., Lindsey, L. F., Januszewski, M., Zheng, Z., Bates, A. S., Taisz, I., Tyka, M., Nichols, M., Li, F., Perlman, E., Maitin-

- Shepard, J., Blakely, T., Leavitt, L., Jefferis, G. S. X. E., Bock, D., and Jain, V. (2019). Automated Reconstruction of a Serial-Section EM Drosophila Brain with Flood-Filling Networks and Local Realignment. *bioRxiv*, page 605634. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Maitin-Shepard, J. B., Jain, V., Januszewski, M., Li, P., and Abbeel, P. (2016). Combinatorial Energy Learning for Image Segmentation. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R., editors, Advances in Neural Information Processing Systems 29, pages 1966–1974. Curran Associates, Inc.
- Meilă, M. (2007). Comparing clusterings—an information based distance. *Journal of Multivariate Analysis*, 98(5):873–895.
- Motta, A., Berning, M., Boergens, K. M., Staffler, B., Beining, M., Loomba, S., Hennig, P., Wissler, H., and Helmstaedter, M. (2019). Dense connectomic reconstruction in layer 4 of the somatosensory cortex. *Science*, 366(6469). Publisher: American Association for the Advancement of Science \_eprint: https://science.sciencemag.org/content/366/6469/eaay3134.full.pdf.
- Nguyen, T., Malin-Mayor, C., Patton, W., and Funke, J. (2020). Daisy: A library for block-wise task scheduling for large nd volumes. in preparation.
- Phelps, J. S., Hildebrand, D. G. C., Graham, B. J., Kuan, A. T., Thomas, L. A., Nguyen, T. M., Buhmann, J., Azevedo, A. W., Sustar, A., Agrawal, S., Liu, M., Shanny, B. L., Funke, J., Tuthill, J. C., and Lee, W.-C. A. (2021). Reconstruction of motor control circuits in adult Drosophila using automated transmission electron microscopy. *Cell*, 0(0). Publisher: Elsevier.
- Plaza, S. M. (2016). Focused Proofreading to Reconstruct Neural Connectomes from EM Images at Scale. In Carneiro, G., Mateus, D., Peter, L., Bradley, A., Tavares, J. M. R. S., Belagiannis, V., Papa, J. P., Nascimento, J. C., Loog, M., Lu, Z., Cardoso, J. S., and Cornebise, J., editors, *Deep Learning and Data Labeling for Medical Applications*, Lecture Notes in Computer Science, pages 249–258, Cham. Springer International Publishing.
- Plaza, S. M. and Funke, J. (2018). Analyzing Image Segmentation for Connectomics. *Frontiers in Neural Circuits*, 12.
- Scheffer, L. K., Xu, C. S., Januszewski, M., Lu, Z., Takemura, S.-y., Hayworth, K. J., Huang, G. B., Shinomiya, K., Maitin-Shepard, J., Berg, S., Clements, J., Hubbard, P., Katz, W., Umayam, L., Zhao, T., Ackerman, D., Blakely, T., Bogovic, J., Dolafi, T., Kainmueller, D., Kawase, T., Khairy, K. A., Leavitt, L., Li, P. H., Lindsey, L., Neubarth, N., Olbris, D. J., Otsuna, H., Trautman, E. T., Ito, M., Goldammer, J., Wolff, T., Svirskas, R., Schlegel, P., Neace, E. R., Knecht, C. J., Alvarado, C. X., Bailey, D. A., Ballinger, S., Borycz, J. A., Canino, B. S., Cheatham, N., Cook, M., Dreher, M., Duclos, O., Eubanks, B., Fairbanks, K., Finley, S., Forknall, N., Francis, A., Hopkins, G. P., Joyce, E. M., Kim, S., Kirk, N. A., Kovalyak, J., Lauchie, S. A., Lohff, A., Maldonado, C., Manley, E. A., McLin, S., Mooney, C., Ndama, M., Ogundeyi, O., Okeoma, N., Ordish, C., Padilla, N., Patrick, C., Paterson, T., Phillips, E. E., Phillips, E. M., Rampally, N., Ribeiro, C., Robertson, M. K., Rymer, J. T., Ryan, S. M., Sammons, M., Scott, A. K., Scott, A. L., Shinomiya, A., Smith, C., Smith, K., Smith, N. L., Sobeski, M. A., Suleiman, A., Swift, J., Takemura, S., Talebi, I., Tarnogorska, D., Tenshaw, E., Tokhi, T., Walsh, J. J., Yang, T., Horne, J. A., Li, F., Parekh, R., Rivlin, P. K., Jayaraman, V., Ito, K., Saalfeld, S., George, R., Meinertzhagen, I. A., Rubin, G. M., Hess, H. F., Jain, V., and Plaza, S. M. (2020). A Connectome and Analysis of the Adult Drosophila Central Brain. bioRxiv, page 2020.04.07.030213. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Schlegel, P., Texada, M. J., Miroschnikow, A., Schoofs, A., Hückesfeld, S., Peters, M., Schneider-Mizell, C. M., Lacin, H., Li, F., Fetter, R. D., Truman, J. W., Cardona, A., and Pankratz, M. J. (2016). Synaptic transmission parallels neuromodulation in a central food-intake circuit. *eLife*, 5:e16799. Publisher: eLife Sciences Publications, Ltd.
- Schneider-Mizell, C. M., Bodor, A. L., Collman, F., Brittain, D., Bleckert, A. A., Dorkenwald, S., Turner, N. L., Macrina, T., Lee, K., Lu, R., Wu, J., Zhuang, J., Nandi, A., Hu, B., Buchanan, J., Takeno, M. M., Torres, R., Mahalingam, G., Bumbarger, D. J., Li,

- Y., Chartrand, T., Kemnitz, N., Silversmith, W. M., Ih, D., Zung, J., Zlateski, A., Tartavull, I., Popovych, S., Wong, W., Castro, M., Jordan, C. S., Froudarakis, E., Becker, L., Suckow, S., Reimer, J., Tolias, A. S., Anastassiou, C., Seung, H. S., Reid, R. C., and Costa, N. M. d. (2020). Chandelier cell anatomy and function reveal a variably distributed but common signal. *bioRxiv*, page 2020.03.31.018952. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Schneider-Mizell, C. M., Gerhard, S., Longair, M., Kazimiers, T., Li, F., Zwart, M. F., Champion, A., Midgley, F. M., Fetter, R. D., Saalfeld, S., and Cardona, A. (2016). Quantitative neuroanatomy for connectomics in Drosophila. *eLife*, 5:e12059. Publisher: eLife Sciences Publications, Ltd.
- Takemura, S.-y., Xu, C. S., Lu, Z., Rivlin, P. K., Parag, T., Olbris, D. J.,
  Plaza, S., Zhao, T., Katz, W. T., Umayam, L., Weaver, C., Hess,
  H. F., Horne, J. A., Nunez-Iglesias, J., Aniceto, R., Chang, L.-A.,
  Lauchie, S., Nasca, A., Ogundeyi, O., Sigmund, C., Takemura, S.,
  Tran, J., Langille, C., Le Lacheur, K., McLin, S., Shinomiya, A.,
  Chklovskii, D. B., Meinertzhagen, I. A., and Scheffer, L. K. (2015).
  Synaptic circuits and their variations within different columns in the
  visual system of Drosophila. Proceedings of the National Academy
  of Sciences of the United States of America, 112(44):13711–13716.
- Tu, Z. and Bai, X. (2010). Auto-Context and Its Application to High-Level Vision Tasks and 3D Brain Image Segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(10):1744–1757. Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- Turaga, S. C., Briggman, K. L., Helmstaedter, M., Denk, W., and Seung, H. S. (2009). Maximin affinity learning of image segmentation. arXiv:0911.5372 [cs]. arXiv: 0911.5372.
- Turaga, S. C., Murray, J. F., Jain, V., Roth, F., Helmstaedter, M., Briggman, K., Denk, W., and Seung, H. S. (2010). Convolutional Networks Can Learn to Generate Affinity Graphs for Image Segmentation. *Neural Computation*, 22(2):511–538. Conference Name: Neural Computation.
- Turner, N. L., Macrina, T., Bae, J. A., Yang, R., Wilson, A. M., Schneider-Mizell, C., Lee, K., Lu, R., Wu, J., Bodor, A. L., Bleckert, A. A., Brittain, D., Froudarakis, E., Dorkenwald, S., Collman, F., Kemnitz, N., Ih, D., Silversmith, W. M., Zung, J., Zlateski, A., Tartavull, I., Yu, S.-c., Popovych, S., Mu, S., Wong, W., Jordan, C. S., Castro, M., Buchanan, J., Bumbarger, D. J., Takeno, M., Torres, R., Mahalingam, G., Elabbady, L., Li, Y., Cobos, E., Zhou, P., Suckow, S., Becker, L., Paninski, L., Polleux, F., Reimer, J., Tolias, A. S., Reid, R. C., Costa, N. M. d., and Seung, H. S. (2020). Multiscale and multimodal reconstruction of cortical structure and function. *bioRxiv*, page 2020.10.14.338681. Publisher: Cold Spring Harbor Laboratory Section: New Results.
- Turner-Evans, D. B. and Jayaraman, V. (2016). The insect central complex. *Current Biology*, 26(11):R453–R457.
- Turner-Evans, D. B., Jensen, K. T., Ali, S., Paterson, T., Sheridan, A., Ray, R. P., Wolff, T., Lauritzen, J. S., Rubin, G. M., Bock, D. D., and Jayaraman, V. (2020). The Neuroanatomical Ultrastructure and Function of a Biological Ring Attractor. *Neuron*, 0(0). Publisher: Elsevier.
- Wolf, S., Pape, C., Bailoni, A., Rahaman, N., Kreshuk, A., Köthe, U., and Hamprecht, F. A. (2018). The Mutex Watershed: Efficient, Parameter-Free Image Partitioning. In Ferrari, V., Hebert, M., Sminchisescu, C., and Weiss, Y., editors, Computer Vision ECCV 2018, volume 11208, pages 571–587. Springer International Publishing, Cham. Series Title: Lecture Notes in Computer Science.
- Yin, W., Brittain, D., Borseth, J., Scott, M. E., Williams, D., Perkins, J., Own, C. S., Murfitt, M., Torres, R. M., Kapner, D., Mahalingam, G., Bleckert, A., Castelli, D., Reid, D., Lee, W.-C. A., Graham, B. J., Takeno, M., Bumbarger, D. J., Farrell, C., Reid, R. C., and da Costa, N. M. (2020). A petascale automated imaging pipeline for mapping neuronal circuits with high-throughput transmission electron microscopy. *Nature Communications*, 11(1):4949. Number: 1 Publisher: Nature Publishing Group.
- Zheng, Z., Lauritzen, J. S., Perlman, E., Robinson, C. G., Nichols, M., Milkie, D., Torrens, O., Price, J., Fisher, C. B., Sharifi, N.,

