

000
001
002003

Guided Proofreading of Automatic Segmentations for Connectomics

004
005
006
007
008
009
010
011

Anonymous CVPR submission

012

Abstract

013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

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.

Paper ID 0947

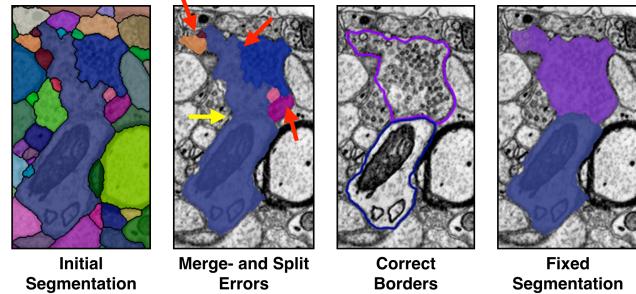
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

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.

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 $\approx 4\text{ nm}$ per pixel to identify cells, synapses, and vesicles. For 40 nm thick sections, a 1 mm^3 volume of brain contains 10^{15} voxels, or 1 petabyte of data. With so much data, manual annotation is infeasible, and automatic annotation methods are needed [12, 22, 25, 17].

Automatic annotation by segmentation and classification of brain tissue is challenging [1] and all the 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-

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.

2. Related Work

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

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

195 All these approaches make good progress; however, in
196 general, proofreading is still required to correct errors.

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

205 Recent papers have attacked the problem of proofreading
206 massive datasets through crowdsourcing with novices [29,
207 4, 9]. One popular platform is EyeWire, by Kim *et al.* [18],
208 where participants earn virtual rewards for merging over-
209 segmented labeling to reconstruct retina cells.

210 Between expert systems and online games sit Mojo and
211 Dojo, by Haehn *et al.* [10, 3], which use simple scribble in-
212 terfaces for error correction. Dojo extends this to distributed
213 proofreading via a minimalistic web-based user interface.
214 The authors define requirements for general proofreading

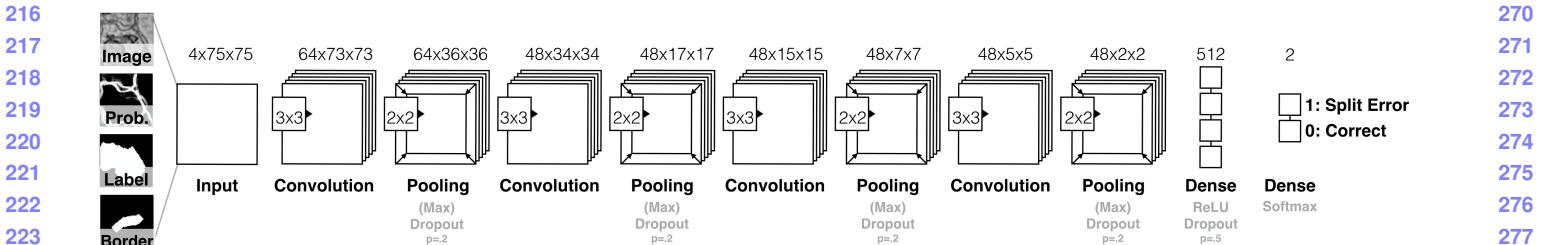


Figure 2. We build the guided proofreading classifiers using a traditional CNN architecture. The network is based on four convolutional layers, each followed by max pooling as well as dropout regularization. The 4-channel input patches are rated as either correct splits or as split errors.

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

We first describe our classifier for detecting split errors which is based on a convolutional neural network (CNN). We detail the CNN architecture, input features and the training method. We then describe how the same classifier can be used to detect merge errors and how we create potential corrections. The classifiers are integrated into an existing proofreading workflow as reported after. Finally, we explore an active label suggestion method which reorders the ranking obtained by our classifiers and maximizes the information gain provided by each potential correction.

