
Disentangling Human Error from Ground Truth in Segmentation of Medical Images

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent years have seen increasing use of supervised learning methods for segmenta-
2 tion tasks. However, the predictive performance of these algorithms depends on the
3 quality of labels. This problem is particularly pertinent in the medical image domain,
4 where both the annotation cost and inter-observer variability are high. In a typical la-
5 bel acquisition process, different human experts provide their estimates of the “true”
6 segmentation labels under the influence of their own biases and competence levels.
7 Treating these noisy labels blindly as the ground truth limits the performance that
8 automatic segmentation algorithms can achieve. In this work, we present a method
9 for jointly learning, from purely noisy observations alone, the reliability of individual
10 annotators and the true segmentation label distributions, using two coupled CNNs.
11 The separation of the two is achieved by encouraging the estimated annotators to
12 be maximally unreliable while achieving high fidelity with the noisy training data.
13 We first define a toy segmentation dataset based on MNIST and study the properties
14 of the proposed algorithm. We then demonstrate the utility of the method on three
15 public medical imaging segmentation datasets with simulated (when necessary) and
16 real diverse annotations: 1) MSLSC (multiple-sclerosis lesions); 2) BraTS (brain
17 tumours); 3) LIDC-IDRI (lung abnormalities). In all cases, our method outperforms
18 competing methods and relevant baselines particularly in cases where the number
19 of annotations is small and the amount of disagreement is large. The experiments
20 also show strong ability to capture the complex spatial characteristics of annotators’
21 mistakes, which could be potentially utilised for the purpose of education.

22 1 Introduction

23 Segmentation of anatomical structures in medical images is known to suffer from high inter-reader
24 variability [1, 2, 3, 4, 5], affecting limiting the performance of downstream supervised machine
25 learning models. This problem is particularly prominent in the medical domain where the labelled
26 data is commonly scarce due to the high cost of annotations. For instance, accurate identification of
27 multiple sclerosis (MS) lesions in MRIs is difficult even for experienced experts due to variability in
28 lesion location, size, shape and anatomical variability across patients [6]. Another example [4] reports
29 the average inter-reader variability in the range 74-85% for glioblastoma (a type of brain tumour)
30 segmentation. Further aggravated by differences in biases and levels of expertise, segmentation
31 annotations of structures in medical images suffer from high annotation variations [7]. In consequence,
32 despite the present abundance of medical imaging data thanks to over two decades of digitisation,
33 the world still remains relatively short of access to data with curated labels [8], that is amenable to
34 machine learning, necessitating intelligent methods to learn robustly from such noisy annotations.

35 To mitigate inter-reader variations, different pre-processing techniques are commonly used to curate
36 segmentation annotations by fusing labels from different experts. The most basic yet popular approach
37 is based on the majority vote where the most representative opinion of the experts is treated as the
38 ground truth (GT). A smarter version that accounts for similarity of classes has proven effective in

39 aggregation of brain tumour segmentation labels [4]. A key limitation of such approaches, however,
40 is that all experts are assumed to be equally reliable. Warfield *et al.* [9] proposed a label fusion method,
41 called STAPLE that explicitly models the reliability of individual experts and uses that information to
42 “weigh” their opinions in the label aggregation step. After consistent demonstration of its superiority
43 over the standard majority-vote pre-processing in multiple applications, STAPLE has become the go-to
44 label fusion method in the creation of public medical image segmentation datasets e.g., ISLES [10],
45 MSSeg [11], Gleason’19 [12] datasets. Asman *et al.* later extended this approach in [13] by accounting
46 for voxel-wise consensus to address the issue of under-estimation of annotators’ reliability. In [14],
47 another extension was proposed in order to model the reliability of annotators across different pixels
48 in images. More recently, within the context of multi-atlas segmentation problems [15] where image
49 registration is used to warp segments from labeled images (“atlases”) onto a new scan, STAPLE has
50 been enhanced in multiple ways to encode the information of the underlying images into the label
51 aggregation process. A notable example is STEP proposed in Cardoso *et al.* [16] who designed a
52 strategy to further incorporate the local morphological similarity between atlases and target images,
53 and different extensions of this approach such as [17, 18] have since been considered. However,
54 these previous label fusion approaches have a common drawback—they critically lack a mechanism
55 to integrate information across different training images. This fundamentally limits the remit of
56 applications to cases where each image comes with a reasonable number of annotations from multiple
57 experts, which can be prohibitively expensive in practice. Moreover, relatively simplistic functions
58 are used to model the relationship between observed noisy annotations, true labels and reliability of
59 experts, which may fail to capture complex characteristics of human annotators.

60 In this work, we introduce the first instance of an end-to-end supervised segmentation method that
61 jointly estimates, from noisy labels alone, the reliability of multiple human annotators and true
62 segmentation labels. The proposed architecture (Fig. 1) consists of two coupled CNNs where one
63 estimates the true segmentation probabilities and the other models the characteristics of individual
64 annotators (e.g., tendency to over-segmentation, mix-up between different classes, etc) by estimating
65 the pixel-wise confusion matrices (CMs) on a per image basis. Unlike STAPLE [9] and its variants,
66 our method models, and disentangles with deep neural networks, the complex mappings from the input
67 images to the annotator behaviours and to the true segmentation label. Furthermore, the parameters
68 of the CNNs are “global variables” that are optimised across different image samples; this enables
69 the model to disentangle robustly the annotators’ mistakes and the true labels based on correlations
70 between similar image samples, even when the number of available annotations is small per image
71 (e.g., a single annotation per image). In contrast, this would not be possible with STAPLE [9] and
72 its variants [14, 16] where the annotators’ parameters are estimated on every target image separately.

73 For evaluation, we first simulate a diverse range of annotator types on the MNIST dataset by performing
74 morphometric operations with Morpho-MNIST framework [19]. Then we demonstrate the potential
75 in several real-world medical imaging datasets, namely (i) MS lesion segmentation dataset (MSLSC)
76 from the ISBI 2015 challenge [20], (ii) Brain tumour segmentation dataset (BraTS) [4] and (iii)
77 Lung nodule segmentation dataset (LIDC-IDRI) [21]. Experiments on all datasets demonstrate
78 that our method consistently leads to better segmentation performance compared to widely adopted
79 label-fusion methods and other relevant baselines, especially when the number of available labels
80 for each image is low and the degree of annotator disagreement is high.

81 **2 Related Work**

82 The majority of algorithmic innovations in the space of *label aggregation for segmentation* have
83 uniquely originated from the medical imaging community, partly due to the prominence of the inter-
84 reader variability problem in the field, and the wide-reaching values of reliable segmentation methods.
85 The aforementioned methods based on the STAPLE-framework such as [9, 13, 14, 16, 22, 17, 17, 18, 23]
86 are based on generative models of human behaviours, where the latent variables of interest are the
87 unobserved true labels and the “reliability” of the respective annotators. Our method can be viewed
88 as an instance of translation of the STAPLE-framework to the supervised learning paradigm. As such,
89 our method produces a model that can segment test images without needing to acquire labels from
90 annotators or atlases unlike STAPLE and its local variants. Another key difference is that our method
91 is jointly trained on many different subjects while the STAPLE-variants are only fitted on a per-subject
92 basis. This means that our method is able to learn from correlations between different subjects, which
93 previous works have not attempted— for example, our method uniquely can estimate the reliability
94 and true labels even when there is only one label available per input image as shown later.

