

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# Guided Proofreading of Automatic Segmentations for Connectomics

Anonymous CVPR submission

Paper ID 0947

## Abstract

Automatic cell image segmentation methods in connectomics produce merge and split errors, which require correction through proofreading. Previous research has identified the visual search for these errors as the bottleneck in interactive proofreading. To aid error correction, we develop two classifiers to recommend candidate merge and split errors and their corrections to the user. These classifiers are informed by training a convolutional neural network with known errors in automatic segmentations against expert-labeled ground truth. Our classifiers detect potentially-erroneous regions by considering a large context region around a segmentation boundary. Corrections can then be performed as yes/no decisions resulting in faster correction times than previous methods. We evaluate our approach on connectomics datasets of different species and compare correction performance of novice and expert users against different existing systems. We report significant improvements compared to pure automatic and pure manual proofreading.

## 1. Introduction

In connectomics, neuroscientists annotate neurons and their connectivity within 3D volumes to gain insight into the functional structure of the brain. Rapid progress in automatic sample preparation and electron microscopy (EM) acquisition techniques has made it possible to image large volumes of brain tissue at nanometer resolution. A typical voxel size is  $4 \times 4 \times 40 \text{ nm}^3$ . For a  $1 \text{ mm}^3$  volume, this resolution results in 1 petabyte of data. With so much data, manual annotation is not feasible, and automatic annotation methods are needed [12, 22, 25, 17].

Automatic annotation by segmentation and classification of brain tissue is challenging [1] and all available methods make errors. This means we are left with large data which needs *proofreading* by humans. This crucial task serves two purposes: 1) to correct errors in the segmentation, and 2) to provide a large body of labeled data to train better automatic segmentation methods. Recent proofread-

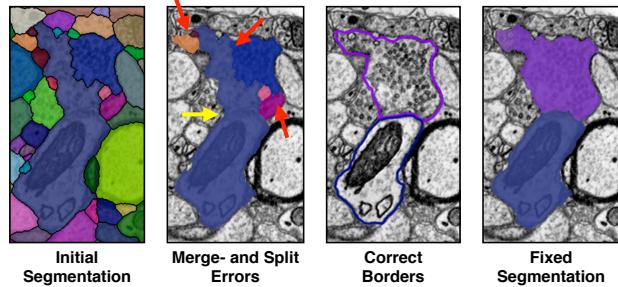


Figure 1. The most common proofreading corrections are fixing split errors (red arrows) and merge errors (yellow arrow). A fixed segmentation matches the cell borders.

ing tools provide intuitive user interfaces to browse segmentation data in 2D and 3D and to identify and manually correct errors [30, 13, 19, 10]. Many kinds of errors exist, such as inaccurate boundaries, but the most common are *split errors*, where a single segment is labeled as two, and *merge errors*, where two segments are labeled as one (Fig. 1). With user interaction, split errors can be joined, and the missing boundary in a merge error can be defined with manually-seeded watersheds [10]. However, even with semi-automatic correction tools, the visual inspection to find errors takes the majority of the time [26].

Our goal is to add automatic detection of split and merge errors to proofreading tools. We design automatic classifiers that detect split and merge errors in segmentations so the user does not need to visually inspect the whole data volume to spot errors. A proofreading tool then recommends regions with a high probability of an error to the user, and suggest corrections to accept or reject. We call this process *guided proofreading*.

In this paper, we introduce classifiers to detect merge-and split errors based on a convolutional neural network (CNN). We believe that this is the first time that deep learning is applied to the task of proofreading. Our classifiers work on top of any existing automatic segmentation method to find potential errors and suggest corrections. Given a membrane segmentation from a fast automatic method, our classifiers operate on the boundaries of whole cell regions.

108 Compared to techniques that must analyze every input pixel,  
109 we reduce the data analysis to the boundaries. First, we  
110 train a CNN to detect only split errors. The output of this  
111 network is a probability whether a boundary between two  
112 segments is valid or not. We then reuse the same network to  
113 also detect merge errors by generating possible boundaries  
114 within a cell and inverting the split error score. We create  
115 corrections for both types of errors which can be accepted  
116 or rejected. This reduces the proofreading operation to sim-  
117 ple yes/no decisions.  
118

119 We further propose a greedy algorithm to perform proof-  
120 reading. Possible erroneous regions are sorted by their score  
121 and the algorithm iteratively suggests a correction for each  
122 region. A user then works through this stream of regions  
123 and corrections. In a forced choice setting, the user either  
124 selects a correction or skips it to advance to the next region.  
125 This choice can be also performed automatically by run-  
126 ning the algorithm until a configurable threshold is reached.  
127 In addition, if ground truth data is available, we can use a  
128 selection oracle to drive the forced choice selection. The or-  
129 acle only accepts corrections which improve the automatic  
130 segmentation. This equals perfect proofreading.

131 We evaluate our method automatically by threshold and  
132 oracle on multiple real-world connectomics datasets. To  
133 evaluate the forced choice setting, we perform a quantita-  
134 tive user study. The study targets non-experts with no pre-  
135 vious experience of proofreading electron microscopy data.  
136 We ask the participants to proofread a small segmentation  
137 volume in a fixed time frame by performing yes/no deci-  
138 sions. The user study is designed as a between-subjects  
139 experiment and compares guided proofreading against two  
140 other methods: a recently published fully interactive proof-  
141 reading tool named *Dojo* by Haehn *et al.* [10] and the semi-  
142 automatic *focused proofreading* approach by Plaza [27]. We  
143 also asked four domain experts to use guided proofreading  
144 and focused proofreading for additional comparison.