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. Our approach was inspired by Bogovic *et al.* [6] but works 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 consideration for the decision making process. For this, we grab a 75×75 pixel patch at the center of an existing boundary. This covers approximately 80% of all boundaries in real-world connectomics data with nanometer resolution. If the boundary edge is not fully covered, we sample up to 10 non-overlapping patches along the boundary and combine the resulting score by weighted averaging based on boundary length coverage per patch. In their paper, Bogovich *et al.* propose to use grayscale image data, corresponding boundary probabilities, and a single binary mask combining the two neighboring labels as features for their recursive neural network [6]. We are building on this set of features as inputs for our CNN and create a stacked pixel patch. However, we observed that the boundary probability information generated from EM images is often misleading due to noise or artefacts in the data. This can result in merge errors within the automatic segmentation. To better direct our classifier to train on the true boundary edge, we extract the border between two segments. We then dilate this border by 5 pixels to consider slight edge ambiguities and use this additional binary mask as another feature to create a 4-channel input patch. Fig. 3 shows examples of correct and erroneous feature patches and their corresponding automatic segmentation and ground truth.

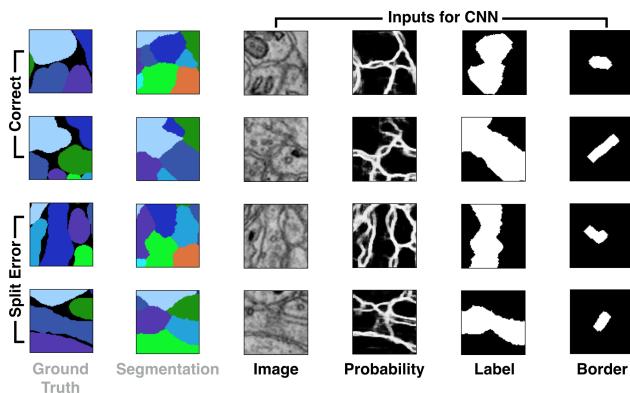


Figure 3. Example inputs for learning correct splits and split errors as reflected in the segmentation relative to the ground truth. Image, membrane probabilities, merged binary labels, and a dilated border mask are combined to 4-channel input patches.

Training. To initially train our network, we use the blue 3-cylinder mouse cortex volume of Kasthuri *et al.* [15]

($2048 \times 2048 \times 300$ voxels). The tissue is dense mammalian 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 dataset¹. 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.

	cost [m]	Val. loss	Val. acc.	Test acc.	Prec./Recall	F1 Score
Guided Proofreading 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
Focused Proofreading 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.

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 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. 4 displays re-

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

ceiver 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.

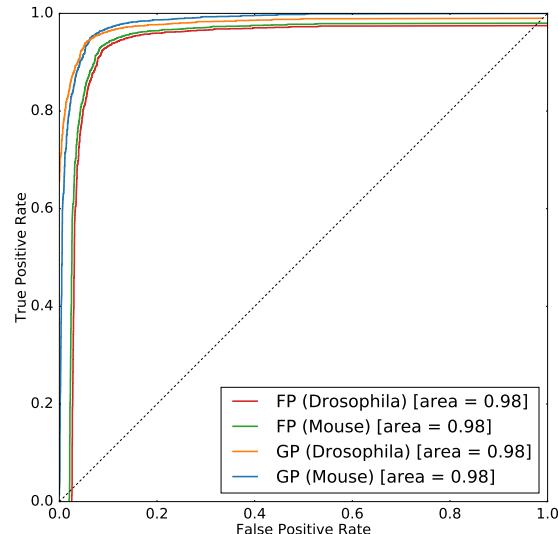


Figure 4. ROC performance of guided proofreading (GP) and focused proofreading (FP) trained separately on mouse and drosophila brain images. The area under the curve indicates better performance for GP.