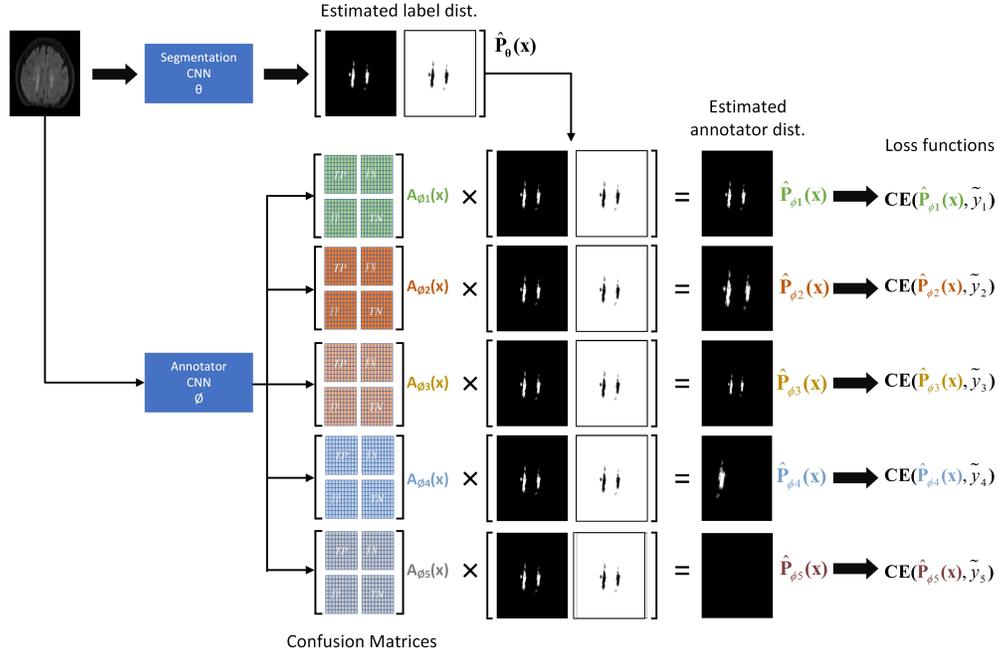


Figure 1: General schematic of the neural fusion network. The method consists of two components: (1) Segmentation network parametrised by θ that generates an estimate of the GT segmentation probabilities, $p_\theta(\mathbf{x})$ for the given input image \mathbf{x} ; (2) Annotator network, parametrised by ϕ , that estimates the CMs $\{\mathbf{A}_\phi^{(r)}(\mathbf{x})\}_{r=1}^n$ of the annotators. The segmentation probabilities of respective annotators $\hat{\mathbf{p}}_\phi^{(r)}(\mathbf{x}) := \mathbf{A}_\phi^{(r)}(\mathbf{x}) \cdot \mathbf{p}_\theta(\mathbf{x})$ are then computed. The model parameters $\{\theta, \phi\}$ are optimized to minimize the sum of five cross-entropy losses between each estimated annotator distribution $\mathbf{p}_\phi^{(r)}(\mathbf{x})$ and the noisy labels $\tilde{\mathbf{y}}^{(r)}$ observed from each annotator.

95 Our work also relates to a recent strand of methods that aim to generate a set of diverse and plausible
 96 segmentation proposals on a given image. Notably, probabilistic U-net [24] and its recent variants,
 97 PHiSeg [25] have shown that the aforementioned inter-reader variations in segmentation labels can be
 98 modelled with sophisticated forms of probabilistic CNNs. Such approaches, however, fundamentally
 99 differ from ours in that variable annotations from many experts in the training data are assumed to
 100 be all realistic instances of the true segmentation; we assume, on the other hand, that there is a single,
 101 unknown, true segmentation map of the underlying anatomy, and each individual annotator produces
 102 a noisy approximation to it with variations that reflect their individual characteristics. The latter
 103 assumption may be reasonable in the context of segmentation problems since there exists only one
 104 true boundary of the physical objects captured in an image while multiple hypothesis can arise from
 105 ambiguities in human interpretations.

106 We also note that, in standard classification problems, a plethora of different works have shown
 107 the utility of modelling the labeling process of human annotators in restoring the true label distri-
 108 bution [26, 27, 28]. Such approaches can be categorized into two groups: (1) *two-stage* approach
 109 [29, 30, 31, 32, 33], and (2) *simultaneous* approach. In the first category, the noisy labels are first curated
 110 through a probabilistic model of annotators, and subsequently, a supervised machine-learning model
 111 is trained on the curated labels. The initial attempt [29] was made in the early 1970s, and numerous
 112 advances such as [30, 31, 32, 33] since built upon this work e.g. by estimating sample difficulty and
 113 human biases. In contrast, models in the second category aim to curate labels and learn a supervised
 114 model jointly in an end-to-end fashion [34, 35, 36, 37, 27, 28] so that the two components inform each
 115 other. Although the evidence still remains limited to the simple classification task, these *simultaneous*
 116 approaches have shown promising improvements over the methods in the first category in terms of the
 117 predictive performance of the supervised model and the sample efficiency (i.e., fewer labels are required
 118 per input). However, to date very little attention has been paid to the same problem in more complicated,
 119 structured prediction tasks where the outputs are high dimensional. In this work, we propose the
 120 first *simultaneous* approach to addressing such a problem for image segmentation, while drawing
 121 inspirations from the STAPLE framework [9] which would fall into the *two-stage* approach category.

122 3 Method

123 3.1 Problem Set-up

124 In this work, we consider the problem of learning a supervised segmentation model from noisy
 125 labels acquired from multiple human annotators. Specifically, we consider a scenario where set of
 126 images $\{\mathbf{x}_n \in \mathbb{R}^{W \times H \times C}\}_{n=1}^N$ (with W, H, C denoting the width, height and channels of the image)
 127 are assigned with noisy segmentation labels $\{\tilde{\mathbf{y}}_n^{(r)} \in \mathcal{Y}^{W \times H}\}_{n=1, \dots, N}^{r \in S(\mathbf{x}_i)}$ from multiple annotators where
 128 $\tilde{\mathbf{y}}_n^{(r)}$ denotes the label from annotator $r \in \{1, \dots, R\}$ and $S(\mathbf{x}_n)$ denotes the set of all annotators who
 129 labelled image \mathbf{x}_i and $\mathcal{Y} = [1, 2, \dots, L]$ denotes the set of classes.

130 Here we assume that every image \mathbf{x} annotated by at least one person i.e., $|S(\mathbf{x})| \geq 1$, and no GT labels
 131 $\{\mathbf{y}_n \in \mathcal{Y}^{W \times H}\}_{n=1, \dots, N}$ are available. The problem of interest here is to *learn the unobserved true*
 132 *segmentation distribution* $p(\mathbf{y} | \mathbf{x})$ from such noisy labelled dataset $\mathcal{D} = \{\mathbf{x}_n, \tilde{\mathbf{y}}_n^{(r)}\}_{n=1, \dots, N}^{r \in S(\mathbf{x}_n)}$ i.e., the
 133 combination of images, noisy annotations and experts’ identities for labels (which label was obtained
 134 from whom).