145 Our first contribution is a classifier for split error detec-  
146 tion based on a convolutional neural network. The clas-  
147 sifier performs well even when trained with little amounts  
148 of training data. This is important since generating ground  
149 truth labels in connectomics requires manually labeling pix-  
150 els and is very time-consuming. Our second contribution is  
151 a mechanism to identify merge-errors by re-using the split  
152 error classifier. Merge errors are usually less common than  
153 split errors in the oversegmented automatic labelings. How-  
154 ever, they require more interaction during correction since  
155 split lines need to be manually drawn. Our method reduces  
156 this to a single click by providing the potential correction.  
157 The split and merge error identification is executed as a  
158 greedy algorithm to correct segmentation volumes, the third  
159 contribution of this paper. The algorithm can be driven au-  
160 tomatically with a threshold, by an oracle based on ground  
161 truth and interactively in a forced choice setting. Our final

162 contribution is our quantitative user study. We present sta-  
163 tistically significant results showing that novice and expert  
164 users of guided proofreading are able to proofread a given  
165 dataset better and faster than with existing interactive and  
166 semi-automatic proofreading tools. As a consequence, we  
167 are able to provide tools to proofread segmentations more  
168 efficiently, and so better tackle large volumes of connec-  
169 toomics imagery.  
170

## 2. Related Work

171 **Automatic Segmentation.** Multi-terabyte EM brain vol-  
172 umes require automatic segmentation [12, 22, 24, 25], but  
173 can be hard to classify due to ambiguous intercellular space:  
174 the 2013 IEEE ISBI neurites 3D segmentation challenge [1]  
175 showed that existing algorithms which learn from expert-  
176 segmented training data still exhibit high error rates.  
177

178 Many works tackle this problem. NeuroProof [2] de-  
179 creases error rates by learning a agglomeration on over-  
180 segmentations of images, based on a random forest classi-  
181 fier. Vazquez-Reina *et al.* [33] take whole EM volumes into  
182 account rather than a per section approach, then solve a fu-  
183 sion problem with a global context. Kaynig *et al.* [16] pro-  
184 pose a random forest classifier coupled with an anisotropic  
185 smoothing prior in a conditional random field framework  
186 with 3D segment fusion. Bogovic *et al.* [6] learn 3D fea-  
187 tures unsupervised, and show that they can be better than  
188 by-hand designs.  
189

190 It is also possible to learn segmentation classification  
191 features directly from images with CNNs. Ronneberger *et*  
192 *al.* [28] use a contracting/expanding CNN path architecture  
193 to enable precise boundary localization with small amounts  
194 of training data. Lee *et al.* [20] recursively train very deep  
195 networks with 2D and 3D filters to detect boundaries.  
196

197 All these approaches make good progress; however, in  
198 general, proofreading is still required to correct errors.  
199

200 **Interactive Proofreading.** While proofreading is very time  
201 consuming, it is fairly easy for humans to perform correc-  
202 tions through splitting and merging segments. One expert  
203 tool is Raveler, introduced by Chklovskii *et al.* [7, 13].  
204 Raveler is used today by professional proofreaders, and it  
205 offers many parameters for tweaking the process. Similar  
206 systems exist as products or plugins to visualization  
207 systems, *e.g.* V3D [26] or AVIZO [30].  
208

209 Recent papers have attacked the problem of proofreading  
210 massive datasets through crowdsourcing with novices [29,  
211 4, 9]. One popular platform is EyeWire, by Kim *et al.* [18],  
212 where participants earn virtual rewards for merging over-  
213 segmentated labeling to reconstruct retina cells.  
214

215 Between expert systems and online games sit Mojo and  
216 Dojo, by Haehn *et al.* [10, 3], which use simple scribble  
217 interfaces for error correction. Dojo extends this to distributed  
218 proofreading via a minimalistic web-based user interface.  
219

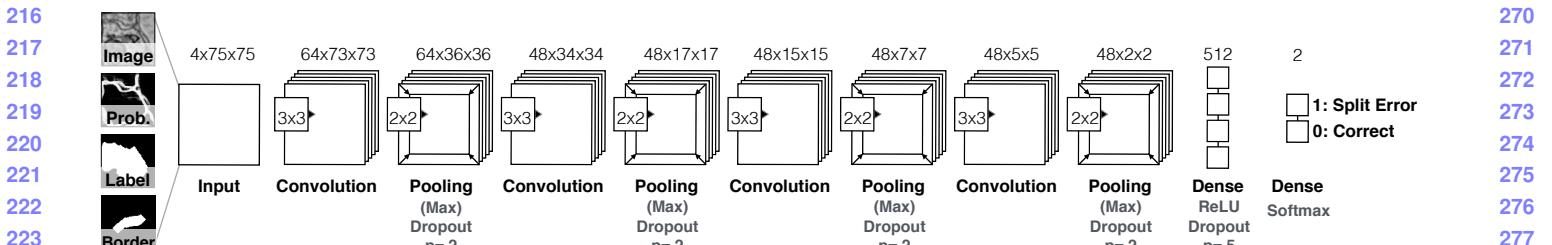


Figure 2. We build the guided proofreading classifiers using a traditional CNN architecture. The network is based on four convolutional layers, each followed by max pooling as well as dropout regularization. The 4-channel input patches are rated as either correct splits or as split errors.

The authors define requirements for general proofreading tools, and then evaluate the accuracy and speed of Raveler, Mojo, and Dojo through a quantitative user study (Sec. 3 and 4) [10]. Dojo had the highest performance. In this paper, we use Dojo as a baseline for interactive proofreading, and so we extend the Haehn *et al.* experiment.

All interactive proofreading solutions require the user to find potential errors manually, which takes the majority of the time [26, 10]. Recent works propose computer-aided proofreading systems which help reduce the time spent in this visual search task.

**Computer-aided Proofreading.** Uzunbas *et al.* showed that potential labeling errors can be found by considering the merge tree of an automatic segmentation method [32]. The authors track uncertainty throughout the automatic labeling by training a conditional random field. This segmentation technique produces uncertainty estimates, which can be used to present potential regions for proofreading to the user. While it works on isotropic volumes, more work is needed to apply it to anisotropic volumes, like most connectomics datasets.

Karimov *et al.* propose guided volume editing [14], which measures the difference in histogram distributions in image data to find potential split and merge errors in the corresponding segmentation. This lets expert users correct labeled computer-tomography datasets, using several interactions per correction. To correct merge errors, the authors create a large number of superpixels within a single segment and then successively group them based on dissimilarities. We were inspired by this approach but generate single watershed boundaries to handle the intracellular variance in high-resolution EM images (Sec. 3).

Most closely related to our approach is the work of Plaza, who proposed *focused proofreading* [27]. This method generates affinity scores by analyzing a region adjacency graph across slices, then finds the largest affinities based on a defined impact score. This yields edges of potential split errors which can be presented to the proofreader. Plaza reports that additional manual work is required to find and correct merge errors. Focused proofreading builds upon

NeuroProof [2] as its agglomerator, and is open source with integration into Raveler. As the closest related work, we wish to use this method as a baseline to evaluate our approach (Sec. 4). However, as Haehn *et al.* showed that Raveler is less performant than Dojo for novice users, we separate the backend affinity score calculation from the expert-level front end, and present our own interface (Sec. 4).