3.2. Merge Error Detection

Identification and correction of merge errors is more challenging, because we must look inside segmentation regions for missing or incomplete boundaries and then propose the correct boundary. However, we can reuse the same trained CNN for this task. Similar to guided volume editing by Karimov *et al.* [14] we generate potential borders within a segment. For each segmentation label, we dilate the label by 20 pixel and generate 30 potential boundaries through the region by randomly placing watershed seed points at opposite sides of the label boundary. For watershed, we use the inverted gray scale EM image as features. This yields 30 corresponding splits. Dilation of the segment prior to watershed is motivated by our observation that the generated split then actually hogs the real membrane boundary. These boundaries are then individually rated using our split error classifier. For this, we invert the probability score meaning that a correct split (previously encoded as $p = 0$) is most likely a candidate for a merge error (now encoded as $p = 1$). In other words, if a generated boundary is ranked as correct, it probably should be in the segmentation. Fig. 5 illustrates this procedure.

3.3. Error Correction

We use the proposed classifiers in combination to perform corrections of split and merge errors in automatic seg-

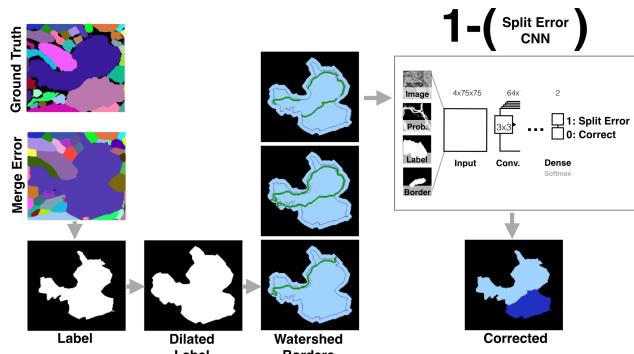


Figure 5. Merge errors are identified by generating randomly seeded watershed borders within a dilated label segment. These borders then are individually rated using the split error CNN by inverting the probability score. This way, a confident rating for a correct split most likely indicates the missing border of the merge error and can be used for correcting the labeling.

mentations. For this, we first perform merge error detection for all existing segments in a data set and store the inverted ranking $1 - p$. We then sort the rankings and loop through all of them in greedy fashion, starting with the most likely error. If the inverted score $1 - p$ of a merge error is higher than a threshold p_t , we mark this merge error for correction. The merge error detection yields the potential new boundary and we can modify the segmentation data to create a new segment accordingly. Depending on the variation of guided proofreading as described in section 4, we either perform the correction directly (automatic GP) or we have a user accept or reject the correction (simulated GP or human novice/expert). Once all merge errors are corrected, we perform split error detection. For this, we perform split error detection and store the ranking p for all existing segments in the for merge errors corrected segmentation. We then sort the rankings and again loop through all of them - the most likely error first. If the score p is higher than a threshold p_t , we mark the split error for correction. Potential split errors are identified at the border of two segments. This reduces the correction to merging the segments and is therefore trivial. Similarly to merge errors, we either perform the correction directly or present a user with a yes/no decision. Our experiments have shown that the threshold p_t is the same for merge and split errors which makes sense for a balanced classifier. The only exception is when a user drives the correction process: we then set p_t for split errors to 0 to let the user inspect every possible split error. Inspecting all merge errors is not possible for users due to the sheer amount of generated borders.

3.4. Application

Guided proofreading is integrated into an existing workflow for large connectomics data. The GP system is web-based and is designed with a minimalistic user interface

540 showing three components. First, we show the outline of
541 the current labeling of a cell boundary and its proposed cor-
542 rection on top of the EM image data. For the user, it is not
543 possible to distinguish the current labeling and the proposed
544 correction to avoid selection bias. Second, we show a solid
545 overlay of the current and proposed labeling. And finally,
546 to provide context, we show a larger area of the EM image
547 where the potentially erroneous region is highlighted. User
548 interaction is simple and involves one mouse click on either
549 the current labeling or the correction. After interaction, the
550 next potential error is shown. We provide a screenshot of
551 the application as part of the supplementary material.
552