Calle-Schuler, S. A., Kmecova, L., Ali, I. J., Karsh, B., Trautman, E. T., Bogovic, J. A., Hanslovsky, P., Jefferis, G. S. X. E., Kazhdan, M., Khairy, K., Saalfeld, S., Fetter, R. D., and Bock, D. D. (2018). A Complete Electron Microscopy Volume of the Brain of Adult Drosophila melanogaster. *Cell*, 174(3):730–743.

# A LOCAL SHAPE DESCRIPTORS

We define the notational shorthand

$$\mathbf{f}_{k}^{i}(v) = v_{k} \, \mathbf{b}^{i}(v) \qquad \qquad k \in \{x, y, z\} \tag{9}$$

$$f_{kl}^{i}(v) = v_k v_l b^{i}(v)$$
  $k, l \in \{x, y, z\},$  (10)

and use those to rewrite (4) as follows:

$$\mathbf{m}_{k}^{i}(v) = (\mathbf{f}_{k}^{i} * \mathbf{w})(v)/\mathbf{s}^{i}(v)$$
  $k \in \{x, y, z\}$  (11)

$$c_{kl}^{i}(v) = (f_{kl}^{i} * w)(v)/s^{i}(v) - m_{k}^{i}(v) m_{l}^{i}(v) \quad k, l \in \{x, y, z\}.$$
 (12)

It can easily be seen that (11) is equal to the local center of mass, limited both by the object mask b and the local window w:

$$\begin{split} \mathbf{m}_{k}^{i}(v) &= (\mathbf{f}_{k}^{i} * \mathbf{w})(v) / \mathbf{s}^{i}(v) \\ &= \frac{1}{\mathbf{s}^{i}(v)} \sum_{v' \in \Omega} \mathbf{f}_{k}^{i}(v') \, \mathbf{w}(v - v') \\ &= \frac{1}{\mathbf{s}^{i}(v)} \sum_{v' \in \Omega} v_{k}' \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') \end{split}$$

Similarly, the computation of the local covariance of voxel coordinates is equivalent to a convolution of the local window w with  $f_{kl}^i(v)$ . The local covariance is defined as:

$$\begin{split} \mathbf{c}_{kl}^{i}(v) &= \frac{1}{\mathbf{s}^{i}(v)} \sum_{v' \in \Omega} (v_{k}' - \bar{v}_{k})(v_{l}' - \bar{v}_{l}) \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') \\ &= \frac{1}{\mathbf{s}^{i}(v)} \sum_{v' \in \Omega} (v_{k}' - \mathbf{m}_{k}^{i}(v))(v_{l}' - \mathbf{m}_{l}^{i}(v)) \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') \\ &= \frac{1}{\mathbf{s}^{i}(v)} \sum_{v' \in \Omega} \left( v_{k}' v_{l}' - v_{k}' \, \mathbf{m}_{k}^{i}(v) - v_{l}' \, \mathbf{m}_{k}^{i}(v) + \mathbf{m}_{k}^{i}(v) \, \mathbf{m}_{l}^{i}(v) \right) \cdot \\ &\cdot \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') \end{split}$$

Rearranging terms reveals that  $c_{kl}^i(v)$  can efficiently be computed via a convolution as well:

$$\begin{split} \mathbf{c}_{kl}^{i}(v) &= \frac{1}{\mathbf{s}^{i}(v)} \Bigg( \sum_{v' \in v} v_{k}' v_{l}' \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') - \\ & \mathbf{m}_{l}^{i}(v) \sum_{v' \in v} v_{k}' \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') - \\ & \mathbf{m}_{k}^{i}(v) \sum_{v' \in v} v_{l}' \, \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') + \\ & \mathbf{m}_{k}^{i}(v) \, \mathbf{m}_{l}^{i}(v) \sum_{v' \in v} \mathbf{b}^{i}(v') \, \mathbf{w}(v - v') \Bigg) \end{split}$$

$$\begin{split} c_{kl}^{i}(v) = & \frac{1}{s^{i}(v)} \Biggl( \sum_{v' \in v} \underbrace{v_{k}' v_{l}' \, b^{i}(v')}_{f_{kl}^{i}(v')} \, w(v - v') - \\ & m_{l}^{i}(v) \underbrace{(f_{k}^{i} * w)(v) - }_{m_{k}^{i}(v) \, s^{i}(v)} \\ & m_{k}^{i}(v) \underbrace{(f_{l}^{i} * w)(v) + }_{m_{l}^{i}(v) \, s^{i}(v)} \\ & m_{k}^{i}(v) \, m_{l}^{i}(v) \, s^{i}(v) \Biggr) \\ = & \frac{1}{s^{i}(v)} \Biggl( (f_{kl}^{i} * w)(v) - s^{i}(v) \, m_{l}^{i}(v) \, m_{k}^{i}(v) \Biggr) \\ = & (f_{kl}^{i} * w)(v) / s^{i}(v) - m_{l}^{i}(v) \, m_{k}^{i}(v) \end{split}$$

#### **B** MIN-CUT METRIC

The Min-Cut Metric (MCM) measures the number of edit operations that need to be performed by a human annotator in a hypothetical proofreading tool that allows to: (1) merge wrongly split segments and (2) split wrongly merged segments by means of a min-cut15. To this end, we assume that the segmentation to interact with results from an agglomeration of fragments (or "supervoxels"). In particular, we assume that a fragment graph G = (V, E, s) is available, where each node  $v \in V$  corresponds to a fragment and edges  $(u, v) \in E$  are introduced between neighboring fragments. Each edge  $e \in E$  has an associated merge score s(e), which denotes under which agglomeration threshold the two incident fragments are to be merged into the same segment. A segmentation of the fragment graph is induced by a merge score threshold  $\theta$ . Let  $E_{\theta} = \{e \in E \mid s(e) \leq \theta\}$  be the set of *filtered* edges. Each connected component in the graph  $G_{\theta} = (V, E_{\theta})$  then corresponds to one segment. We will refer to the segment ID of a fragment under a given threshold  $\theta$  as  $l_{\theta}(v)$ .

For the MCM metric, we assume ground-truth is available in the form of *skeletons*. Let T be the set of ground-truth skeletons, with each  $t \in T$  being a set of skeleton nodes. We will refer to the skeleton ID of a skeleton node a as s(a) and the fragment underlying the skeleton node as  $l_{\theta}(a)$ .

Given a segmentation  $l_{\theta}$ , MCM first simulates the splitting of all wrongly merged structures as given by ground-truth skeletons. For that, we first identify *merging* segments, *i.e.*, segments that contain nodes from more than one skeleton. For each merging segment, we iteratively perform a series of min-cuts through the fragment graph, until the skeletons are separated. For that, we repeatedly find a pair of skeleton nodes a and b, such that

- 1)  $s(a) \neq s(b)$  (the skeleton nodes belong to different skeletons),
- 2)  $l_{\theta}(a) = l_{\theta}(b)$  (the underlying fragments belong to the same segment), and
- 3) the Euclidean distance between a and b is minimized.