135 We also emphasise that *the goal at inference time is to segment a given unlabelled test image but not*
 136 *to fuse multiple available labels as is typically done in multi-atlas segmentation approaches* [15].

137 3.2 Probabilistic Model and Proposed Architecture

138 Here we describe the probabilistic model of the observed noisy labels from multiple annotators. We
 139 make two key assumptions: (1) annotators are statistically independent, (2) annotations over different
 140 pixels are independent given the input image. Under these assumptions, the probability of observing
 141 noisy labels $\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})}$ on \mathbf{x} factorises as:

$$p(\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})} | \mathbf{x}) = \prod_{r \in S(\mathbf{x})} p(\tilde{\mathbf{y}}^{(r)} | \mathbf{x}) = \prod_{r \in S(\mathbf{x})} \prod_{\substack{w \in \{1, \dots, W\} \\ h \in \{1, \dots, H\}}} p(\tilde{y}_{wh}^{(r)} | \mathbf{x}) \quad (1)$$

142 where $\tilde{y}_{wh}^{(r)} \in [1, \dots, L]$ denotes the $(w, h)^{\text{th}}$ elements of $\tilde{\mathbf{y}}^{(r)} \in \mathcal{Y}^{W \times H}$. Now we rewrite the probability
 143 of observing each noisy label on each pixel (w, h) as:

$$p(\tilde{y}_{wh}^{(r)} | \mathbf{x}) = \sum_{y_{wh}=1}^L p(\tilde{y}_{wh}^{(r)} | y_{wh}, \mathbf{x}) \cdot p(y_{wh} | \mathbf{x}) \quad (2)$$

144 where $p(y_{wh} | \mathbf{x})$ denotes the GT label distribution over the $(w, h)^{\text{th}}$ pixel in the image \mathbf{x} , and
 145 $p(\tilde{y}_{wh}^{(r)} | y_{wh}, \mathbf{x})$ describes the noisy labelling process by which annotator r corrupts the true seg-
 146 mentation label. In particular, we refer to the $L \times L$ matrix whose each $(i, j)^{\text{th}}$ element is defined by the
 147 second term $\mathbf{a}^{(r)}(\mathbf{x}, w, h)_{ij} := p(\tilde{y}_{wh}^{(r)} = i | y_{wh} = j, \mathbf{x})$ as the CM of annotator r at pixel (w, h) in image \mathbf{x} .

148 We introduce a CNN-based architecture which models the different constituents in the above joint
 149 probability distribution $p(\{\tilde{\mathbf{y}}^{(r)}\}_{r \in S(\mathbf{x})} | \mathbf{x})$ as illustrated in Fig. 1. The model consists of two
 150 components: (1) *Segmentation Network*, parametrised by θ , which estimates the GT segmentation
 151 probability map, $\hat{\mathbf{p}}_{\theta}(\mathbf{x}) \in \mathbb{R}^{W \times H \times L}$ whose each $(w, h, i)^{\text{th}}$ element approximates $p(y_{wh} = i | \mathbf{x})$; (2)
 152 *Annotator Network*, parametrised by ϕ , that generate estimates of the pixel-wise CMs of respective
 153 annotators as a function of the input image, $\{\hat{\mathbf{A}}_{\phi}^{(r)}(\mathbf{x}) \in [0, 1]^{W \times H \times L \times L}\}_{r=1}^R$ whose each $(w, h, i, j)^{\text{th}}$
 154 element approximates $p(\tilde{y}_{wh}^{(r)} = i | y_{wh} = j, \mathbf{x})$. Each product $\hat{\mathbf{p}}_{\phi}^{(r)}(\mathbf{x}) := \hat{\mathbf{A}}_{\phi}^{(r)}(\mathbf{x}) \cdot \hat{\mathbf{p}}_{\theta}(\mathbf{x})$ represents the
 155 estimated segmentation probability map of the corresponding annotator. Note that here “ \cdot ” denotes
 156 the element-wise matrix multiplications in the spatial dimensions W, H . At inference time, we use
 157 the output of the segmentation network $\hat{\mathbf{p}}_{\theta}(\mathbf{x})$ to segment test images.

158 We note that each spatial CM $\hat{\mathbf{A}}_{\phi}^{(r)}(\mathbf{x})$ contains WHL^2 variables, and calculating the corresponding
 159 annotator’s prediction $\hat{\mathbf{p}}_{\phi}^{(r)}(\mathbf{x})$ requires $WH(2L-1)L$ floating-point operations, potentially incurring
 160 a large time/space cost when the number of classes is large. Although not the focus of this work (as we
 161 are concerned with medical imaging applications for which the number of classes are mostly limited
 162 to less than 10), we also consider a low-rank approximation (rank=1) scheme to alleviate this issue
 163 wherever appropriate. More details are provided in the supplementary.

164 3.3 Learning Spatial Confusion Matrices and True Segmentation

165 Next, we describe how we jointly optimise the parameters of segmentation network, θ and the param-
 166 eters of annotator network, ϕ . In short, we minimise the negative log-likelihood of the probabilistic model
 167 plus a regularisation term via stochastic gradient descent. A detailed description is provided below.

168 Given training input $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$ and noisy labels $\tilde{\mathbf{Y}}^{(r)} = \{\tilde{\mathbf{y}}_n^{(r)} : r \in \mathcal{S}(\mathbf{x}_n)\}_{n=1}^N$ for
 169 $r = 1, \dots, R$, we optimize the parameters $\{\theta, \phi\}$ by minimizing the negative log-likelihood
 170 (NLL), $-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X})$. From eqs. (1) and (2), this optimization objective equates to the sum
 171 of cross-entropy losses between the observed noisy segmentations and the estimated annotator label
 172 distributions:

$$-\log p(\tilde{\mathbf{Y}}^{(1)}, \dots, \tilde{\mathbf{Y}}^{(R)} | \mathbf{X}) = \sum_{n=1}^N \sum_{r=1}^R \mathbb{1}(\tilde{\mathbf{y}}_n^{(r)} \in \mathcal{S}(\mathbf{x}_n)) \cdot \text{CE}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}) \cdot \hat{\mathbf{p}}_\theta(\mathbf{x}_n), \tilde{\mathbf{y}}_n^{(r)}) \quad (3)$$

173 Minimizing the above encourages each annotator-specific prediction $\hat{\mathbf{p}}^{(r)}(\mathbf{x}) := \hat{\mathbf{A}}_\phi^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$ to be
 174 as close as possible to the true noisy label distribution of the annotator $\mathbf{p}^{(r)}(\mathbf{x})$. However, this loss
 175 function alone is not capable of separating the annotation noise from the true label distribution; there
 176 are many combinations of pairs $\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})$ and segmentation model $\hat{\mathbf{p}}_\theta(\mathbf{x})$ such that $\hat{\mathbf{p}}^{(r)}(\mathbf{x})$ perfectly
 177 matches the true annotator’s distribution $\mathbf{p}^{(r)}(\mathbf{x})$ for any input \mathbf{x} (e.g., permutation of rows in the
 178 CMs). To combat this problem, inspired by Tanno *et al.*[28], which addressed an analogous issue
 179 for the classification task, we add the trace of the estimated CMs to the loss function in Eq. (3) as a
 180 regularisation term (see Sec 3.4). We thus optimize the combined loss:

$$\sum_{n=1}^N \sum_{r=1}^R \mathbb{1}(\tilde{\mathbf{y}}_n^{(r)} \in \mathcal{S}(\mathbf{x}_i)) \cdot \left[\text{CE}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}) \cdot \hat{\mathbf{p}}_\theta(\mathbf{x}_n), \tilde{\mathbf{y}}_n^{(r)}) + \lambda \cdot \text{tr}(\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x}_n)) \right] \quad (4)$$