553 3.5. Active Label Suggestion

554 In an interactive setting, one way to present patches to
555 the user for proofreading is to order them by the confidence
556 probability of the GP classifier. However, in an active learn-
557 ing setting, where the network is retrained repeatedly on
558 new label evidence, this approach is less likely to decrease
559 segmentation error as, with the new labels, we are only rein-
560 forcing what the network already has a high confidence in.
561 Instead, we apply active label suggestion to guide the user
562 into labeling patches which will be more informative to re-
563 training, and so overall decrease VI faster within the proof-
564 reading cycle of label → train → label. For each patch, we
565 remove the softmax classification layer and look at the ac-
566 tivation weights associated with the last dense layer. These
567 become a high-dimensional feature vector. Then, we adapt
568 Anon *et al.* [5] to provide label suggestions based on fea-
569 tures from the learned CNN, which is based on maximiz-
570 ing the average information gain provided by a candidate
571 patch to label. A second consideration is that each patch
572 labeled by the user provides evidence to other patches, e.g.,
573 correcting a split error redefines an entire boundary, from
574 which multiple candidate patch labelings could have been
575 drawn. As such, when the user labels a patch, we consider
576 all ‘knock-on’ effect patches as also being labeled, and feed
577 these into the active label suggestion system similarly. In
578 section 4, we report the difference in performance from us-
579 ing active label suggestion rather than confidence ordering
580 when presenting patches to the user. These results are with-
581 out retraining the network after new labelings: this should
582 improve results, but would have to be batched to reduce
583 computational load; hence, we leave this for future work.
584

585 4. Evaluation

586 We evaluate guided proofreading (GP) on multiple real-
587 world connectomics datasets of different species. In partic-
588 ular, we evaluate GP on two datasets of mouse brain and
589 three datasets of fruitfly brain (*drosophila*). For compari-
590 son, we choose the fully interactive proofreading software
591 *Dojo* by Haehn *et al.* [10] as well as the aided proofreading
592 framework *focused proofreading* (FP) by Plaza [27]. We
593

594 first describe the evalution on mouse brain data and then the
595 evaluation on drosophila brain.
596

597 4.1. Mouse Brain

598 Mouse brain is a common target for connectomics
599 research because the structural proportions are similar to
600 human brains [21]. For our first experiment we recruited
601 novice and expert participants as part of a quantitative user
602 study. Our second experiment is performed on a larger
603 dataset and we evaluate a simulated user.
604

605 **User study.** Recently, Haehn *et al.* evaluated the in-
606 teractive proofreading tools Raveler, Mojo, and Dojo as
607 part of an experiment with novice users [10]. The partic-
608 ipants corrected an automatic segmentation with merge
609 and split errors. The dataset was the most representative
610 sub-volume (based on object size histograms) of a larger
611 connectomics dataset and 400x400x10 voxels in size. The
612 participants were given a fixed time frame of 30 minutes
613 to perform the correction interactively. While participants
614 clearly struggled with the proofreading task, the best
615 performing tool in their evaluation was Dojo. The dataset
616 including manually labeled ground truth and the results of
617 Haehn *et al.* are publicly available. This means we are able
618 to use their findings as a baseline for comparison of GP for
619 novices. In particular, we use the best performing user of
620 Dojo who was truly an outlier as reported by Haehn *et al.*
621

622 Since interactive proofreading most likely yields lower
623 performance than aided proofreading, we also compare
624 against FP by Plaza [27] which is integrated in Raveler and
625 freely available. For FP we consulted an expert to obtain
626 the best possible parameters as shown in table 1. Besides
627 performance by novices, we are also interested in expert
628 proofreading performance. Therefore, we design between-
629 subjects experiments for 20 novice users and separately, for
630 6 expert users using the exact same conditions as Haehn *et*
631 *al.* The recruiting, consent and debriefing process is further
632 described in the supplementary material. We randomly as-
633 sign 10 novices to GP with active label suggestion (GP*)
634 and 10 novices to FP. For the expert experiment, we assign
635 accordingly. In addition to human performance, we also
636 evaluate automatic GP, automatic GP with active label sug-
637 gestion (GP*) and automatic FP. Due to the automatic na-
638 ture, we do not enforce the 30 minute time limit but we stop
639 once our probability threshold of $p_t = .95$ is reached. This
640 value was observed as stable in previous experiments using
641 automatic GP (see supplementary material). To measure
642 proofreading performance in comparison to ground truth,
643 we use the adapted Rand error (aRE) metric [31]. aRE is a
644 measure of dissimilarity, related to introduced errors, mean-
645 ing lower scores are better.
646