### 3. Method

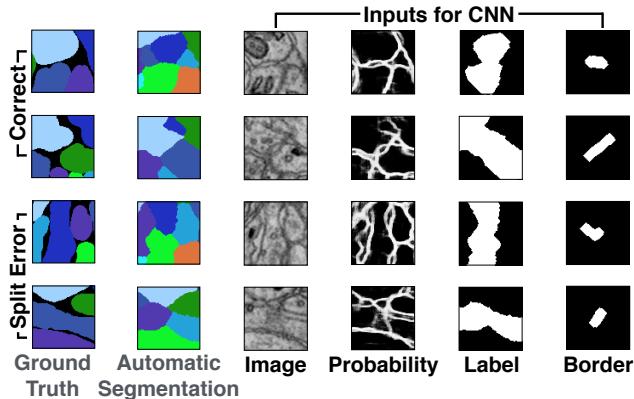
#### 3.1. Split Error Detection

We build a split error classifier with output  $p$  using a convolutional neural network (CNN) to check whether an edge within an existing automatic segmentation is valid ( $p = 0$ ) or not ( $p = 1$ ). Rather than analyzing every input pixel, the classifier operates only on segment boundaries which requires less pixel context and is faster. In contrast to Bogovic *et al.* [6], we work with 2D slices rather than 3D volumes. This enables proofreading prior or in parallel to an expensive alignment of individual EM images.

**Convolutional Neural Network Architecture.** Split error detection of a given boundary is really a binary classification task since the boundary is either correct or erroneous. However, in reality the score  $p$  is between 0 and 1. The classification complexity arises from hundreds of different cell types in connectomics data rather than from the classification decision. Intuitively, this yields a wider (meaning more filters) rather than a deeper (meaning more layers) architecture. We explored different architectural configurations - including residual networks [11] - by performing a brute force parameter search and comparing precision and recall (see supplementary materials). Our final CNN configuration for split error detection is composed of four convolutional layers, each followed by max pooling as well as dropout regularization to prevent overfitting due to limited training data. Fig. 2 shows the CNN architecture for split error detection.

**Classifier Inputs.** To train the CNN for split error detection, we take boundary context information into

324 consideration for the decision making process. For this,  
 325 we use a  $75 \times 75$  pixel patch at the center of an existing  
 326 boundary. This covers approximately 80% of all bound-  
 327 aries in real-world connectomics data with nanometer  
 328 resolution. If the boundary edge is not fully covered,  
 329 we sample up to 10 non-overlapping patches along the  
 330 boundary and combine the resulting score by weighted  
 331 averaging based on boundary length coverage per patch.  
 332 Similar to Bogovich *et al.* [6], we use grayscale image data,  
 333 corresponding boundary probabilities, and a single binary  
 334 mask combining the two neighboring labels as features  
 335 for our CNN. However, we observed that the boundary  
 336 probability information generated from EM images is often  
 337 misleading due to noise or artifacts in the data. This can  
 338 result in merge errors within the automatic segmenta-  
 339 tion. To better direct our classifier to train on the true boundary  
 340 edge, we extract the border between two segments. We  
 341 then dilate this border by 5 pixels to consider slight edge  
 342 ambiguities and use this binary mask as an additional fea-  
 343 ture to create a stacked 4-channel input patch. Fig. 3 shows  
 344 examples of correct and erroneous feature patches and their  
 345 corresponding automatic segmentation and ground truth.



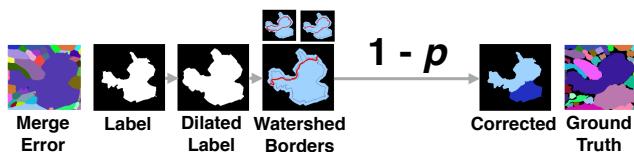
360 Figure 3. Example inputs for learning correct splits and split errors  
 361 as reflected in the segmentation relative to the ground truth. Im-  
 362 age, membrane probabilities, merged binary labels, and a dilated  
 363 border mask are combined to 4-channel input patches.

### 3.2. Merge Error Detection

367 Identification and correction of merge errors is more  
 368 challenging than finding and fixing split errors, because we  
 369 must look inside segmentation regions for missing or in-  
 370 complete boundaries and then propose the correct boundary.  
 371 However, we can reuse the same trained CNN for this task.  
 372 Similar to guided volume editing by Karimov *et al.* [14] we  
 373 generate potential borders within a segment. For each seg-  
 374 mentation label, we dilate the label by 20 pixel and gener-  
 375 ate 50 potential boundaries through the region by randomly  
 376 placing watershed seed points at opposite sides of the label  
 377 boundary. For watershed, we use the inverted gray scale

378 EM image as features. This yields 50 candidate splits.  
 379

380 Dilation of the segment prior to watershed is motivated  
 381 by our observation that the generated split likely attaches  
 382 to the real membrane boundary. These boundaries are then  
 383 individually rated using our split error classifier. For this,  
 384 we invert the probability score such that a correct split (pre-  
 385 viously encoded as  $p = 0$ ) is most likely a candidate for a  
 386 merge error (now encoded as  $p = 1$ ). In other words, if a  
 387 generated boundary is ranked as correct, it probably should  
 388 be in the segmentation. Fig. 4 illustrates this procedure.



390 Figure 4. Merge errors are identified by generating randomly  
 391 seeded watershed borders within a dilated label segment. These  
 392 borders then are individually rated using the split error CNN by  
 393 inverting the probability score. This way, a confident rating for a  
 394 correct split most likely indicates the missing border of the merge  
 395 error and can be used for correcting the labeling.

### 3.3. Error Correction

401 We use the proposed classifiers in combination to  
 402 perform corrections of split and merge errors in automatic  
 403 segmentations. For this, we first perform merge error  
 404 detection for all existing segments in a dataset and store  
 405 the inverted rankings  $1 - p$  as well as potential corrections.  
 406 After that, we perform split error detection and store the  
 407 ranking  $p$  for all neighboring segments in the segmentation.  
 408 We then sort the merge and split error rankings individually  
 409 from highest to lowest. For error correction, we first loop  
 410 through the potential merge error regions and then through  
 411 the potential split error regions. During this process, each  
 412 error is now subject to a yes/no decision which can be  
 413 provided in different ways:

414 **Selection oracle.** If ground data is available, the selection  
 415 oracle *knows* whether a possible correction improves  
 416 an automatic segmentation. This is realized by simply  
 417 applying the correction and comparing the outcome using  
 418 a defined measure. The oracle only accepts corrections  
 419 which improve the automatic segmentation - others get  
 420 discarded. This equals perfect proofreading.