We then perform a min-cut on the fragment graph, with u and v as the source and sink, respectively, and the capacity c(e) of edges e in the fragment graph set proportional to -s(e), such that edges with a high merge score are cheaper to cut. Once the min-cut is found, all edges of the cut are removed from  $E_{\theta}$  and the segmentation is updated accordingly. For a visualization of this procedure, see Supplemental Fig. 22.A. This procedure aims to mimic a proofreader, who identified a merge and consequently picked two close locations on either side of the merge to perform a split operation.

In some cases, a min-cut can fail to separate all nodes of the merged skeletons with a single cut from two selected nodes. In this case, the procedure is repeated until all skeletons are separated, leading to additional split errors (see Supplemental Fig. 22.B for an example).

After all skeletons are separated, the remaining split errors are counted. For that, we assume that each split requires one merge operation to be fixed. More generally, we identify the segments underlying each skeleton t: Let  $L(t) = \{l_{\theta}(a) \mid a \in t\}$  be the set of segments underlying a skeleton. The number of required merge operations is then recorded as 1 - |L(t)|.

<sup>15</sup> Also referred to as "cleaving".

# C ZEBRAFINCH

# C.1 Training

## C.1.1 Data

33 volumes of densely labeled neurons<sup>6</sup> were used for training. Each volume was padded with raw data. 30 volumes had raw dimensions of  $\sim$ 7, 4.95, 4.95µm (zyx) and label dimensions of  $\sim$ 3, 1.35µm. The remaining 3 volumes had raw dimensions of  $\sim$ 6.6, 5.9, 5.9µm and label dimensions of  $\sim$ 2.6, 2.3, 2.3µm. Some regions containing glial processes were already set to zero and incorporated during network training (Supplemental Fig. 17.A). A labels mask (1 inside labels RoI, null outside) was generated and used during training.

#### C.1.2 Networks

All methods used the U-Net architecture described in Funke et al. (2019). Networks consisted of three layers and were downsampled by a factor of [1,3,3] in the first two layers and [3,3,3] in the last layer. The reverse was done for the upsampling path. 12 initial feature maps were used and features were multiplied by a factor of 5 between layers. The resulting data was further convolved and passed through a sigmoid activation to get from 12 output feature maps <sup>16</sup> to either 3 (affinities) or 10 feature maps (LSDs). All networks used an MSE loss, minimized with an Adam optimizer. The Malis network was trained to 10k iterations using MSE to initialize affinities and was then switched to Malis loss for the remainder of training.

Non auto-context networks had an input shape (raw) of [84,268,268] and output shape (labels, LSDs, affinities) of [48,56,56] (voxels, zyx). Auto-context networks had an input shape (raw) of [120,484,484], an intermediate shape (predicted LSDs) of [84,268,268], and an output shape (labels, affinities) of [48,56,56] (see Supplemental Fig. 19.A for visualization of auto-context training shapes on Fib-25). The predicted LSDs used in the intermediate shape were taken from a pre-trained network which predicted LSDs from raw. Non auto-context networks were trained to 400k iterations. Auto-context networks were trained to ~200k iterations following the 400k iterations of LSD training. See Supplemental Table 3 for a breakdown of the MTLsD network as an example.

All networks used a single voxel affinity neighborhood [1,1,1]. The LR network used three additional neighborhood steps of [3,3,3], [5,9,9] and [15,27,27]. The computed LSDs used a sigma of 120 nm and a downsampling factor of two.

# C.1.3 Pipeline

Each training batch was randomly picked from one of the 33 training volumes. For each batch, the raw data was first normalized and padded with zeros. Labels were padded with the maximum padding required to contain at least 50% of ground-truth data assuming a worst case rotation of 45°. Data was randomly sampled from each dataset using a labels mask to ensure every batch contained at least 50% of ground-truth data. Data was then augmented with elastic transformations, random mirrors + transposes, and intensities (see Supplemental Table 3 for augmentation hyperparameters used for example MrLsp network). The following was done to the respective networks:

 $^{16}\mbox{MrLsp}$  network had 14 output feature maps to account for 13 final feature maps from the affinities (3) and LSDs (10)

- BASELINE, LR Label boundaries were first eroded by a single voxel. Ground truth affinities were calculated on the labels using the pre-defined affinity neighborhoods and a scale array was created to balance loss between class labels. Training: [raw + gt affs] → pred affs.
- MALIS Label boundaries were eroded. If training loss was in MALIS phase (*i.e.* after 10k iterations), connected components were relabelled before calculating ground-truth affinities. If training loss was in MSE phase (*i.e.* before 10k iterations), labels were subsequently balanced. Training: [raw + gt affs] → pred affs.
- LSDs Ground truth LSDs were calculated on the labels using the pre-defined sigma and downsampling factor. Training:
   [raw + gt LSDs ] → pred LSDs.
- MTLsD Label boundaries were eroded. Ground truth LSDs were calculated followed by ground-truth affinities. Labels were then balanced. Training: [raw + gt LSDs + gt affs] → [pred LSDs + pred affs].
- AcLsd, AcRLsd Label boundaries were eroded, ground-truth affinities were calculated, and labels were balanced. LSDs were then predicted in a slightly larger region and used as input to train the affinities. Training: [raw + gt LSDs ] → pred LSDs → pred affs. For AcRLsd, cropped raw was incorporated in the second pass, in addition to predicted LSDs.

This process was repeated for a pre-defined number of iterations (generally until loss convergence).

## C.2 Prediction

Prediction was done in a block-wise fashion restricted to the BENCHMARK Roi. Individual workers used Gunpowder<sup>4</sup> to predict output data (i.e. affinities or LSDs) and were distributed throughout the volume with DaisyNguyen et al. (2020). Block size was chosen with respect to how much data could fit in GPU memory. Most networks had a smaller block size in order to fit on 2080 RTX GPUs (~12 GB RAM). While this increased the total number of blocks to process, the amount of workers available to use was sufficient to minimize total processing time. The autocontext networks were too large to fit on 2080 RTX GPUs and were therefore run on Tesla V100 GPUs. While there were less V100's available, the block size could be greatly increased (~32) GB RAM), decreasing the total amount of blocks to process. LSDs were physically written to file before use in the auto-context networks. This was done for visualization; prediction could be adapted to generate LSDs on the fly. All networks wrote affinities to file and then subsequently used them for segmentation.

#### C.3 Segmentation

## C.3.1 Watershed

Seeded watershed<sup>17</sup> was done on the affinities generated during prediction. Both non-masked and neuropil-masked supervoxels were produced. Due to data anisotropy, supervoxels were extracted for each section separately. An epsilon agglomeration was used to agglomerate fragments to a predefined threshold (0.1). This was done to decrease the number of RAG nodes during the full agglomeration step. Supervoxels which had an average affinity

value lower than a pre-defined threshold (0.05) were filtered out of the RAG and set to zero in the resulting datasets. A block size of  $3.6\mu\text{m}^3$  and context of [12,27,27] voxels (zyx) were used.

# C.3.2 Agglomeration

Supervoxels were agglomerated using hierarchical region agglomeration 17 in which edges with lower affinity scores are merged earlier. We empirically chose to use both 50 and 75 quantile merge functions since they produced the best results in Funke et al. (2019). The same block size/context as watershed were used.

#### C.3.3 Segment

The center point of the Benchmark Roi was used to grow sub RoIs. The first sub RoI was created by growing the center point in each dimension (positive and negative) by the block size used during watershed and agglomeration. This resulted in  $10.8\mu m$  edge lengths  $(3.6\mu m + (3.6 \text{ x }2))$ . This RoI was again grown by the block size to produce an RoI with  $18\mu m$  edge lengths  $(10.8\mu m + (3.6 \text{ x }2))$ . This was repeated for a total of 10 RoIs (in addition to the Benchmark Roi). Segmentations were created for each RoI by cropping the RAG and relabelling connected components on the graph  $^{18}$ . This was done over a range of thresholds for each network (threshold range = [0 - 1], step size = 0.02, total thresholds = 50). Segmentations were created for both masked / non-masked data and both merge functions.

#### C.4 Evaluation

Manually traced skeletons<sup>6</sup> (12 validation, 50 testing) were used for evaluation. For each sub RoI, skeletons were cropped, either masked to neuropil or not masked, and connected components were relabelled<sup>18</sup>. For affinity-based methods, fragment ids were first mapped to skeleton ids in each block. This mapping was then used to assign segment ids to skeleton ids for each threshold, using the lookup tables generated in Supplemental Section C.3.3. A site mask was used to restrict segments to the skeleton nodes. The resulting node - segment mapping was used to compute<sup>19</sup> ERL, NERL and VoI.

Additionally, on the first three sub RoIs, the MCM was calculated using masked skeletons and segmentations. For the FFN, a single segmentation<sup>6</sup> was used to generate the node - segment mappings. The full segmentation was downloaded<sup>20</sup> and cropped to each sub RoI, either masked or not masked, and connected components were relabelled. Only ERL, NERL and VoI were calculated as there were no supervoxels to use for the MCM. For all affinity-based methods, we repeated these steps using the validation skeletons on the benchmark RoI. The optimal thresholds indicated in test set plots were determined by the thresholds which minimized VoI (for VoI and MCM plots) and maximized ERL (for ERL plots).