647 The results of our comparisons are shown in the first
648 row of Fig. 6. In all cases, GP* is able to correct the
649

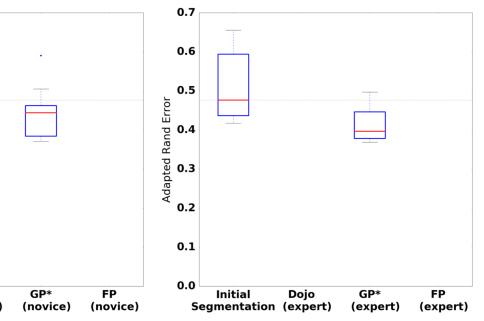
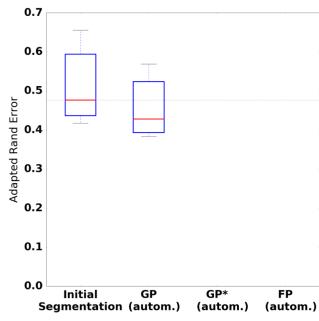
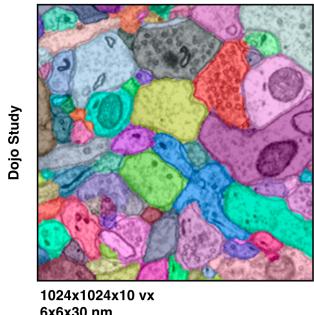
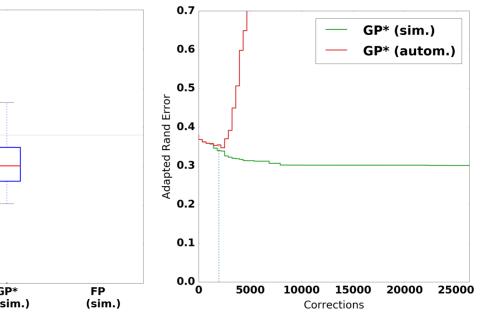
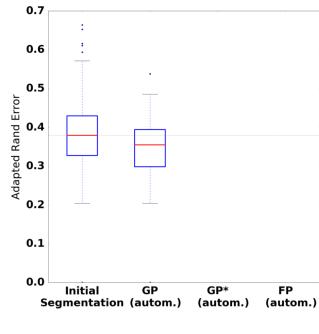
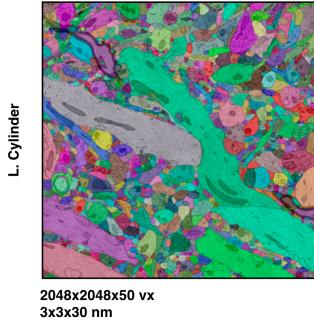
648
649
650
651
652
653
654
655
656
657
658702
703
704
705
706
707
708
709
710
711
712659
660
661
662
663
664
665
666
667
668
669713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

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.

677

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 2048x2048x50 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

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
electron microscopy images (CREMI)². 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 datasets. We use 300 slices of the A,B,C samples for training and validation, and 75 slices for testing. This results in YYY correct and ZZZ split error patches (respectively, XXX and YYY for testing). The architecture and all parameters of our classifiers stay the same. The trained GP classifier exhibits a reasonable performance on the testing data as seen in Fig. 4.

We then use the trained GP* and FP classifiers to evaluate proofreading automatically. Since ground truth labeling is not available, the evaluation is performed by submitting our results to the CREMI leaderboard. Again, we use adapted Rand error to quantify the performance. Fig. 7 shows the results for each of the A+,B+, and C+ datasets. The performance of GP* is significantly better than FP and places us XXnd on the CREMI leaderboard.