181 where $\mathcal{S}(\mathbf{x})$ denotes the set of all labels available for image \mathbf{x} , and $\text{tr}(\mathbf{A})$ denotes the trace of matrix
 182 \mathbf{A} . The mean trace represents the average probability that a randomly selected annotator provides an
 183 accurate label. Intuitively, minimising the trace encourages the estimated annotators to be maximally
 184 unreliable while minimising the cross entropy ensures fidelity with observed noisy annotators. We
 185 minimise this combined loss via stochastic gradient descent to learn both $\{\theta, \phi\}$.

186 3.4 Justification for the Trace Norm

187 Here we provide a further justification for using the trace regularisation. Tanno *et al.*[28] showed that if
 188 the average CM of annotators is *diagonally dominant*, and the cross-entropy term in the loss function is
 189 zero, minimising the trace of the estimated CMs uniquely recovers the true CMs. However, their results
 190 concern properties of the average CMs of both the annotators and the classifier over the data population,
 191 rather than individual data samples. We show a similar but slightly weaker result in the sample-specific
 192 regime, which is more relevant as we estimate CMs of respective annotators on every input image.

193 First, let us set up the notations. For brevity, for a given input image $\mathbf{x} \in \mathbb{R}^{W \times H \times C}$, we denote the
 194 estimated CM of annotator r at (i, j) th pixel by $\hat{\mathbf{A}}^{(r)} := [\mathbf{A}^{(r)}(\mathbf{x})_{ij}] \in [0, 1]^{L \times L}$. We also define the
 195 mean CM $\mathbf{A}^* := \sum_{r=1}^R \pi_r \hat{\mathbf{A}}^{(r)}$ and its estimate $\hat{\mathbf{A}}^* := \sum_{r=1}^R \pi_r \hat{\mathbf{A}}^{(r)}$ where $\pi_r \in [0, 1]$ is the probability
 196 that the annotator r labels image \mathbf{x} . Lastly, as we stated earlier, we assume there is a single GT
 197 segmentation label per image — thus the true L -dimensional probability vector at pixel (i, j) takes
 198 the form of a one-hot vector i.e., $\mathbf{p}(\mathbf{x}) = \mathbf{e}_k$ for, say, class $k \in [1, \dots, L]$. Then, the followings result
 199 motivates the use of the trace regularisation:

200 **Theorem 1.** *If the annotator’s segmentation probabilities are perfectly modelled by the model*
 201 *for the given image i.e., $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}) = \mathbf{A}^{(r)} \mathbf{p}(\mathbf{x}) \forall r = 1, \dots, R$, and the average true CM \mathbf{A}^* at a*
 202 *given pixel and its estimate $\hat{\mathbf{A}}^*$ are diagonally dominant ($a_{ii}^* > a_{ij}^*$, $\hat{a}_{ii}^* > \hat{a}_{ij}^*$ for all $i \neq j$), then*
 203 *$\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(R)} = \text{argmin}_{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}} [\text{tr}(\hat{\mathbf{A}}^*)]$ and such solutions are **unique** up to the k^{th} column where*
 204 *k is the correct pixel class.*

205 The corresponding proof is provided in the supplementary material. The above result shows that if
 206 each estimated annotator’s distribution $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x})$ is very close to the true noisy distribution $\mathbf{p}^{(r)}(\mathbf{x})$

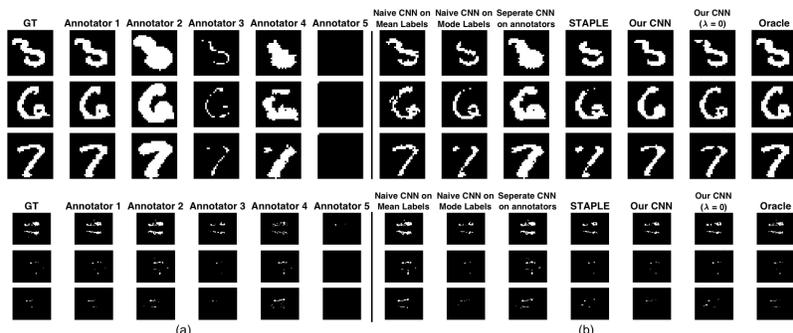


Figure 2: CMs of 5 simulated annotators on MNIST dataset (Best viewed in colour: white is the true positive, green indicates the false negative, red is the false positive and black is the true negative).

207 (which is encouraged by minimizing the cross-entropy loss), and for a given pixel, the average CM
 208 has diagonal entries larger than any other entries in each row ¹, then minimizing its trace will drive
 209 the estimates of the k^{th} (‘correct class’) columns in the respective annotator’s CMs to match the true
 210 values. Although this result is weaker than what was shown in [28] for the population setting rather
 211 than the individual samples, the single-ground-truth assumption means that the remaining values of
 212 the CMs are uniformly equal to $1/L$, and thus it suffices to recover the column of the correct class.

213 To encourage $\{\hat{A}^{(1)}, \dots, \hat{A}^{(R)}\}$ to be also diagonally dominant, we initialize them with identity matrices
 214 by training the *annotation network* to maximise the trace for sufficient iterations as a warm-up period. In-
 215 tuitively, the combination of the trace term and cross-entropy separates the true distribution from the an-
 216 notation noise by finding the maximal amount of confusion which explains the noisy observations well.

217 4 Experiments

218 We evaluate our method on a variety of datasets including both synthetic and real-world scenarios: 1)
 219 for MNIST segmentation and ISBI2015 MS lesion segmentation challenge dataset [38], we apply
 220 morphological operations to generate synthetic noisy labels in binary segmentation tasks; 2) for BraTS
 221 2019 dataset [4], we apply similar simulation to create noisy labels in a multi-class segmentation task;
 222 3) we also consider the LIDC-IDRI dataset which contains multiple annotations per input acquired
 223 from different clinical experts as the evaluation in practice. Details of noisy label simulation can be
 224 found in Appendix A.1.

225 Our experiments are based on the assumption that no ground-truth (GT) label is not known a priori,
 226 hence, we compare our method against multiple label fusion methods. IN particular, we consider four
 227 label fusion baselines: a) mean of all of the noisy labels; b) mode labels by taking the ‘majority vote’;
 228 c) label fusion via the original STAPLE method [9]; d) Spatial STAPLE, a more recent extension of c)
 229 that accounts for spatial variations in CMs. After curating the noisy annotations via above methods, we
 230 train the segmentation network and report the results. For c) and d), we used the toolkit². In addition,
 231 we also include a recent method called Probabilistic U-net as another baseline, which has been shown
 232 to capture inter-reader variations accurately. The details are presented in Appendix A.2.