# D FIB-SEM VOLUMES

# D.1 Training

#### D.1.1 Data

**HEMI-BRAIN**: 8 volumes of densely labeled neurons<sup>7</sup> were used for training. Volumes were taken from various neuropils<sup>21</sup> contained

within the dataset generated in (Scheffer et al., 2020). The Lobula Plate and Lateral Horn volumes contained  ${\sim}4\mu\text{m}^3$  of raw data and  ${\sim}2\mu\text{m}^3$  of labeled data, while the others contained  ${\sim}6\mu\text{m}^3$  of raw data and  ${\sim}4\mu\text{m}^3$  of labeled data (Supplemental Fig. 17.B).

**FIB-25**: 4 volumes of densely labeled neurons<sup>7</sup> were used for training. The labels were not padded with raw as was done in the Zebrafinch and Hemi-brain volumes. Two volumes contained  $\sim 4 \mu \text{m}^3$  of raw / labeled data, and two volumes contained  $2 \mu \text{m}^3$  of raw / labeled data (Supplemental Fig. 17.C). Label masks were generated for all volumes, as done in the Zebrafinch.

#### D.1.2 Networks

Networks consisted of same architecture as Zebrafinch networks except downsampling was isotropic with a factor of [2,2,2] in the first two layers and [3,3,3] in the last layer. Features were multiplied by a factor of 6 between layers.

Non auto-context networks had an input shape (raw) of [196,196,196] and output shape (labels, LSDs, affinities) of [92,92,92]. Auto-context networks had an input shape (raw) of [304,304,304], an intermediate shape (predicted LSDs) of [196,196,196], and an output shape (labels, affinities) of [92,92,92]. Non auto-context networks were trained to 400k iterations. Auto-context networks were trained to ~300k iterations following 400k iterations of LSD training. See Supplemental Table 4 for a breakdown of the MTLsD network as an example.

All networks used a single voxel affinity neighborhood [1,1,1]. The LR network used three additional neighborhood steps of [3,3,3], [5,5,5] and [13,13,13]. The computed LSDs used a sigma of 80 nm and a downsampling factor of two.

#### D.1.3 Pipeline

All networks were trained following the same pipeline as the Zebrafinch networks using either 8 (Hemi-brain) or 4 (Fib-25) ground-truth volumes. The augmentations were computed isotropically, in contrast to the Zebrafinch networks (see Supplemental Table 4 for augmentation hyper-parameters used for example MtLsd network). For Fib-25 training, affinities and LSDs were masked at the boundaries to ensure that prediction on the irregularly shaped Fib-25 volume did not include boundary artifacts (Supplemental Fig. 19.B).

# **D.2** Prediction

For the Hemi-brain, prediction was restricted to the three RoIs described in Section 3.3.2. For Fib-25, prediction was done on the full Fib-25 volume (including the background). The process was the same as in the Zebrafinch.

#### **D.3 Segmentation**

## D.3.1 Watershed

Fragment extraction was performed isotropically and used no epsilon agglomeration step or mean affinity filtering, in contrast to the Zebrafinch. A block size of 3µm³ and context of 31 voxels were used. For the Hemi-brain, watershed was done on each predicted RoI and restricted using an Ellipsoid Body mask³. For Fib-25, an irregularly shaped tissue mask6 was used22.

<sup>18</sup>https://github.com/funkelab/funlib.segment

<sup>19</sup>https://github.com/funkelab/funlib.evaluate

<sup>&</sup>lt;sup>20</sup>https://github.com/seung-lab/cloud-volume

 $<sup>^{21}\</sup>mbox{Ellipsoid}$  Body: 2, Protocerebral Bridge: 2, Fan-Shaped Body: 2, Lobula Plate: 1, Lateral Horn: 1

<sup>&</sup>lt;sup>22</sup>Mask contains some background and cell bodies

# D.3.2 Agglomeration

Agglomeration was done using the same merge functions on the Zebrafinch. The same block size and context from watershed were used.

## D.3.3 Segment

For the Hemi-brain, segmentations were created for the three processed RoIs. For Fib-25, segmentations were created for the full Fib-25 RoI and two sub RoIs. The same threshold range from the Zebrafinch was used.

#### D.4 Evaluation

HEMI-BRAIN: dense ground-truth<sup>7</sup> and an FFN segmentation<sup>6</sup> were available for the entire Hemi-brain. Both volumes were downloaded<sup>23,24</sup> and cropped to the three established RoIs. The cropped datasets were then constrained to the Ellipsoid Body. The ground-truth was filtered using a whitelist of proofread ids. Connected components were relabelled and boundaries were slightly eroded. Fib-25: dense ground-truth<sup>7</sup> and an FFN segmentation<sup>6</sup> were already cropped to the testing RoI. The ground-truth was already filtered with a whitelist and boundaries were already eroded. Both volumes were further cropped to the two sub RoIs and connected components were relabelled. For affinity-based methods, VoI was calculated between the consolidated ground-truth and segmentations over all thresholds on each RoI. For the FFN, VoI was calculated on the single segmentation for each RoI.

# **E** THROUGHPUT

For each affinity-based network, we calculated the amount of floating point operations (FLOPs) for the processing of one block  $^{25}$  using TensorFlow's Profiler  $^{26}$  (see Supplemental Table 5 for a breakdown by operation). From the computed FLOPs and the block size, we derived FLOPs/ $\mu m$ .

For FFN, FLOPs are reported for the full Zebrafinch RoI in Januszewski et al. (2018), which we divided by the full RoI volume to get FLOPs/ $\mu m$ .

<sup>&</sup>lt;sup>23</sup>https://github.com/janelia-flyem/dvid

<sup>&</sup>lt;sup>24</sup>https://github.com/janelia-flyem/neuclease

<sup>&</sup>lt;sup>25</sup>One block is defined as the largest output volume that can be predicted by a network in one pass on the respective GPU it was evaluated on.

<sup>&</sup>lt;sup>26</sup>https://www.tensorflow.org/guide/profiler

Parameter	Value
Input feature maps	12
Layer fmap scale	5
Downsampling factors	[[1,3,3],[1,3,3],[3,3,3]]
Output feature maps	14
Input shape	[84,268,268]
Output shape	[48,56,56]
Loss	MSE
Optimizer	Adam
Learning rate	$0.5 \times 10^{-4}$
$\beta_1$	0.95
$\beta_2$	0.999
$\epsilon$	$1 \times 10^{-8}$
Iterations	400,000

Augmentation	Parameter	Value
Elastic	control point spacing	(4,4,10)
	jitter sigma	(0, 2, 2)
	subsample	8
Rotation	axis	x,y,z
	angle	in $[0, 2\pi]$
Section Defects	slip probability	0.05
	shift probability	0.05
	max misalign	10
Mirror	axes	x,y,z
Transpose	axes	x, y
Intensity	scale	in [0.9, 1.1]
	shift	in [-0.1, 0.1]

Table 3: Training parameters and augmentations<sup>4</sup> of MTLsD network on Zebrafinch dataset.

Parameter	Value
Input feature maps	12
Layer fmap scale	6
Downsampling factors	[[2,2,2],[2,2,2],[3,3,3]]
Output feature maps	14
Input shape	[196,196,196]
Output shape	[92,92,92]
Loss	MSE
Optimizer	Adam
Learning rate	$0.5 \times 10^{-4}$
$\beta_1$	0.95
$oldsymbol{eta}_2$	0.999
$\epsilon$	$1 \times 10^{-8}$
Iterations	400,000

Augmentation	Parameter	Value
Elastic 1	control point spacing	(40,40,40)
	jitter sigma	(0, 0, 0)
	subsample	8
Rotation 1	axis	x,y,z
	angle	in $[0, 2\pi]$
Section Defects 1	slip probability	0
	shift probability	0
	max misalign	0
Mirror	axes	x,y,z
Transpose	axes	x,y,z
Elastic 2	control point spacing	(40,40,40)
	jitter sigma	(2, 2, 2)
	subsample	8
Rotation 2	axis	x,y,z
	angle	in $[0, 2\pi]$
Section Defects 2	slip probability	0.01
	shift probability	0.01
	max misalign	1
Intensity	scale	in [0.9, 1.1]
	shift	in [-0.1, 0.1]

Table 4: Training parameters and augmentations<sup>4</sup> of MTLsD network on FIB-SEM datasets.

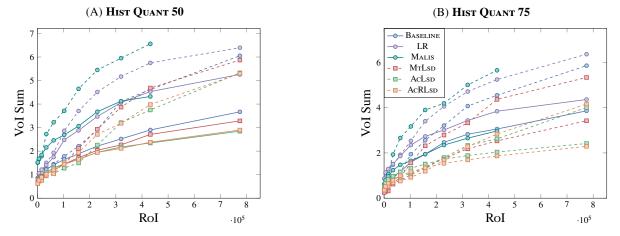


Figure 11: VoI Sum vs RoI on masked (solid) and non-masked (dashed) Zebrafinch data.