421 **Automatic selection with threshold.** The decision  
 422 whether to accept or reject a potential correction is done  
 423 by comparing rankings to a threshold  $p_t$ . If the inverted  
 424 score  $1 - p$  of a merge error is higher than a threshold  
 425  $1 - p_t$ , the correction is accepted. Similarly, a correction  
 426 is accepted for a split error if the ranking  $p$  is higher than  
 427  $p_t$ . Our experiments have shown that the threshold  $p_t$  is the

432 same for merge and split errors which makes sense for a  
 433 balanced classifier.  
 434

435 **Forced choice setting.** A user is presented with the  
 436 choices of either accepting or rejecting a correction. This  
 437 way, all potential split errors are looked at. Inspecting  
 438 all merge errors is not possible for users due to the sheer  
 439 amount of generated borders. We therefor only present  
 440 merge errors which satisfy  $1 - p_t$ .  
 441

442 In all cases, a decision has to be made to advance to  
 443 the next possible erroneous region. If a merge error  
 444 correction was accepted, the newly found boundary is  
 445 added to the segmentation data. This partially updates  
 446 the merge error and split error ranking in respect of the  
 447 new segment. If a split error correction was accepted, two  
 448 segments are merged in the segmentation data and the  
 449 disappearing segment is removed from all error rankings.  
 450 We then perform merge error detection on the now larger  
 451 segment and update the ranking. We also update the split  
 452 error rankings to include all new neighbors and re-sort. The  
 453 error with the next highest ranking is now subject to the  
 454 yes/no decision.  
 455

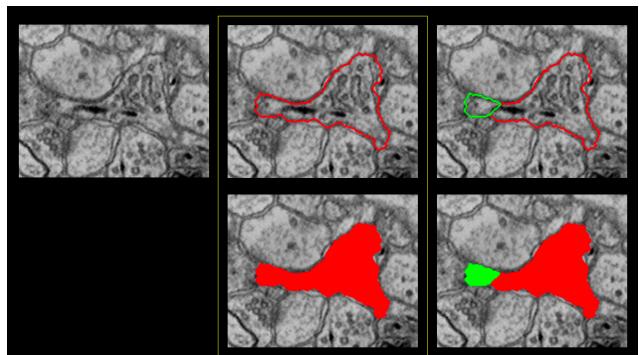
### 456 3.4. User Interface

457 Guided proofreading is integrated into an existing workflow  
 458 for large connectomics data. The system is web-based  
 459 and is designed with a minimalistic user interface showing  
 460 three components. We show the outline of the current  
 461 labeling of a cell boundary and its proposed correction on top  
 462 of the EM image data. For the user, it is not possible to  
 463 distinguish the current labeling and the proposed correction  
 464 to avoid selection bias. We also show a solid overlay of  
 465 the current and the proposed labeling. In addition, we show  
 466 the image without overlays to provide an unoccluded view.  
 467 User interaction is simple and involves one mouse click on  
 468 either the current labeling or the correction. After interaction,  
 469 the next potential error is shown. Figure 5 shows the  
 470 user interface.  
 471

## 472 4. Evaluation

473 We evaluate guided proofreading (GP) on multiple real-  
 474 world connectomics datasets of different species. In partic-  
 475 ular, we evaluate GP on two datasets of mouse brain and  
 476 three datasets of fruitfly brain (*drosophila*). For compari-  
 477 son, we choose the fully interactive proofreading software  
 478 *Dojo* by Haehn *et al.* [10] as well as the aided proofreading  
 479 framework *focused proofreading* (FP) by Plaza [27]. We  
 480 first describe the evaluation on mouse brain data and then the  
 481 evaluation on *drosophila* brain.  
 482

483 **Training.** To initially train our network, we use the  
 484 blue 3-cylinder mouse cortex volume of Kasthuri *et al.* [15]  
 485 (2048 × 2048 × 300 voxels). The tissue is dense mam-



486  
 487  
 488  
 489  
 490  
 491  
 492  
 493  
 494  
 495  
 496  
 497  
 498  
 499  
 500  
 501  
 502  
 503  
 504  
 505  
 506  
 507  
 508  
 509  
 510  
 511  
 512  
 513  
 514  
 515  
 516  
 517  
 518  
 519  
 520  
 521  
 522  
 523  
 524  
 525  
 526  
 527  
 528  
 529  
 530  
 531  
 532  
 533  
 534  
 535  
 536  
 537  
 538  
 539

Figure 5. Segmentation correction with the guided proofreading user interface. An image without overlays is shown on the left. A split error (right) and its correction (center) are the possible choices for the user. Hovering highlights the current selection with a yellow border and a mouse click confirms it and advances to the next potential error.

malian neuropil from layers 4 and 5 of the S1 primary somatosensory cortex of a healthy mouse. The resolution of our dataset is 3 nm per pixel, and the section thickness is 30 nm. The image data and a manually-labeled expert segmentation is publicly available as ground truth for the entire dataset<sup>1</sup>. We use the first 250 sections of the data for training and validation and the last 50 for testing. We use a state-of-the-art method to create a dense automatic segmentation of the data. To generate training data, we identify correct regions and split errors in the automatic segmentation by intersection with ground truth regions. This is required since extracellular space is not labeled in the ground truth but in our dense automatic segmentation. From these regions, we sample 120,000 correct and 120,000 split error patches with 4-channels as described above. The patches are normalized and to further augment our training data, we rotate patches within each mini-batch by  $k * 90$  degrees with randomly chosen  $k$ . The training parameters such as filter size, number of filters, learning rate, and momentum are the result of intuition and experience, studying recent machine learning research as well as a brute force parameter search within a limited range (see supplementary material). The final parameters and training results are listed in table 1. For baseline comparison, we also list the parameters and training results of focused proofreading in this table but elaborate on these further in section 4. Our CNN configuration results in approximately 170,000 learnable parameters. We assume that training has converged if the validation loss does not decrease for 50 epochs.