233 For evaluation metrics, we use: 1) root-MSE between estimated CMs and real CMs; 2) Dice coefficient
 234 (DICE) between estimated segmentation and true segmentation; 3) The generalized energy distance
 235 proposed in [24] to measure the quality of the estimated annotator’s labels.

236 4.1 MNIST and MS lesion segmentation datasets

237 MNIST dataset consists of 60,000 training and 10,000 testing examples, all of which are 28×28
 238 grayscale images of digits from 0 to 9, and we derive the segmentation labels by thresholding the
 239 intensity values at 0.5. The MS dataset is publicly available and comprises 21 3D scans from 5 subjects.
 240 All scans are split into 10 for training and 11 for testing. We hold out 20% of training images as a
 241 validation set for both datasets. On both datasets, our proposed model achieves a higher dice similarity
 242 coefficient than STAPLE on the dense label case and, even more prominently, on the single label
 243 (i.e., 1 label per image) case (shown in Tables. 1&2 and Fig. 2). In addition, our model outperforms

¹For the standard ‘majority vote’ or the mean label to capture the correct true labels, one requires each diagonal element in the average CM to be larger than the sum of the remaining elements in the same row, which is a more strict condition.

²<https://www.nitrc.org/projects/masi-fusion/>

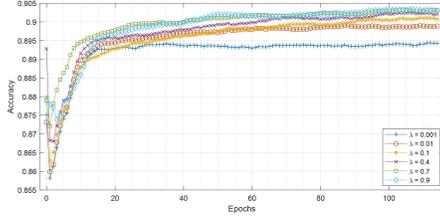


Figure 3: Curves of validation accuracy during training of our model for a range of hyperparameters. For our method, the scaling of trace regularizer is varied in [0.001, 0.01, 0.1, 0.4, 0.7, 0.9].)

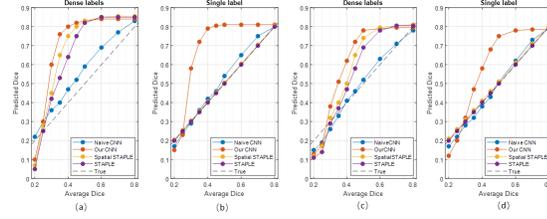


Figure 4: Segmentation accuracy of different models on MNIST (a, b) and MS (c, d) dataset for a range of annotation noise (measured in averaged Dice with respect to GT).

244 STAPLE without or with trace norm, in terms of CM estimation, specifically, we could achieve an
 245 increase at 6.3%. Additionally, we include the performance on different regularisation coefficient,
 246 which is presented in Fig. 3. Fig. 4 compares the segmentation accuracy on MNIST and MS lesion
 247 for a range of average dice where labels are generated by a group of 5 simulated annotators. Fig. 5
 248 illustrates our model can capture the patterns of mistakes for each annotator.

Models	MNIST DICE (%) (testing)	MNIST CM estimation (validation)	MSLesion DICE (%) (testing)	MSLesion CM estimation (validation)
Naive CNN on mean labels	38.36 ± 0.41	n/a	46.55 ± 0.53	n/a
Naive CNN on mode labels	62.89 ± 0.63	n/a	47.82 ± 0.76	n/a
Probabilistic U-net [24]	65.12 ± 0.83	n/a	46.15 ± 0.59	n/a
Separate CNNs on annotators	70.44 ± 0.65	n/a	46.84 ± 1.24	n/a
STAPLE [9]	78.03 ± 0.29	0.1241 ± 0.0011	55.05 ± 0.53	0.1502 ± 0.0026
Spatial STAPLE [14]	78.96 ± 0.22	0.1195 ± 0.0013	58.37 ± 0.47	0.1483 ± 0.0031
Ours without Trace	79.63 ± 0.53	0.1125 ± 0.0037	65.77 ± 0.62	0.1342 ± 0.0053
Ours	82.92 ± 0.19	0.0893 ± 0.0009	67.55 ± 0.31	0.0811 ± 0.0024
Oracle (Ours but with known CMs)	83.29 ± 0.11	0.0238 ± 0.0005	78.86 ± 0.14	0.0415 ± 0.0017

Table 1: Comparison of segmentation accuracy and error of CM estimation for different methods with dense labels (mean ± standard deviation).

Models	MNIST DICE (%) (testing)	MNIST CM estimation (validation)	MSLesion DICE (%) (testing)	MSLesion CM estimation (validation)
Naive CNN on mean & mode labels	32.79 ± 1.13	n/a	27.41 ± 1.45	n/a
STAPLE [9]	54.07 ± 0.68	0.2617 ± 0.0064	35.74 ± 0.84	0.2833 ± 0.0081
Spatial STAPLE [14]	56.73 ± 0.53	0.2384 ± 0.0061	38.21 ± 0.71	0.2591 ± 0.0074
Ours without Trace	74.48 ± 0.37	0.1538 ± 0.0029	54.76 ± 0.66	0.1745 ± 0.0044
Ours	76.48 ± 0.25	0.1329 ± 0.0012	56.43 ± 0.47	0.1542 ± 0.0023

Table 2: Comparison of segmentation accuracy and error of CM estimation for different methods with one label per image (mean ± standard deviation).

Generalised Energy Distance (Dice)	MNIST	MS	BraTS	LIDC-IDRI
Probabilistic U-net [24]	1.46 ± 0.04	1.91 ± 0.03	3.23 ± 0.07	1.97 ± 0.03
Ours	1.24 ± 0.02	1.67 ± 0.03	3.14 ± 0.05	1.87 ± 0.04

Table 3: Comparison of Generalised Energy Distance on different datasets (mean ± standard deviation). The distance metric used here is Dice.

249 4.2 BraTS Dataset and LIDC-IDRI Dataset

250 We also evaluate our model on a multi-class segmentation task, using all of the 259 high grade glioma
 251 (HGG) cases in training data from 2019 multi-modal Brain Tumour Segmentation Challenge (BraTS).
 252 We extract each slice as 2D images and split them at case-wise to have, 1600 images for training, 300
 253 for validation and 500 for testing. Pre-processing includes: concatenation of all of available modalities;
 254 centre cropping to 192 x 192; normalisation for each case at each modality. To create synthetic
 255 noisy labels in multi-class scenario, we first choose a target class and then apply morphological
 256 operations on the provided GT mask to create 4 synthetic noisy labels at different patterns, namely,
 257 over-segmentation, under-segmentation, wrong segmentation and good segmentation. Details of noisy
 258 label simulation are in Appendix A.3.