Operation	FLOPs	Operation
gradients/mean_squared_error_1/S	940584960000	unet_layer_2_left_1/convolution
mean_squared_e	786542400000	unet_layer_2_right_0/convolution
gradients/mean_squared_err	419904000000	unet_layer_3_left_1/convolution
gradients/mean_squared_e	416320819200	unet_layer_1_left_1/convolution
affs_0/Bia	338411520000	unet_layer_2_right_1/convolution
gradients/mean_squared_error_1/	226738828800	unet_layer_1_right_0/convolution
mean_squared_e	208980000000 164826316800	unet_layer_2_left_0/convolution unet_layer_0_left_1/convolution
gradients/affs_0/BiasAdd	123832800000	unet_layer_3_left_0/convolution
unet_layer_2_right_0/kernel/In	105305702400	unet_layer_1_right_1/convolution
unet_layer_3_left gradients/unet_layer_3_left_0/B	87354028800	unet_layer_1_left_0/convolution
unet_layer_3_left	84455094528	unet_layer_0_right_0/convolution
gradients/unet_layer_3_left_1/B	83980800000	gradients/unet_up_3_to_2/conv3d_transpose_grad/Conv3D
unet_layer_2_right_1/kernel/In	46982799360	unet_layer_0_right_1/convolution
unet_layer_2_left_1/kernel/Ini	22560768000	gradients/unet_up_2_to_1/conv3d_transpose_grad/Conv3D
unet_layer_2_left_0/kernel/Ini	14151319488	unet_layer_0_left_0/convolution
unet_layer_1_right_0/kernel/In	7020380160	gradients/unet_up_1_to_0/conv3d_transpose_grad/Conv3D
unet_up_2_to_1/kernel/Initia	1242931200	embedding_0/convolution
unet_layer_1_right_1/kernel/In	372879360	affs_0/convolution
unet_layer_1_left_1/kernel/Ini	262061472	unet_layer_0_left_0/BiasAdd
unet_layer_1_left_0/kernel/Ini	262061460	gradients/unet_layer_0_left_0/BiasAdd_grad/BiasAddGrad
unet_layer_0_right_0/kernel/In	254361600 254361600	gradients/AddN_3
unet_up_1_to_0/kernel/Initi	254361588	unet_layer_0_left_1/BiasAdd gradients/unet_layer_0_left_1/BiasAdd_grad/BiasAddGrad
unet_layer_0_right_1/kernel/In	134805600	unet_layer_1_left_0/BiasAdd
unet_layer_0_left_1/kernel/Ini	134805540	gradients/unet_layer_1_left_0/BiasAdd_grad/BiasAddGrad
unet_layer_0_left_0/kernel/Ini	128494080	unet_layer_1_left_1/BiasAdd
embedding_0/kernel/Initial	128494080	gradients/AddN_2
affs_0/kernel/Initializer	128494020	gradients/unet_layer_1_left_1/BiasAdd_grad/BiasAddGrad
gradients/Slice_	60750000	unet_layer_3_left_1/kernel/Initializer/random_uniform
gradients/Slice_1 gradients/Slice_1	44390399	mean_squared_error/num_present
gradients/Slice_2	88780800	mean_squared_error/SquaredDifference
gradients/Slice_z	65165968	unet_layer_0_right_0/BiasAdd
gradients/Slice	65165954	gradients/unet_layer_0_right_0/BiasAdd_grad/BiasAddGrad
mean_squared_err	64500000	unet_layer_2_left_0/BiasAdd
unet_up_1_to	64499700	gradients/unet_layer_2_left_0/BiasAdd_grad/BiasAddGrad
unet_up_3_to_	62146560	unet_layer_0_right_1/BiasAdd
unet_up_3_to_	62146560	gradients/AddN
unet_up_3_to	62146546 58503168	gradients/unet_layer_0_right_1/BiasAdd_grad/BiasAddGrad unet_up_1_to_0/BiasAdd
unet_up_3_to_	58503156	gradients/unet_up_1_to_0/BiasAdd_grad/BiasAddGrad
unet_up_3_to_	58060800	unet_layer_2_left_1/BiasAdd
unet_up_3_to	58060800	gradients/AddN_1
unet_up_2_to_	58060500	gradients/unet_layer_2_left_1/BiasAdd_grad/BiasAddGrad
unet_up_2_to_	44390400	embedding_0/BiasAdd
unet_up_2_to	44390400	gradients/mean_squared_error/Mul_grad/mul
unet_up_2_to_	44390400	gradients/mean_squared_error/Mul_grad/mul_1
unet_up_2_to_	44390400	gradients/mean_squared_error/SquaredDifference_grad/Neg
unet_up_2_to	44390400	gradients/mean_squared_error/SquaredDifference_grad/mul
unet_up_1_to_	44390400	mean_squared_error/Mul
unet_up_1_to_ gradients/mean_squared_erro	44390400	gradients/mean_squared_error/SquaredDifference_grad/mul_1
Adam/mı	44390400	gradients/mean_squared_error/SquaredDifference_grad/sub
add	44390399	mean_squared_error/Sum
gradients/mean_squared_	44390390	gradients/embedding_0/BiasAdd_grad/BiasAddGrad
gradients/mean_squared_er	39951360	gradients/mean_squared_error/Mul_grad/Sum_1
gradients/mean_squared_err	37601280	unet_up_2_to_1/BiasAdd gradients/unet_up_2_to_1/BiasAdd_grad/BiasAddGrad
gradients/mean squared err	37601220 34990560	unet_layer_1_right_0/BiasAdd
gradients/mean_squared_	34990500	gradients/unet_layer_1_right_0/BiasAdd_grad/BiasAddGrad
gradients/mean_squared_e	32501760	unet_layer_1_right_1/BiasAdd
gradients/mean_squared_err	32501700	gradients/unet_layer_1_right_1/BiasAdd_grad/BiasAddGrad
gradients/mean_squared_erro	27993600	unet_up_3_to_2/BiasAdd
unet_up_1_to_	27993300	gradients/unet_up_3_to_2/BiasAdd_grad/BiasAddGrad
gradients/mean_squared_e	26634240	mean_squared_error_1/SquaredDifference
mean_squared_e	13317119	mean_squared_error_1/num_present
mean_squared_e	12150000	unet_layer_3_left_0/kernel/Initializer/random_uniform
mean_squared_	12150000	unet_up_3_to_2/kernel/Initializer/random_uniform
mean_squared_er	24276000	unet_layer_2_right_0/BiasAdd
Adam/n	24275700	gradients/unet_layer_2_right_0/BiasAdd_grad/BiasAddGrad
mean_squared_e	20889600	unet_layer_2_right_1/BiasAdd
unet_up_1_to	20889300	gradients/unet_layer_2_right_1/BiasAdd_grad/BiasAddGrad
unet_up_1_to_ Total	13317120	gradients/mean_squared_error_1/SquaredDifference_grad/mul gradients/mean_squared_error_1/SquaredDifference_grad/Neg
	13317120	aramenismean squared error (/Squared)litterence grad/Neg