For performance comparison on data of a different species, in particular on fruitfly brain (*drosophila*), we re-train our network. The training procedure is according

<sup>1</sup>The Kasthuri 3-cylinder mouse cortex volume is available at <https://software.rc.fas.harvard.edu/lichtman/vast/>

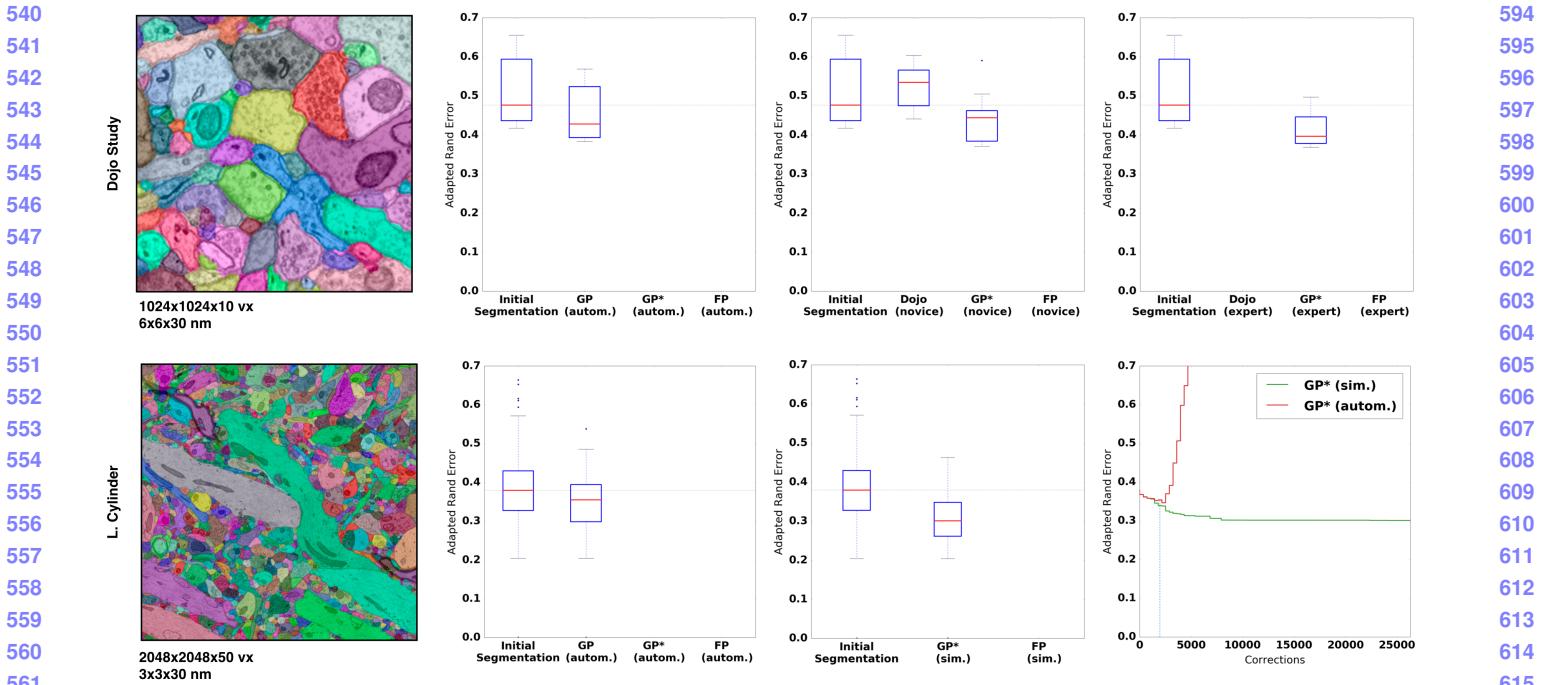


Figure 6. Performance evaluation of the classifiers on two mouse brain datasets measured as adapted Rand error (lower scores are better). We compare guided proofreading (GP), guided proofreading with active label suggestion (GP\*) and focused proofreading. Proofreading is performed automatically (autom., with probability threshold  $p_t = .95$ ), simulated as a perfect user (sim.), or by novice and expert users as indicated. The first row of images shows the results of a user study and includes comparisons to the interactive proofreading software Dojo by Haehn *et al.* [10]. GP\* is able to correct the segmentation further than other methods. The second row shows the results of the simulated user compared to automatic GP\* and FP performance. The bottom right graph compares automatic GP\* and simulated GP\* per individual correction. The blue dashed line here indicates the moment the probability threshold  $p_t$  is reached. The simulated user is able to correct the initial segmentation beyond this threshold while automatic GP\* then introduces errors.

	cost [m]	Val. loss	Val. acc.	Test acc.	Prec./Recall	F1 Score
<b>Guided Proofreading</b> Filter size: 3x3 No. Filters 1: 64 No. Filters 2-4: 48 Dense units: 512 Learning rate: 0.03-0.00001 Momentum: 0.9-0.999 Mini-Batchsize: 128	383	0.0845	0.969	0.94	0.94/0.94	0.94
<b>Focused Proofreading</b> Iterations: 3 Learning strategy: 2 Mito agglomeration: Off Threshold: 0.2	43	?	?	0.839	??	?

Table 1. Training parameters, cost and results of our guided proofreading classifier versus focused proofreading by Plaza [27]. Both methods were trained on the same mouse brain dataset using the same hardware (Tesla K40 graphics card). While the training of our classifier is more expensive, testing accuracy is superior.

to our initial training and network architecture as well as parameters are not changed. We further elaborate on the drosophila datasets in section 4. Fig. 7 displays receiver operating characteristics (ROC) for guided proofreading trained on mouse and drosophila data, as well as our comparison baseline focused proofreading trained on these datasets respectively.

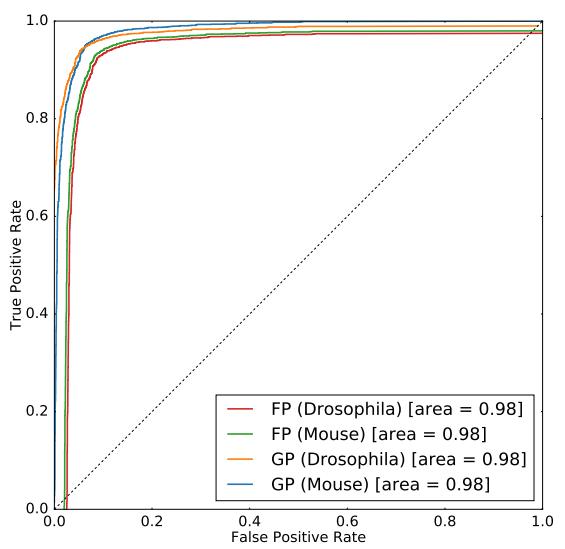


Figure 7. ROC performance of guided proofreading (GP) and focused proofreading (FP) trained separately on mouse and drosophila brain images. The area under the curve indicates better performance for GP.