259 The LIDC-IDRI dataset contains 1018 lung CT scans from 1010 lung patients with manual lesion
 260 segmentations from four experts. For each scan, 4 radiologists provided annotation masks for lesions
 261 that they independently detected and considered to be abnormal. For our experiments, we use the same
 262 method in [24] to pre-process all scans. We split the dataset at case-wise into a training (722 patients),

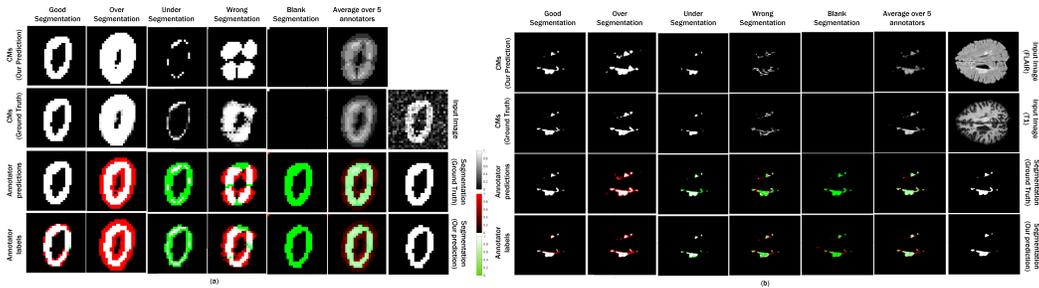


Figure 5: Visualisation of estimated true labels and confusion matrices on MNIST/MS datasets (Best viewed in colour: white is the true positive, green is the false negative, red is the false positive and black is the true negative).

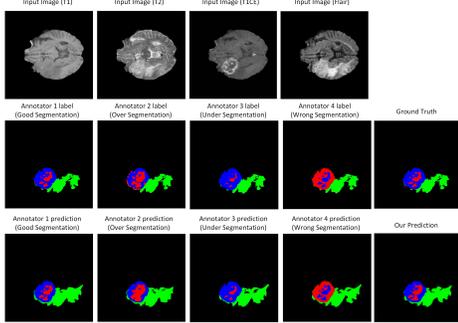


Figure 6: The final segmentation of our model on BraTS and each annotator network predictions visualization. (Best viewed in colour: the target label is red.)

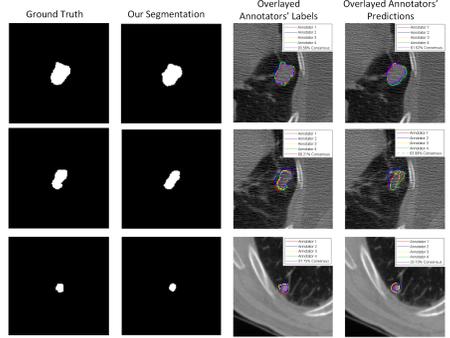


Figure 7: Segmentation results on LIDC-IDRI dataset and the visualization of each annotator contours and the consensus.

263 validation (144 patients) and testing (144 patients). We then resampled the CT scans to $1mm \times 1mm$
 264 in-plane resolution. We also centre cropped 2D images (180×180 pixels) around lesion positions, in
 265 order to focus on the annotated lesions. The lesion positions are those where at least one of the experts
 266 segmented a lesion. We hold 5000 images in the training set, 1000 images in the validation set and
 267 1000 images in the test set.

268 On both BraTS and LIDC-IDRI dataset, our proposed model achieves a higher dice similarity coefficient
 269 than STAPLE and Spatial STAPLE on both of the dense labels and single label scenarios (shown in Ta-
 270 ble. 4 and Table. 5 in Appendix A.3). In addition, our model (with trace) outperforms STAPLE in terms
 271 of CM estimation by a large margin at 14.4% on BraTS. In Fig. 6, we visualized the segmentation results
 272 on BraTS and the corresponding annotators' predictions. Fig. 7 presents three examples of the segmen-
 273 tation results and the corresponding four annotator contours, as well as the consensus. As shown in both
 274 figures, our model successfully predicts the both the segmentation of lesions and the variations of each
 275 annotator in different cases. Additionally, as shown in Table.3, our model consistently outperforms Prob-
 276 abilistic U-Net on generalized energy distance across the four test different datasets, which indicates that
 277 our method is better at capturing the inter-annotator variability than the baseline Probabilistic U-Net.

278 5 Conclusion

279 We introduced the first learning method based on CNNs for simultaneously recovering the label noise of
 280 multiple annotators and the GT label distribution for supervised segmentation problems. We demon-
 281 strated this method on real-world datasets with synthetic annotations and real-world annotations. Our
 282 method is capable of estimating individual annotators and thereby improving robustness against label
 283 noise. Experiments have shown our model achieves considerable improvement over the traditional label
 284 fusion approaches including averaging, the majority vote and the widely used STAPLE framework and
 285 spatially varying versions, in terms of both segmentation accuracy and the quality of CM estimation.

286 In the future, we plan to accommodate meta-information of annotators (e.g., number of years of
 287 experience), and non-image data (e.g., genetics) that may influence the pattern of the underlying
 288 segmentation label such as lesion appearance, in our framework. We are also interested in assessing
 289 the downstream utility of our approach in active data collection schemes where the segmentation
 290 model $\hat{\mathbf{p}}_{\theta}(\mathbf{x})$ is used to select which samples to annotate (“active learning”), and the annotator models
 291 $\{\hat{\mathbf{A}}_{\phi}^{(r)}(\mathbf{x})\}_{r=1}^R$ are used to decide which experts to label them (“active labelling”).

292 **Boarder Impact Statement**

293 *Image segmentation* has been one of the main challenges in modern medical image analysis, and
294 describes the process of assigning each pixel or voxel in images with biologically meaningful discrete
295 labels, such as anatomical structures and tissue types (e.g. pathology and healthy tissues). The task
296 is required in many clinical and research applications, including surgical planning [39, 40], and the
297 study of disease progression, aging or healthy development [41, 42, 43]. However, there are often
298 cases in practice where the correct delineation of structures is challenging; this is also reflected in
299 the well-known presence of high inter- and intra-reader variability in segmentation labels obtained
300 from trained experts [9, 23, 5].

301 Although expert manual annotations of lesions is feasible in practice, this task is time consuming.
302 It usually takes 1.5 to 2 hours to label a MS patient with average 3 visit scans. Meanwhile, the
303 long-established gold standard for segmentation of medical images has been manual voxel-by-voxel
304 labeling by an expert anatomist. Unfortunately, this process is fraught with both interand intra-rater
305 variability (e.g., on the order of approximately 10% by volume [44, 45]). Thus, developing an automatic
306 segmentation technique to fix the variability among inter- and intra-readers could be meaningful not
307 only in terms of the accuracy in delineating MS lesions but also in the related reductions in time and
308 economic costs derived from manual lesion labeling. The lack of consistency in labelling is also
309 common to see in other medical imaging applications, e.g., in lung abnormalities segmentation from
310 CT images. A lesion might be clearly visible by one annotator, but the information about whether it
311 is cancer tissue or not might not be clear to others. However, a potential point of criticism could be that
312 our work in the current form has only been demonstrated on medical images. We would like to convince
313 AC/PCs that the medical imaging domain alone offers a considerably broad range of opportunities for
314 impact; e.g., diagnosis/prognosis in radiology, surgical planning and study of disease progression and
315 treatment, etc. In addition, the annotator information could be potentially utilised for the purpose of
316 education. Another potential opportunity is to integrate such information into the data/label acquisition
317 scheme in order to train reliable segmentation algorithms in a data-efficient manner.