5. Quantitative Results

6. Conclusions

The task of automatic cell boundary segmentation is difficult, and trying to improve such segmentations automatically as a post-process through merge and split error correction is, in principle, no different than trying to improve the underlying cell boundary segmentation. Due to the task difficulty, manual proofreading of connectomics segmentations is necessary, but it is a time consuming and error-prone task. Humans are the bottleneck and minimizing the manual labor is the goal. We have addressed this problem through training a convolutional neural network to detect ambiguous regions from labeled data—in effect, by finding a non-linear mapping between image and segmentation data. This allows us to identify merge and split errors with better performance than existing systems. Our experiments have shown that guided proofreading has the potential to reduce the bottleneck in the analysis of large connectomics datasets. To encourage testing of our proposed architecture and replicate our experiments, we provide our framework and data as free and open research at (link omitted for review).

References

- [1] IEEE ISBI challenge: SNEMI3D - 3D segmentation of neurites in EM images. <http://brainiac2.mit.edu/>

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

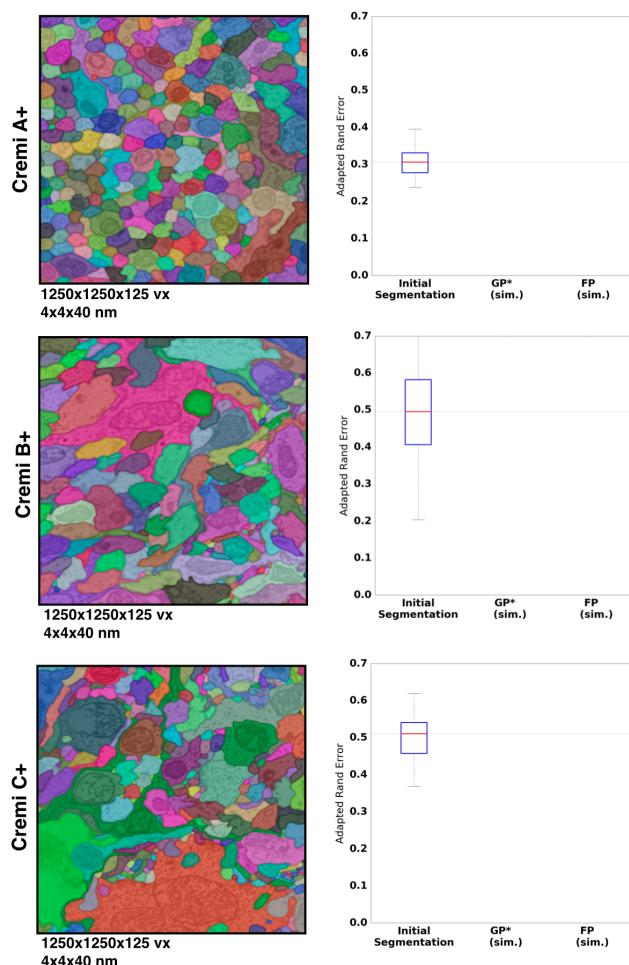


Figure 7. Results of guided proofreading with active label suggestion (GP*) and focused proofreading performed automatically on three drosophila datasets. The datasets are part of the MICCAI 2016 CREMI challenge and publicly available. We measure performance as adapted Rand error (the lower, the better). GP* is able to correct the initial segmentation further than FP. Our GP* scores places us XXnd on the CREMI leaderboard.