648

## 4.1. Mouse Brain

649

Mouse brain is a common target for connectomics research because the structural proportions are similar to human brains [21]. For our first experiment we recruited novice and expert participants as part of a quantitative user study. Our second experiment is performed on a larger dataset and we evaluate a simulated user.

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

**User study.** Recently, Haehn *et al.* evaluated the interactive proofreading tools Raveler, Mojo, and Dojo as part of an experiment with novice users [10]. The participants corrected an automatic segmentation with merge and split errors. The dataset was the most representative sub-volume (based on object size histograms) of a larger connectomics dataset and  $400 \times 400 \times 10$  voxels in size. The participants were given a fixed time frame of 30 minutes to perform the correction interactively. While participants clearly struggled with the proofreading task, the best performing tool in their evaluation was Dojo. The dataset including manually labeled ground truth and the results of Haehn *et al.* are publicly available. This means we are able to use their findings as a baseline for comparison of GP for novices. In particular, we use the best performing user of Dojo who was truly an outlier as reported by Haehn *et al.*

Since interactive proofreading most likely yields lower performance than aided proofreading, we also compare against FP by Plaza [27] which is integrated in Raveler and freely available. For FP we consulted an expert to obtain the best possible parameters as shown in table 1. Besides performance by novices, we are also interested in expert proofreading performance. Therefore, we design between-subjects experiments for 20 novice users and separately, for 6 expert users using the exact same conditions as Haehn *et al.* The recruiting, consent and debriefing process is further described in the supplementary material. We randomly assign 10 novices to GP with active label suggestion (GP\*) and 10 novices to FP. For the expert experiment, we assign accordingly. In addition to human performance, we also evaluate automatic GP, automatic GP with active label suggestion (GP\*) and automatic FP. Due to the automatic nature, we do not enforce the 30 minute time limit but we stop once our probability threshold of  $p_t = .95$  is reached. This value was observed as stable in previous experiments using automatic GP (see supplementary material). To measure proofreading performance in comparison to ground truth, we use the adapted Rand error (aRE) metric [31]. aRE is a measure of dissimilarity, related to introduced errors, meaning lower scores are better.

The results of our comparisons are shown in the first row of Fig. 6. In all cases, GP\* is able to correct the segmentation further than other methods (aRE measures: automatic GP XX, GP\* XX, FP XX, novice Dojo XX, GP\* XX, FP XX, expert Dojo XX, GP\* XX, FP XX). This is

not surprising since guided proofreading works for both merge and split errors while FP does not and in interactive Dojo the majority of time is spent finding errors which is minimized for aided proofreading solutions. In fact, the average correction time for novices is for GP\* 3.6 (expert X), for FP Y (expert YY), and for Dojo 30 (expert ZZ) seconds.

**Simulated experiment.** For our second experiment with mouse brain data, we proofread the last 50 slices of the blue 3-cylinder mouse cortex volume of Kasthuri *et al.* [15] which we also used for testing in section 3. The data was not seen by the network before and includes  $2048 \times 2048 \times 50$  voxels with a total number of 17,560 labeled objects. Since an interactive evaluation of such a large dataset would consume a significant amount of time, we restrict our experiment to a simulated (perfect) user and to automatic corrections, both with GP, GP\* and FP. Similar to our comparison study, the simulated user assess a stream of errors by comparing the adapted Rand error measure before and after each performed correction. The simulated user is designed to be perfect and only accepts corrections if the measure is reduced. This time, we do not enforce a time limit to see the lower bound of possible corrections. For automatic GP and GP\*, we use our defined probability threshold  $p_t = .95$ .

The results of this experiment are shown in the second row of Fig. 6. GP\* is again able to correct the segmentation further than other methods (aRE measures: automatic GP XX, GP\* XX, FP XX, simulated GP\* XX, FP XX). Again, the results are not surprising since GP\* can correct merge and split errors.

## 4.2. Drosophila Brain

The drosophila brain is analyzed by connectomics researchers because of its small size and hence, a reasonable target to obtain a complete wiring diagram. Despite the size, fruit flies exhibit complex behaviors and are in general well studied. We evaluate the performance of our guided proofreading classifiers on three different datasets of adult fly brain. The datasets are publicly available as part of the MICCAI 2016 challenge on circuit reconstruction from electron microscopy images (CREMI)<sup>2</sup>. Each dataset consists of  $1250 \times 1250 \times 125$  voxels of training data (A,B,C) as well as testing data (A+,B+,C+) of the same dimensions. Manually labeled ground truth is also available for A,B, and C but not for the testing data.

Since drosophila brain exhibits different cell structures than mouse brain, we retrain the guided proofreading classifiers (and our automatic segmentation pipeline) as well as focused proofreading combined on the three training

<sup>2</sup>The MICCAI CREMI challenge data is available at <http://www.creml.org>