318 **References**

- 319 [1] Elizabeth Lazarus, Martha B Mainiero, Barbara Schepps, Susan L Koelliker, and Linda S
320 Livingston. Bi-rads lexicon for us and mammography: interobserver variability and positive
321 predictive value. *Radiology*, 239(2):385–391, 2006.
- 322 [2] Takeyuki Watadani, Fumikazu Sakai, Takeshi Johkoh, Satoshi Noma, Masanori Akira, Kiminori
323 Fujimoto, Alexander A Bankier, Kyung Soo Lee, Nestor L Müller, Jae-Woo Song, et al.
324 Interobserver variability in the ct assessment of honeycombing in the lungs. *Radiology*,
325 266(3):936–944, 2013.
- 326 [3] Andrew B Rosenkrantz, Ruth P Lim, Mershad Haghghi, Molly B Somberg, James S Babb, and
327 Samir S Taneja. Comparison of interreader reproducibility of the prostate imaging reporting and
328 data system and likert scales for evaluation of multiparametric prostate mri. *American Journal*
329 *of Roentgenology*, 201(4):W612–W618, 2013.
- 330 [4] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani,
331 Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The
332 multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical*
333 *imaging*, 34(10):1993–2024, 2014.
- 334 [5] Leo Joskowicz, D Cohen, N Caplan, and J Sosna. Inter-observer variability of manual contour
335 delineation of structures in ct. *European radiology*, 29(3):1391–1399, 2019.
- 336 [6] Huahong Zhang, Alessandra M Valcarcel, Rohit Bakshi, Renxin Chu, Francesca Bagnato,
337 Russell T Shinohara, Kilian Hett, and Ipek Oguz. Multiple sclerosis lesion segmentation with
338 tiramisu and 2.5 d stacked slices. In *International Conference on Medical Image Computing*
339 *and Computer-Assisted Intervention*, pages 338–346. Springer, 2019.
- 340 [7] Eytan Kats, Jacob Goldberger, and Hayit Greenspan. A soft staple algorithm combined
341 with anatomical knowledge. In *International Conference on Medical Image Computing and*
342 *Computer-Assisted Intervention*, pages 510–517. Springer, 2019.

- 343 [8] Hugh Harvey and Ben Glocker. A standardised approach for preparing imaging data for machine
344 learning tasks in radiology. In *Artificial Intelligence in Medical Imaging*, pages 61–72. Springer,
345 2019.
- 346 [9] Simon K Warfield, Kelly H Zou, and William M Wells. Simultaneous truth and performance level
347 estimation (staple): an algorithm for the validation of image segmentation. *IEEE transactions*
348 *on medical imaging*, 23(7):903–921, 2004.
- 349 [10] Stefan Winzeck, Arsany Hakim, Richard McKinley, José AADSR Pinto, Victor Alves, Carlos
350 Silva, Maxim Pisov, Egor Krivov, Mikhail Belyaev, Miguel Monteiro, et al. Isles 2016 and
351 2017-benchmarking ischemic stroke lesion outcome prediction based on multispectral mri.
352 *Frontiers in neurology*, 9:679, 2018.
- 353 [11] Olivier Commowick, Audrey Istace, Michael Kain, Baptiste Laurent, Florent Leray, Mathieu
354 Simon, Sorina Camarasu Pop, Pascal Girard, Roxana Ameli, Jean-Christophe Ferré, et al.
355 Objective evaluation of multiple sclerosis lesion segmentation using a data management and
356 processing infrastructure. *Scientific reports*, 8(1):1–17, 2018.
- 357 [12] Gleason 2019 challenge. <https://gleason2019.grand-challenge.org/Home/>. Ac-
358 cessed: 2020-02-30.
- 359 [13] Andrew J Asman and Bennett A Landman. Robust statistical label fusion through consensus
360 level, labeler accuracy, and truth estimation (collate). *IEEE transactions on medical imaging*,
361 30(10):1779–1794, 2011.
- 362 [14] Andrew J Asman and Bennett A Landman. Formulating spatially varying performance in the
363 statistical fusion framework. *IEEE transactions on medical imaging*, 31(6):1326–1336, 2012.
- 364 [15] Juan Eugenio Iglesias, Mert Rory Sabuncu, and Koen Van Leemput. A unified framework for cross-
365 modality multi-atlas segmentation of brain mri. *Medical image analysis*, 17(8):1181–1191, 2013.
- 366 [16] M Jorge Cardoso, Kelvin Leung, Marc Modat, Shiva Keihaninejad, David Cash, Josephine
367 Barnes, Nick C Fox, Sebastien Ourselin, Alzheimer’s Disease Neuroimaging Initiative, et al.
368 Steps: Similarity and truth estimation for propagated segmentations and its application to
369 hippocampal segmentation and brain parcelation. *Medical image analysis*, 17(6):671–684, 2013.
- 370 [17] Andrew J Asman and Bennett A Landman. Non-local statistical label fusion for multi-atlas
371 segmentation. *Medical image analysis*, 17(2):194–208, 2013.
- 372 [18] Alireza Akhondi-Asl, Lennox Hoyte, Mark E Lockhart, and Simon K Warfield. A logarithmic
373 opinion pool based staple algorithm for the fusion of segmentations with associated reliability
374 weights. *IEEE transactions on medical imaging*, 33(10):1997–2009, 2014.
- 375 [19] Daniel C. Castro, Jeremy Tan, Bernhard Kainz, Ender Konukoglu, and Ben Glocker. Morpho-
376 MNIST: Quantitative assessment and diagnostics for representation learning. *Journal of Machine*
377 *Learning Research*, 20, 2019.
- 378 [20] Martin Styner, Joohwi Lee, Brian Chin, M Chin, Olivier Commowick, H Tran, S Markovic-Plese,
379 V Jewells, and S Warfield. 3d segmentation in the clinic: A grand challenge ii: Ms lesion
380 segmentation. *Midas Journal*, 2008:1–6, 2008.
- 381 [21] Samuel G Armato III, Geoffrey McLennan, Luc Bidaut, Michael F McNitt-Gray, Charles R Meyer,
382 Anthony P Reeves, Binsheng Zhao, Denise R Aberle, Claudia I Henschke, Eric A Hoffman,
383 et al. The lung image database consortium (lidc) and image database resource initiative (idri): a
384 completed reference database of lung nodules on ct scans. *Medical physics*, 38(2):915–931, 2011.
- 385 [22] Neil I Weisenfeld and Simon K Warfield. Learning likelihoods for labeling (l3): a general
386 multi-classifier segmentation algorithm. In *International Conference on Medical Image*
387 *Computing and Computer-Assisted Intervention*, pages 322–329. Springer, 2011.
- 388 [23] Leo Joskowicz, D Cohen, N Caplan, and Jacob Sosna. Automatic segmentation variability
389 estimation with segmentation priors. *Medical image analysis*, 50:54–64, 2018.