- [SNEMI3D](#), 2013. Accessed on 11/01/2016. 1, 2
- [2] Neuroproof: Flyem tool, hhmi / janelia farm research campus. <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. Lichtman, H. Pfister, and M. Hadwiger. Neuroblocks - visual tracking of segmentation and proofreading for large connectomics projects. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):738–746, Jan 2016. 2
- [4] J. Anderson, S. Mohammed, B. Grimm, B. Jones, P. Koshevoy, T. Tasdizen, R. Whitaker, and R. Marc. The Viking Viewer for connectomics: Scalable multi-user annotation and summarization of large volume data sets. *Journal of Microscopy*, 241(1):13–28, 2011. 2
- [5] ANON. Anon. ANON, 2016. 6

- 864 [6] J. A. Bogovic, G. B. Huang, and V. Jain. Learned versus
865 hand-designed feature representations for 3d agglomeration.
866 *CoRR*, abs/1312.6159, 2013. 2, 3, 4
- 867 [7] D. B. Chklovskii, S. Vitaladevuni, and L. K. Scheffer. Semi-
868 automated reconstruction of neural circuits using electron
869 microscopy. *Current Opinion in Neurobiology*, 20(5):667 –
870 675, 2010. Neuronal and glial cell biology New technologies. 2
- 871 [8] D. C. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmid-
872 huber. Deep neural networks segment neuronal membranes
873 in electron microscopy images. In *NIPS*, 2012.
- 874 [9] R. J. Giuly, K.-Y. Kim, and M. H. Ellisman. DP2: Dis-
875 tributed 3D image segmentation using micro-labor work-
876 force. *Bioinformatics*, 29(10):1359–1360, 2013. 2
- 877 [10] D. Haehn, S. Knowles-Barley, M. Roberts, J. Beyer,
878 N. Kasthuri, J. Lichtman, and H. Pfister. Design and eval-
879 uation of interactive proofreading tools for connectomics.
880 *IEEE Transactions on Visualization and Computer Graph-
881 ics (Proc. IEEE SciVis 2014)*, 20(12):2466–2475, 2014. 1, 2,
882 3, 6, 7
- 883 [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning
884 for image recognition. In *The IEEE Conference on Computer
885 Vision and Pattern Recognition (CVPR)*, June 2016. 3
- 886 [12] V. Jain, B. Bollmann, M. Richardson, D. Berger,
887 M. Helmstädtter, K. Briggman, W. Denk, J. Bowden,
888 J. Mendenhall, W. Abraham, K. Harris, N. Kasthuri, K. Hay-
889 worth, R. Schalek, J. Tapia, J. Lichtman, and S. Seung.
890 Boundary learning by optimization with topological con-
891 straints. In *Proc. IEEE CVPR 2010*, pages 2488–2495, 2010.
892 1, 2
- 893 [13] Janelia Farm. Raveler. <https://openwiki.janelia.org/wiki/display/flyem/Raveler>, 2014.
894 Accessed on 11/01/2016. 1, 2
- 895 [14] A. Karimov, G. Mistelbauer, T. Auzinger, and S. Bruck-
896 ner. Guided volume editing based on histogram dissimilarity.
897 *Computer Graphics Forum*, 34(3):91–100, May 2015. 3, 5
- 898 [15] N. Kasthuri, K. J. Hayworth, D. R. Berger, R. L. Schalek,
899 J. A. Conchello, S. Knowles-Barley, D. Lee, A. Vázquez-
900 Reina, V. Kaynig, T. R. Jones, et al. Saturated reconstruction
901 of a volume of neocortex. *Cell*, 162(3):648–661, 2015. 4, 7
- 902 [16] V. Kaynig, T. Fuchs, and J. Buhmann. Neuron geometry
903 extraction by perceptual grouping in sstem images. In *Proc.
904 IEEE CVPR*, pages 2902–2909, 2010. 2
- 905 [17] V. Kaynig, A. Vazquez-Reina, S. Knowles-Barley,
906 M. Roberts, T. R. Jones, N. Kasthuri, E. Miller, J. Lichtman,
907 and H. Pfister. Large-scale automatic reconstruction of
908 neuronal processes from electron microscopy images.