Operation	FLOPs
gradients/mean_squared_error_1/SquaredDifference_grad/mul_1	13317120
mean_squared_error_1/Mul	13317120
gradients/mean_squared_error_1/Mul_grad/mul_1	13317120
gradients/mean_squared_error_1/Mul_grad/mul	13317120
affs_0/BiasAdd	13317120
gradients/mean_squared_error_1/SquaredDifference_grad/sub mean_squared_error_1/Sum	13317120 13317119
gradients/affs_0/BiasAdd_grad/BiasAddGrad	13317119
unet_layer_2_right_0/kernel/Initializer/random_uniform	4860000
unet_layer_3_left_0/BiasAdd	7644000
gradients/unet_layer_3_left_0/BiasAdd_grad/BiasAddGrad	7642500
unet_layer_3_left_1/BiasAdd	5184000
gradients/unet_layer_3_left_1/BiasAdd_grad/BiasAddGrad	5182500
unet_layer_2_right_1/kernel/Initializer/random_uniform	2430000
unet_layer_2_left_1/kernel/Initializer/random_uniform unet_layer_2_left_0/kernel/Initializer/random_uniform	2430000 486000
unet_layer_1_right_0/kernel/Initializer/random_uniform	194400
unet_up_2_to_1/kernel/Initializer/random_uniform	162000
unet_layer_1_right_1/kernel/Initializer/random_uniform	97200
unet_layer_1_left_1/kernel/Initializer/random_uniform	97200
unet_layer_1_left_0/kernel/Initializer/random_uniform	19440
unet_layer_0_right_0/kernel/Initializer/random_uniform	9072
unet_up_1_to_0/kernel/Initializer/random_uniform	6480
unet_layer_0_right_1/kernel/Initializer/random_uniform	5292
unet_layer_0_left_1/kernel/Initializer/random_uniform	3888 324
unet_layer_0_left_0/kernel/Initializer/random_uniform embedding_0/kernel/Initializer/random_uniform	324 140
affs_0/kernel/Initializer/random_uniform	42
gradients/Slice_1_grad/sub	5
gradients/Slice_1_grad/sub_1	5
gradients/Slice_2_grad/sub	5
gradients/Slice_2_grad/sub_1	5
gradients/Slice_grad/sub	5
gradients/Slice_grad/sub_1	5
mean_squared_error_1/Greater	1
unet_up_1_to_0/mul unet_up_3_to_2/mul_2	1
unet_up_3_to_2/mul_1	1
unet_up_3_to_2/mul	1
unet_up_3_to_2/add_2	1
unet_up_3_to_2/add_1	1
unet_up_3_to_2/add	1
unet_up_2_to_1/mul_2	1
unet_up_2_to_1/mul_1	1
unet_up_2_to_1/mul unet_up_2_to_1/add_2	1
unet_up_2_to_1/add_2 unet_up_2_to_1/add_1	1
unet_up_2_to_1/add	1
unet_up_1_to_0/mul_2	1
unet_up_1_to_0/mul_1	1
gradients/mean_squared_error_1/div_grad/RealDiv_2	1
Adam/mul_1	1
add	1
gradients/mean_squared_error/div_grad/Neg	1 1
gradients/mean_squared_error/div_grad/RealDiv	1
gradients/mean_squared_error/div_grad/RealDiv_1 gradients/mean_squared_error/div_grad/RealDiv_2	1
gradients/mean_squared_error/div_grad/mul	1
gradients/mean_squared_error_1/div_grad/Neg	1
gradients/mean_squared_error_1/div_grad/RealDiv	1
gradients/mean_squared_error_1/div_grad/RealDiv_1	1
unet_up_1_to_0/add_2	1
gradients/mean_squared_error_1/div_grad/mul	1 1
mean_squared_error/Equal mean_squared_error/Greater	1
mean squared error/div	1
mean_squared_error_1/Equal	1
Adam/mul	1
mean_squared_error_1/div	1
unet_up_1_to_0/add	1
unet_up_1_to_0/add_1	1
Total	4083673765073

Table 5: FLOPs breakdown by operation for MTLsD on Zebrafinch dataset.

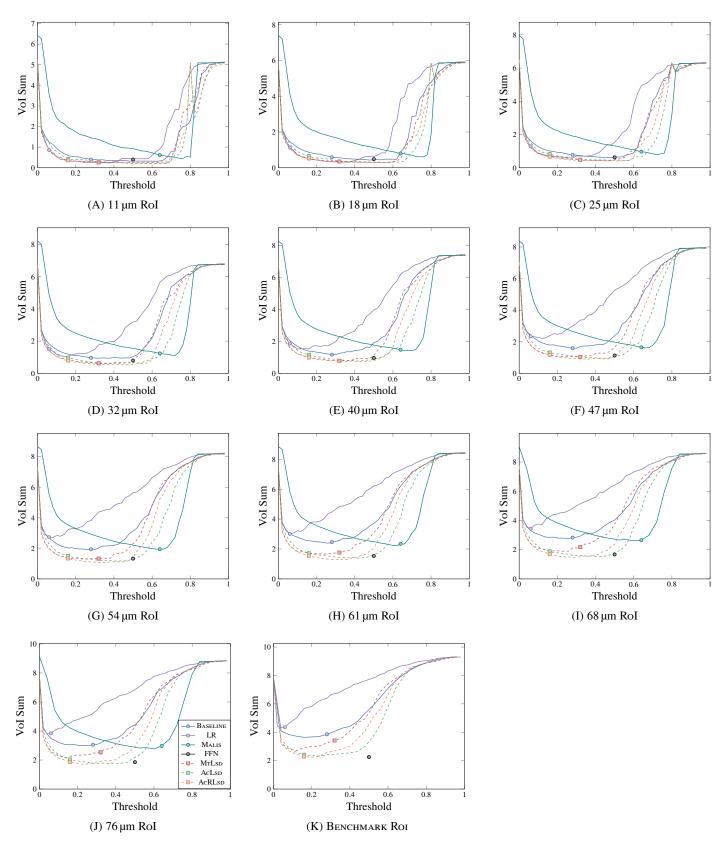


Figure 12: VoI Sum vs threshold across Zebrafinch RoIs. Points correspond to thresholds which minimized VoI Sum on the validation dataset.

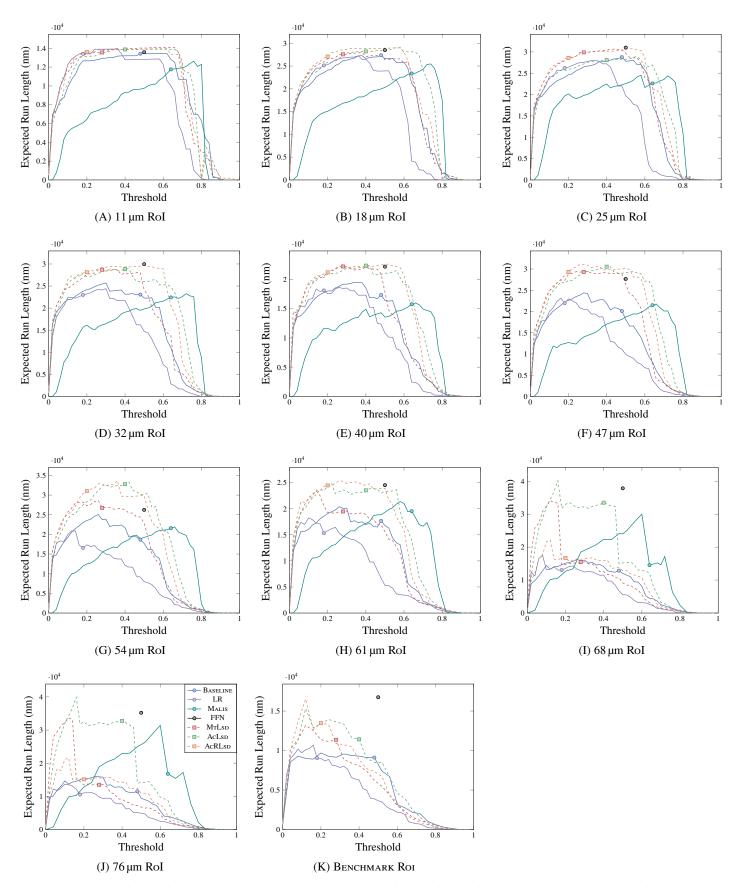


Figure 13: ERL vs threshold across Zebrafinch RoIs. Points correspond to thresholds which maximized ERL on the validation dataset.

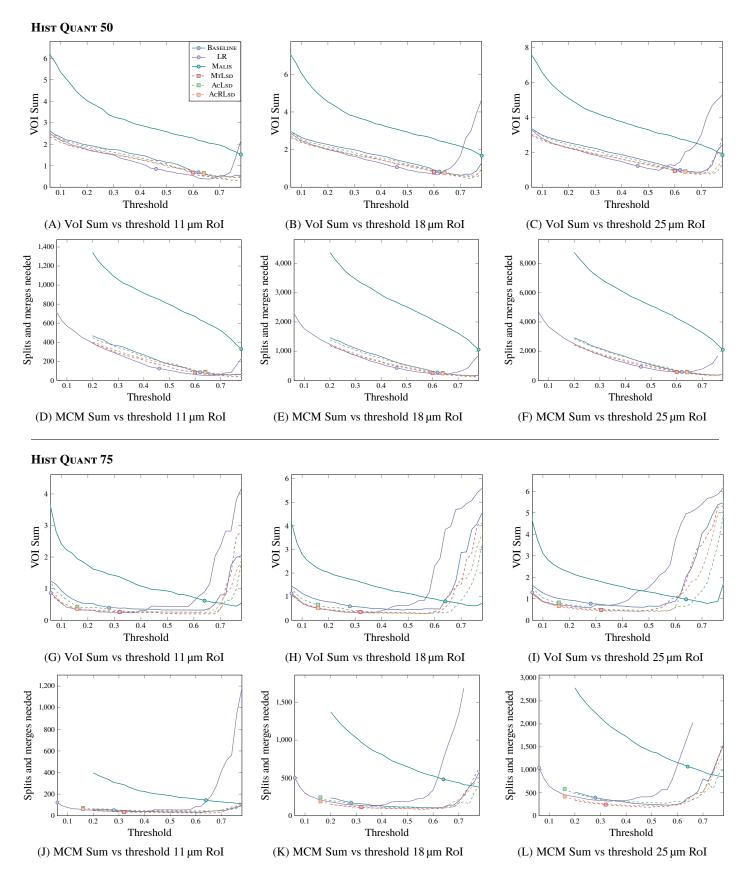


Figure 14: VoI serves as a reasonable proxy for evaluating large volumes. Comparison between VoI and MCM on Zebrafinch on first three RoIs. Similarities are consistent on both Hist Quant 50 (top two rows) and Hist Quant 75 merge functions (bottom two rows). VoI plots are cropped to the threshold range used in the MCM.