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

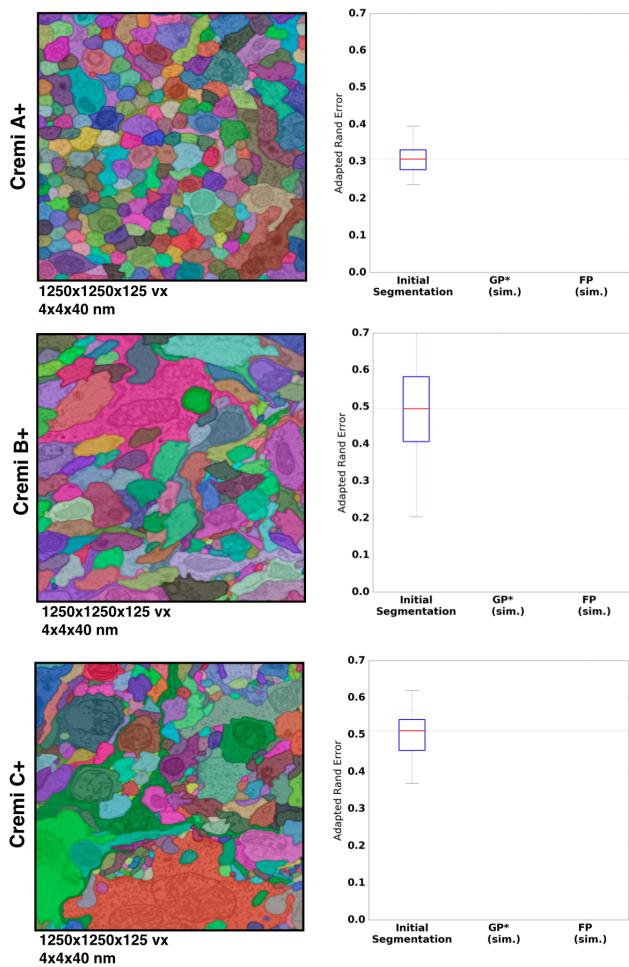
752

753

754

755

756 datasets. We use 300 slices of the A,B,C samples for training  
 757 and validation, and 75 slices for testing. This results  
 758 in YYY correct and ZZZ split error patches (respectively,  
 759 XXX and YYY for testing). The architecture and all pa-  
 760 rameters of our classifiers stay the same. The trained GP  
 761 classifier exhibits a reasonable performance on the testing  
 762 data as seen in Fig. 7.  
 763



795 Figure 8. Results of guided proofreading with active label sug-  
 796 gestion (GP\*) and focused proofreading performed automatically on  
 797 three drosophila datasets. The datasets are part of the MICCAI  
 798 2016 CREMI challenge and publicly available. We measure per-  
 799 formance as adapted Rand error (the lower, the better). GP\* is able  
 800 to correct the initial segmentation further than FP. Our GP\* scores  
 801 places us XXnd on the CREMI leaderboard.

802 We then use the trained GP\* and FP classifiers to eval-  
 803 uate proofreading automatically. Since ground truth label-  
 804 ing is not available, the evaluation is performed by sub-  
 805 mitting our results to the CREMI leaderboard. Again, we  
 806 use adapted Rand error to quantify the performance. Fig. 8  
 807 shows the results for each of the A+,B+, and C+ datasets.  
 808 The performance of GP\* is significantly better than FP and  
 809 places us XXnd on the CREMI leaderboard.

## 5. Quantitative Results

### 6. Conclusions

The task of automatic cell boundary segmentation is dif-  
 810 ficult, and trying to improve such segmentations automatic-  
 811 ally as a post-process through merge and split error cor-  
 812 rection is, in principle, no different than trying to improve  
 813 the underlying cell boundary segmentation. Due to the task  
 814 difficulty, manual proofreading of connectomics segmen-  
 815 tations is necessary, but it is a time consuming and error-prone  
 816 task. Humans are the bottleneck and minimizing the manual  
 817 labor is the goal. We have addressed this problem through  
 818 training a convolutional neural network to detect ambiguous  
 819 regions from labeled data—in effect, by finding a non-linear  
 820 mapping between image and segmentation data. This al-  
 821 lows us to identify merge and split errors with better per-  
 822 formance than existing systems. Our experiments have shown  
 823 that guided proofreading has the potential to reduce the bot-  
 824 tleneck in the analysis of large connectomics datasets. To  
 825 encourage testing of our proposed architecture and replicate  
 826 our experiments, we provide our framework and data as free  
 827 and open research at (link omitted for review).

## References

- [1] IEEE ISBI challenge: SNEMI3D - 3D segmentation of neu-  
 835 rites in EM images. [http://brainiac2.mit.edu/  
 836 SNEMI3D](http://brainiac2.mit.edu/SNEMI3D), 2013. Accessed on 11/01/2016. 1, 2
- [2] Neuroproof: Flyem tool, hhmi / janelia farm research  
 837 campus. [https://github.com/janelia-flyem/  
 838 NeuroProof](https://github.com/janelia-flyem/NeuroProof), 2013. Accessed on 03/15/2106. 2, 3
- [3] A. K. Al-Awami, J. Beyer, D. Haehn, N. Kasthuri, J. W.  
 841 Lichtman, H. Pfister, and M. Hadwiger. Neuroblocks - vi-  
 842 sual tracking of segmentation and proofreading for large con-  
 843 nectomics projects. *IEEE Transactions on Visualization and  
 844 Computer Graphics*, 22(1):738–746, Jan 2016. 2
- [4] J. Anderson, S. Mohammed, B. Grimm, B. Jones, P. Ko-  
 845 shlevoy, T. Tasdizen, R. Whitaker, and R. Marc. The Viking  
 846 Viewer for connectomics: Scalable multi-user annotation  
 847 and summarization of large volume data sets. *Journal of Mi-  
 848 croscopy*, 241(1):13–28, 2011. 2
- [5] ANON. Anon. ANON, 2016.
- [6] J. A. Bogovic, G. B. Huang, and V. Jain. Learned versus  
 851 hand-designed feature representations for 3d agglomeration.  
 852 *CoRR*, abs/1312.6159, 2013. 2, 3, 4
- [7] D. B. Chklovskii, S. Vitaladevuni, and L. K. Scheffer. Semi-  
 854 automated reconstruction of neural circuits using electron  
 855 microscopy. *Current Opinion in Neurobiology*, 20(5):667 –  
 856 675, 2010. Neuronal and glial cell biology New technolo-  
 857 gies. 2
- [8] D. C. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmid-  
 858 huber. Deep neural networks segment neuronal membranes  
 859 in electron microscopy images. In *NIPS*, 2012.
- [9] R. J. Giuly, K.-Y. Kim, and M. H. Ellisman. DP2: Dis-  
 861 tributed 3D image segmentation using micro-labor work-  
 862 force. *Bioinformatics*, 29(10):1359–1360, 2013. 2