909 *Medical image analysis*, 22(1):77–88, 2015. 1
- 910 [18] J. S. Kim, M. J. Greene, A. Zlateski, K. Lee, M. Richardson,
911 S. C. Turaga, M. Purcaro, M. Balkam, A. Robinson, B. F. Be-
912 habadi, M. Campos, W. Denk, H. S. Seung, and EyeWirers.
913 Space-time wiring specificity supports direction selectivity
914 in the retina. *Nature*, 509(7500):331336, May 2014. 2
- 915 [19] S. Knowles-Barley, M. Roberts, N. Kasthuri, D. Lee, H. Pfis-
916 ter, and J. W. Lichtman. Mojo 2.0: Connectome annotation
917 tool. *Frontiers in Neuroinformatics*, (60), 2013. 1
- 918 [20] K. Lee, A. Zlateski, A. Vishwanathan, and H. S. Seung. Re-
919 cursive training of 2d-3d convolutional networks for neu-
920 ronal boundary detection. *arXiv preprint arXiv:1508.04843*,
921 2015. 2
- 922 [21] J. W. Lichtman and W. Denk. The big and the small:
923 Challenges of imaging the brain’s circuits. *Science*,
924 334(6056):618–623, 2011. 6
- 925 [22] T. Liu, C. Jones, M. Seyedhosseini, and T. Tasdizen. A mod-
926 ular hierarchical approach to 3D electron microscopy image
927 segmentation. *Journal of Neuroscience Methods*, 226(0):88
928 – 102, 2014. 1, 2
- 929 [23] J. Masci, A. Giusti, D. C. Ciresan, G. Fricout, and J. Schmid-
930 huber. A fast learning algorithm for image segmentation with
931 max-pooling convolutional networks. In *ICIP*, 2013.
- 932 [24] J. Nunez-Iglesias, R. Kennedy, T. Parag, J. Shi, and D. B.
933 Chklovskii. Machine learning of hierarchical clustering to
934 segment 2D and 3D images. *PLoS ONE*, 8(8):e71715+,
935 2013. 2
- 936 [25] J. Nunez-Iglesias, R. Kennedy, S. M. Plaza, A. Chakraborty,
937 and W. T. Katz. Graph-based active learning of agglomera-
938 tion (GALA): A python library to segment 2D and 3D neu-
939 roimages. *Frontiers in Neuroinformatics*, 8(34), 2014. 1,
940 2
- 941 [26] H. Peng, F. Long, T. Zhao, and E. Myers. Proof-editing is the
942 bottleneck of 3D neuron reconstruction: The problem and
943 solutions. *Neuroinformatics*, 9(2-3):103–105, 2011. 1, 2, 3
- 944 [27] S. M. Plaza. Focused Proofreading: Efficiently Extracting
945 Connectomes from Segmented EM Images, Sept. 2014. 2, 3,
946 4, 6
- 947 [28] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolu-
948 tional networks for biomedical image segmentation. *CoRR*,
949 abs/1505.04597, 2015. 2
- 950 [29] S. Saalfeld, A. Cardona, V. Hartenstein, and P. Tomančák.
951 CATMAID: collaborative annotation toolkit for massive
952 amounts of image data. *Bioinformatics*, 25(15):1984–1986,
953 2009. 2
- 954 [30] R. Sicat, M. Hadwiger, and N. J. Mitra. Graph abstraction for
955 simplified proofreading of slice-based volume segmentation.
956 In *EUROGRAPHICS Short Paper*, 2013. 1, 2
- 957 [31] R. Unnikrishnan, C. Pantofaru, and M. Hebert. A measure
958 for objective evaluation of image segmentation algorithms.
959 pages 34–, 2005. 6
- 960 [32] M. G. Uzunbas, C. Chen, and D. Metaxas. An efficient con-
961 ditional random field approach for automatic and interactive
962 neuron segmentation. *Medical Image Analysis*, 27:31 – 44,
963 2016. Discrete Graphical Models in Biomedical Image Anal-
964 ysis. 3
- 965 [33] A. Vázquez-Reina, M. Gelbart, D. Huang, J. Lichtman,
966 E. Miller, and H. Pfister. Segmentation fusion for connec-
967 toomics. In *Proc. IEEE ICCV*, pages 177–184, Nov 2011. 2
- 968 969 970 971