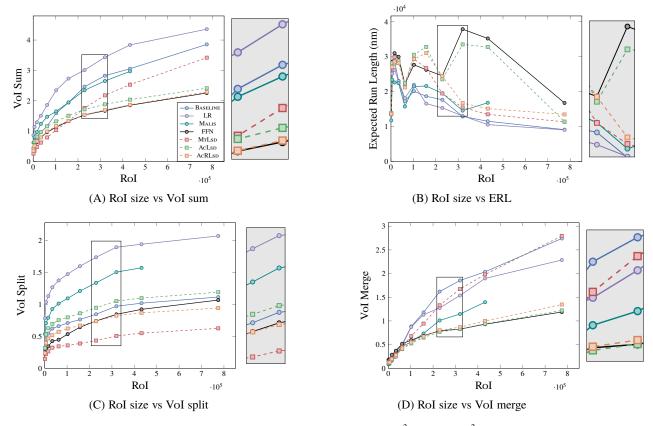


Figure 15: Demonstration of ERL sensitivity. When transitioning from a  $\sim$ 61 $\mu$ m<sup>3</sup> to a  $\sim$ 68 $\mu$ m<sup>3</sup> RoI, the VoI Sum increases as expected (**A**). This is not the case when considering ERL. All networks, with the exception of FFN and AcLsD, show decreases (**B**). Notably, AcRLsD shows a pattern not consistent with VoI Sum. Breaking VoI Sum into false splits (**C**.) and merges (**D**.) shows that there is a slight increase in the AcRLsD merge score between these two RoIs. Even if it is just a single added false merge, the resulting ERL is drastically decreased.

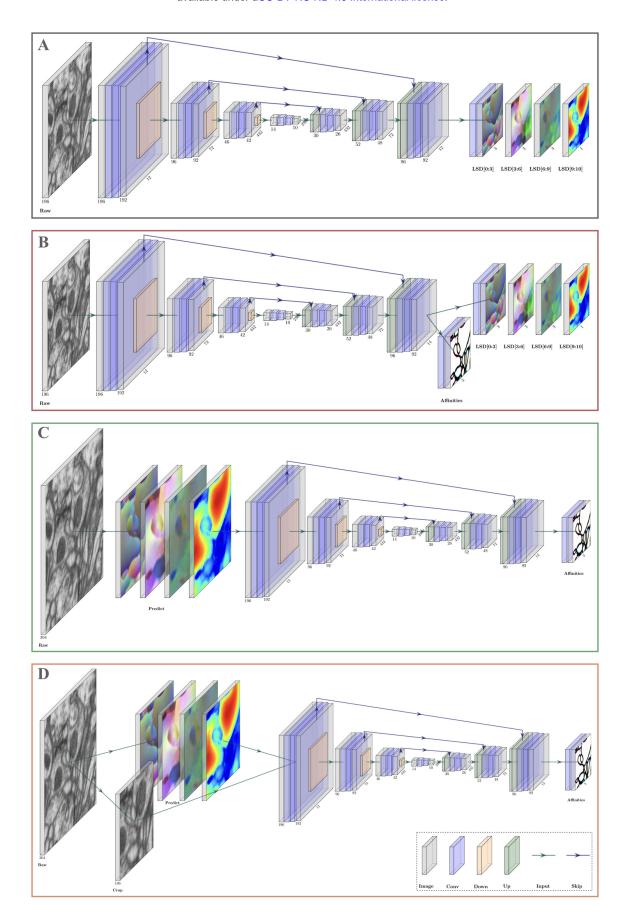


Figure 16: Network architectures used on FIB-SEM datasets (2D representations). All architectures use the U-Net proposed in Funke et al. (2019). Legend on bottom right. A. Network to generate the LSDs used in auto-context setups. An extra convolution is used to get to 10 feature maps for the embedding. B. MTLsD network - both affinities and LSDs are learnt. Number of output feature maps is increased from 12 to 14 to account for the 13 feature maps needed for the affinities and embedding. C. AcLsD network - the output from A is used to predict embedding from raw input. The predicted embedding is then passed in to learn affinities. D. AcRLsD network - same as C, but incorporates cropped raw as input in addition to the embedding.

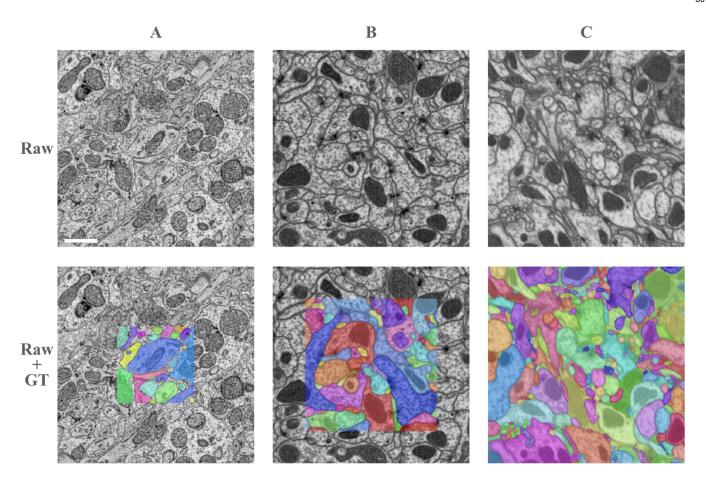


Figure 17: Example training data. A. Zebrafinch labels were heavily padded with raw data and some glia were set to zero (scale bar =  $\sim 1 \mu m$ ). B. Padding was used to a lesser extent for Hemi-brain volumes. Example taken from Ellipsoid Body. C. No padding was used for Fib-25 volumes.

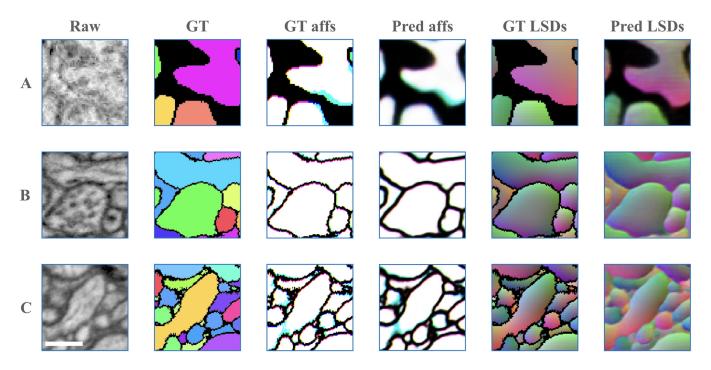


Figure 18: Example training batches. Masked-out regions were factored into the training loss on the Zebrafinch dataset (**A**). Conversely, the Hemi-brain and Fib-25 volumes used no masking during training (**B,C**, respectively). Scale bar =  $\sim$ 300 nm.

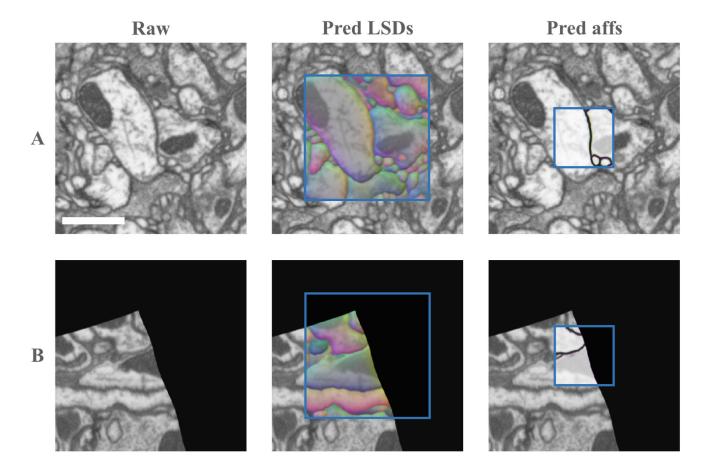


Figure 19: Example auto-context training batches on Fib-25. **A**. Batch in which no boundary masking is needed. The first pass predicts LSDs in an intermediate RoI to provide context for affinity prediction in the second pass. **B**. Batch requires boundary masking. The combination of elastic deformation and zero padding simulates the tissue irregularities seen in the full Fib-25 volume. In these background areas, LSDs and affinities are taught to predict zero. Scale bar =  $\sim 1 \, \mu m$ .