- 864 [10] D. Haehn, S. Knowles-Barley, M. Roberts, J. Beyer,  
865 N. Kasthuri, J. Lichtman, and H. Pfister. Design and eval-  
866 uation of interactive proofreading tools for connectomics.  
867 *IEEE Transactions on Visualization and Computer Graph-  
868 ics (Proc. IEEE SciVis 2014)*, 20(12):2466–2475, 2014. 1, 2,  
869 3, 5, 6, 7
- 870 [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning  
871 for image recognition. In *The IEEE Conference on Computer  
872 Vision and Pattern Recognition (CVPR)*, June 2016. 3
- 873 [12] V. Jain, B. Bollmann, M. Richardson, D. Berger,  
874 M. Helmstädt, K. Briggman, W. Denk, J. Bowden,  
875 J. Mendenhall, W. Abraham, K. Harris, N. Kasthuri, K. Hay-  
876 worth, R. Schalek, J. Tapia, J. Lichtman, and S. Seung.  
877 Boundary learning by optimization with topological con-  
878 straints. In *Proc. IEEE CVPR 2010*, pages 2488–2495, 2010.  
879 1, 2
- 880 [13] Janelia Farm. Raveler. <https://openwiki.janelia.org/wiki/display/flyem/Raveler>, 2014. Ac-  
881 cessed on 11/01/2016. 1, 2
- 882 [14] A. Karimov, G. Mistelbauer, T. Auzinger, and S. Bruck-  
883 ner. Guided volume editing based on histogram dissimilarity.  
884 *Computer Graphics Forum*, 34(3):91–100, May 2015. 3, 4
- 885 [15] N. Kasthuri, K. J. Hayworth, D. R. Berger, R. L. Schalek,  
886 J. A. Conchello, S. Knowles-Barley, D. Lee, A. Vázquez-  
887 Reina, V. Kaynig, T. R. Jones, et al. Saturated reconstruction  
888 of a volume of neocortex. *Cell*, 162(3):648–661, 2015. 5, 7
- 889 [16] V. Kaynig, T. Fuchs, and J. Buhmann. Neuron geometry  
890 extraction by perceptual grouping in sstem images. In *Proc.  
891 IEEE CVPR*, pages 2902–2909, 2010. 2
- 892 [17] V. Kaynig, A. Vazquez-Reina, S. Knowles-Barley,  
893 M. Roberts, T. R. Jones, N. Kasthuri, E. Miller, J. Lichtman,  
894 and H. Pfister. Large-scale automatic reconstruction of  
895 neuronal processes from electron microscopy images.  
896 *Medical image analysis*, 22(1):77–88, 2015. 1
- 897 [18] J. S. Kim, M. J. Greene, A. Zlateski, K. Lee, M. Richardson,  
898 S. C. Turaga, M. Purcaro, M. Balkam, A. Robinson, B. F. Be-  
899 habadi, M. Campos, W. Denk, H. S. Seung, and EyeWirers.  
900 Space-time wiring specificity supports direction selectivity  
901 in the retina. *Nature*, 509(7500):331336, May 2014. 2
- 902 [19] S. Knowles-Barley, M. Roberts, N. Kasthuri, D. Lee, H. Pfis-  
903 ter, and J. W. Lichtman. Mojo 2.0: Connectome annotation  
904 tool. *Frontiers in Neuroinformatics*, (60), 2013. 1
- 905 [20] K. Lee, A. Zlateski, A. Vishwanathan, and H. S. Seung. Re-  
906 cursive training of 2d-3d convolutional networks for neu-  
907 ronal boundary detection. *arXiv preprint arXiv:1508.04843*,  
908 2015. 2
- 909 [21] J. W. Lichtman and W. Denk. The big and the small:  
910 Challenges of imaging the brain’s circuits. *Science*,  
911 334(6056):618–623, 2011. 7
- 912 [22] T. Liu, C. Jones, M. Seyedhosseini, and T. Tasdizen. A mod-  
913 ular hierarchical approach to 3D electron microscopy image  
914 segmentation. *Journal of Neuroscience Methods*, 226(0):88  
915 – 102, 2014. 1, 2
- 916 [23] J. Masci, A. Giusti, D. C. Ciresan, G. Fricout, and J. Schmid-  
917 huber. A fast learning algorithm for image segmentation with  
max-pooling convolutional networks. In *ICIP*, 2013.
- 918 [24] J. Nunez-Iglesias, R. Kennedy, T. Parag, J. Shi, and D. B.  
919 Chklovskii. Machine learning of hierarchical clustering to  
920 segment 2D and 3D images. *PLoS ONE*, 8(8):e71715+,  
921 2013. 2
- 922 [25] J. Nunez-Iglesias, R. Kennedy, S. M. Plaza, A. Chakraborty,  
923 and W. T. Katz. Graph-based active learning of agglomera-  
924 tion (GALA): A python library to segment 2D and 3D neu-  
925 roimages. *Frontiers in Neuroinformatics*, 8(34), 2014. 1,  
926 2
- 927 [26] H. Peng, F. Long, T. Zhao, and E. Myers. Proof-editing is the  
928 bottleneck of 3D neuron reconstruction: The problem and  
929 solutions. *Neuroinformatics*, 9(2-3):103–105, 2011. 1, 2, 3
- 930 [27] S. M. Plaza. Focused Proofreading: Efficiently Extracting  
931 Connectomes from Segmented EM Images, Sept. 2014. 2, 3,  
932 5, 6, 7
- 933 [28] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolu-  
934 tional networks for biomedical image segmentation. *CoRR*,  
935 abs/1505.04597, 2015. 2
- 936 [29] S. Saalfeld, A. Cardona, V. Hartenstein, and P. Tomancák.  
937 CATMAID: collaborative annotation toolkit for massive  
938 amounts of image data. *Bioinformatics*, 25(15):1984–1986,  
939 2009. 2
- 940 [30] R. Sicat, M. Hadwiger, and N. J. Mitra. Graph abstraction for  
941 simplified proofreading of slice-based volume segmentations.  
942 In *EUROGRAPHICS Short Paper*, 2013. 1, 2
- 943 [31] R. Unnikrishnan, C. Pantofaru, and M. Hebert. A measure  
944 for objective evaluation of image segmentation algorithms.  
945 pages 34–, 2005. 7
- 946 [32] M. G. Uzunbas, C. Chen, and D. Metaxas. An efficient con-  
947 ditional random field approach for automatic and interactive  
948 neuron segmentation. *Medical Image Analysis*, 27:31 – 44,  
949 2016. Discrete Graphical Models in Biomedical Image Anal-  
950 ysis. 3
- 951 [33] A. Vázquez-Reina, M. Gelbart, D. Huang, J. Lichtman,  
952 E. Miller, and H. Pfister. Segmentation fusion for connec-  
953 toomics. In *Proc. IEEE ICCV*, pages 177–184, Nov 2011. 2
- 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971