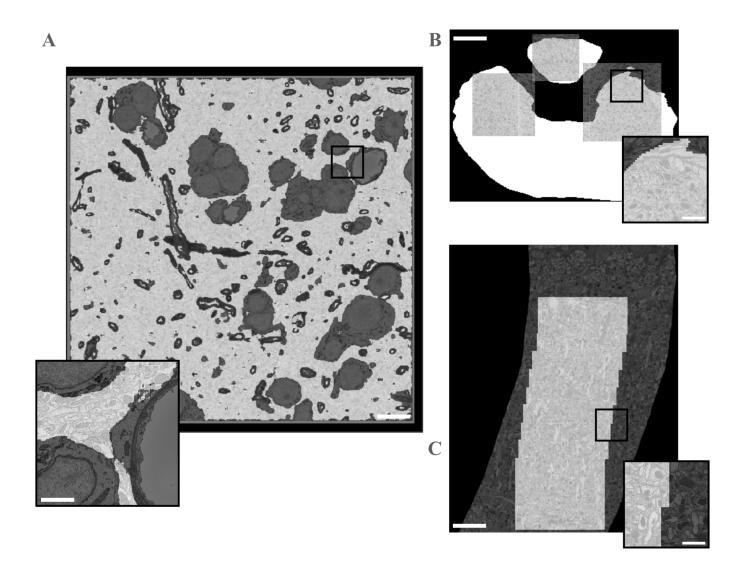


Figure 20: Masks used in study. A. Zebrafinch mask removed cell bodies, myelin, blood vessels and background. Scale bar =  $\sim$ 10  $\mu$ m. Inset scale bar =  $\sim$ 3  $\mu$ m. B. Hemi-brain mask restricted volumes to Ellipsoid Body neuropil. Scale bar =  $\sim$ 10  $\mu$ m. Inset scale bar =  $\sim$ 2  $\mu$ m. C. Fib-25 used an irregularly shaped tissue mask, mostly limited to neuropil. Scale bar =  $\sim$ 8  $\mu$ m. Inset scale bar =  $\sim$ 1  $\mu$ m.

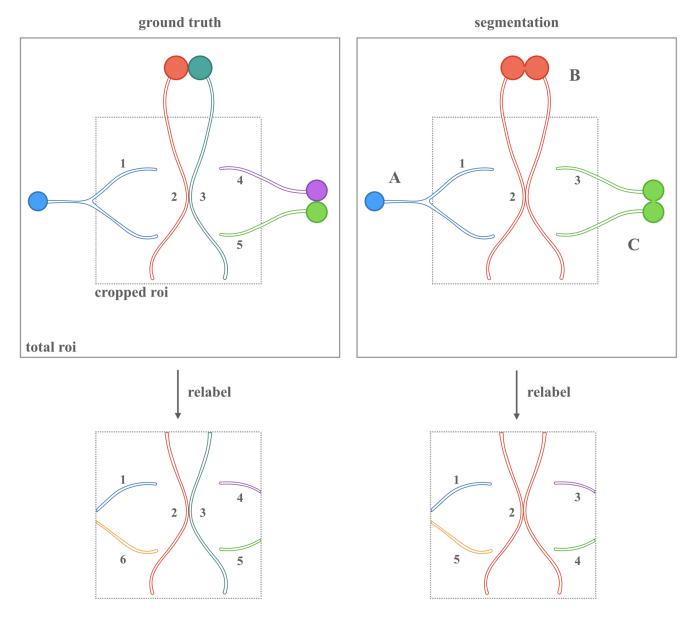


Figure 21: Potential side effects of cropping data. Ground truth (left) shows 5 correctly labeled neurons. Example Segmentation (right) shows 3 labeled neurons as a result of false merges. Bottom squares show results from relabeling connected components inside the cropped RoI. A. Correctly segmented neuron in total RoI would be counted as a false merge inside cropped RoI. This is fixed by relabelling connected components and the merge/split scores are unaffected. B. A falsely merged neuron inside the cropped RoI is caused by a false merge outside of the cropped RoI and should not be counted. Relabelling doesnt resolve the touching boundaries and the merge score is subsequently overestimated. C. Incorrectly segmented neuron in total RoI is counted as correct in cropped RoI after relabelling. Merge score is underestimated.

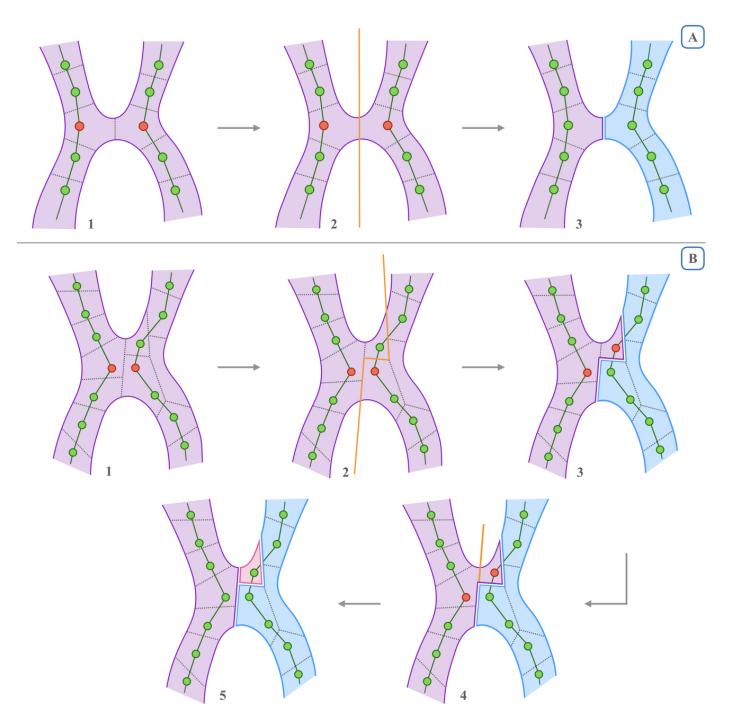


Figure 22: Overview of the proposed Min-cut Metric (MCM). A. Simple case. Two ground-truth skeletons are contained inside an erroneously merged segment. Dashed lines represent supervoxel boundaries and the closest skeleton nodes need to be split to resolve the merge (1). A min-cut is performed (2), resulting in a new segment (3). B. Complex case. Two skeletons are contained in a falsely merged segment as before (1), but the supervoxels are more fragmented. A min-cut is performed (2), resulting in a new segment (3). However, two nodes contained within the original segment need to be split. A second min-cut is performed (4), which produces another segment (5). This results in an additional split error caused by the original cut.

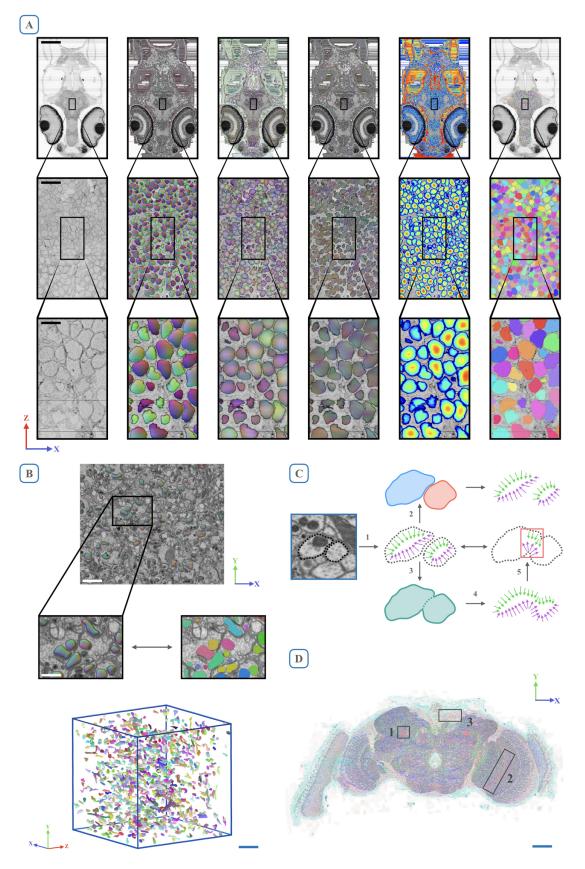


Figure 23: Other potential uses of LSDs. **A.** Nuclei segmentation on full zebrafish brain (Hildebrand et al., 2017). Columns from left to right: raw, LSD offset vectors, LSD direction vectors (covariance), LSD direction vectors (Pearson's), size, resulting segmentation. Scale bars from top to bottom:  $\sim 150\mu m$ ,  $20\mu m$ ,  $5\mu m$ . **B.** Mitochondria segmentation on cropout from Fib-25. Inset shows LSD predictions and corresponding segmentation. Bottom image shows 3D reconstructions of a random sample (n=1000) in predicted RoI. Scale bars from top to bottom:  $\sim 3\mu m$ , 750nm,  $4\mu m$ . **C.** Error mapping. Example predicted LSDs between two neurons (1). If the resulting segmentation is correct (2), segmentation LSDs do not differ from predicted LSDs. If the resulting segmentation is incorrect (3), segmentation LSDs (4) might differ from the predicted LSDs. The difference (5) could expose errors in a segmentation. **D.** Predicted direction vectors (covariance) on single section of full adult fly brain. Mushroom body pedunculi (1), optic chiasm (2), cell rind (3) highlight directionality. Scale bar =  $\sim 150\mu m$ .