- 390 [24] Simon Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R Ledsam,
391 Klaus Maier-Hein, SM Ali Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A
392 probabilistic u-net for segmentation of ambiguous images. In *Advances in Neural Information*
393 *Processing Systems*, pages 6965–6975, 2018.
- 394 [25] Christian F Baumgartner, Kerem C Tezcan, Krishna Chaitanya, Andreas M Hötter, Urs J
395 Muehlematter, Khoschy Schawkat, Anton S Becker, Olivio Donati, and Ender Konukoglu.
396 Phiseg: Capturing uncertainty in medical image segmentation. In *International Conference on*
397 *Medical Image Computing and Computer-Assisted Intervention*, pages 119–127. Springer, 2019.
- 398 [26] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Gerardo Hermosillo Valadez, Charles Florin,
399 Luca Bogoni, and Linda Moy. Learning from crowds. *Journal of Machine Learning Research*,
400 11(Apr):1297–1322, 2010.
- 401 [27] Ashish Khetan, Zachary C Lipton, and Anima Anandkumar. Learning from noisy singly-labeled
402 data. *arXiv preprint arXiv:1712.04577*, 2017.
- 403 [28] Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C Alexander, and Nathan
404 Silberman. Learning from noisy labels by regularized estimation of annotator confusion. *arXiv*
405 *preprint arXiv:1902.03680*, 2019.
- 406 [29] Alexander Philip Dawid and Allan M Skene. Maximum likelihood estimation of observer
407 error-rates using the em algorithm. *Applied statistics*, pages 20–28, 1979.
- 408 [30] Padhraic Smyth, Usama M Fayyad, Michael C Burl, Pietro Perona, and Pierre Baldi. Inferring
409 ground truth from subjective labelling of venus images. In *Advances in neural information*
410 *processing systems*, pages 1085–1092, 1995.
- 411 [31] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier R Movellan, and Paul L Ruvolo. Whose
412 vote should count more: Optimal integration of labels from labelers of unknown expertise. In
413 *Advances in neural information processing systems*, pages 2035–2043, 2009.
- 414 [32] Peter Welinder, Steve Branson, Pietro Perona, and Serge J Belongie. The multidimensional wis-
415 dom of crowds. In *Advances in neural information processing systems*, pages 2424–2432, 2010.
- 416 [33] Filipe Rodrigues, Francisco Pereira, and Bernardete Ribeiro. Learning from multiple annotators:
417 distinguishing good from random labelers. *Pattern Recognition Letters*, 34(12):1428–1436, 2013.
- 418 [34] Vikas C Raykar, Shipeng Yu, Linda H Zhao, Anna Jerebko, Charles Florin, Gerardo Hermosillo
419 Valadez, Luca Bogoni, and Linda Moy. Supervised learning from multiple experts: whom to
420 trust when everyone lies a bit. In *Proceedings of the 26th Annual international conference on*
421 *machine learning*, pages 889–896. ACM, 2009.
- 422 [35] Yan Yan, Rómer Rosales, Glenn Fung, Mark Schmidt, Gerardo Hermosillo, Luca Bogoni, Linda
423 Moy, and Jennifer Dy. Modeling annotator expertise: Learning when everybody knows a bit
424 of something. In *AISTATS*, pages 932–939, 2010.
- 425 [36] Steve Branson, Grant Van Horn, and Pietro Perona. Lean crowdsourcing: Combining humans
426 and machines in an online system. In *Proceedings of the IEEE Conference on Computer Vision*
427 *and Pattern Recognition*, pages 7474–7483, 2017.
- 428 [37] Grant Van Horn, Steve Branson, Scott Loarie, Serge Belongie, Cornell Tech, and Pietro Perona.
429 Lean multiclass crowdsourcing. In *Proceedings of the IEEE Conference on Computer Vision*
430 *and Pattern Recognition*, pages 2714–2723, 2018.
- 431 [38] Andrew Jesson and Tal Arbel. Hierarchical mrf and random forest segmentation of ms lesions
432 and healthy tissues in brain mri. *Proceedings of the 2015 Longitudinal Multiple Sclerosis Lesion*
433 *Segmentation Challenge*, pages 1–2, 2015.
- 434 [39] David T Gering, Arya Nabavi, Ron Kikinis, Noby Hata, Lauren J O’Donnell, W Eric L Grimson,
435 Ferenc A Jolesz, Peter M Black, and William M Wells III. An integrated visualization system
436 for surgical planning and guidance using image fusion and an open mr. *Journal of Magnetic*
437 *Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance*
438 *in Medicine*, 13(6):967–975, 2001.

- 439 [40] Gloria P Mazzara, Robert P Velthuizen, James L Pearlman, Harvey M Greenberg, and Henry
440 Wagner. Brain tumor target volume determination for radiation treatment planning through
441 automated mri segmentation. *International Journal of Radiation Oncology* Biology* Physics*,
442 59(1):300–312, 2004.
- 443 [41] Bruce Fischl, David H Salat, Evelina Busa, Marilyn Albert, Megan Dieterich, Christian
444 Haselgrove, Andre Van Der Kouwe, Ron Killiany, David Kennedy, Shuna Klaveness, et al.
445 Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain.
446 *Neuron*, 33(3):341–355, 2002.
- 447 [42] Marcel Prastawa, John H Gilmore, Weili Lin, and Guido Gerig. Automatic segmentation of mr
448 images of the developing newborn brain. *Medical image analysis*, 9(5):457–466, 2005.
- 449 [43] Alex P Zijdenbos, Reza Forghani, Alan C Evans, et al. Automatic" pipeline" analysis of
450 3-d mri data for clinical trials: application to multiple sclerosis. *IEEE Trans. Med. Imaging*,
451 21(10):1280–1291, 2002.
- 452 [44] Edward A Ashton, Chihiro Takahashi, Michel J Berg, Andrew Goodman, Saara Totterman, and
453 Sven Ekholm. Accuracy and reproducibility of manual and semiautomated quantification of ms
454 lesions by mri. *Journal of Magnetic Resonance Imaging: An Official Journal of the International*
455 *Society for Magnetic Resonance in Medicine*, 17(3):300–308, 2003.
- 456 [45] Bonnie N Joe, Melanie B Fukui, Carolyn Cidis Meltzer, Qing-shou Huang, Roger S Day, Phil J
457 Greer, and Michael E Bozik. Brain tumor volume measurement: comparison of manual and
458 semiautomated methods. *Radiology*, 212(3):811–816, 1999.
- 459 [46] Sergi Valverde, Mostafa Salem, Mariano Cabezas, Deborah Pareto, Joan C. Vilanova, Lluís
460 Ramió-Torrentà, Àlex Rovira, Joaquim Salvi, Arnau Oliver, and Xavier Lladó. One-shot
461 domain adaptation in multiple sclerosis lesion segmentation using convolutional neural networks.
462 *NeuroImage: Clinical*, page 101638, 2018.
- 463 [47] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training
464 convolutional networks with noisy labels. *arXiv preprint arXiv:1406.2080*, 2014.

465 **A Additional results**

466 **A.1 Data Simulation**

467 We generate synthetic annotations from an assumed GT to generate efficacy of the approach in
 468 an idealised situation where the GT is known. We simulate a group of 5 annotators of disparate
 469 characteristics by performing morphological transformations (e.g., thinning, thickening, fractures, etc)
 470 on the ground-truth (GT) segmentation labels, using Morpho-MNIST software [19]. In particular, the
 471 first annotator provides faithful segmentation (“good-segmentation”) with approximate GT, the second
 472 tends over-segment (“over-segmentation”), the third tends to under-segment (“under-segmentation”),
 473 the fourth is prone to the combination of small fractures and over-segmentation (“wrong-segmentation”)
 474 and the fifth always annotates everything as the background (“blank-segmentation”). We create
 475 training data by deriving labels from the simulated annotators.

476 **A.2 MNIST and MS Dataset**

477 We examine the ability of our method to learn the CMs of annotators and the true label distribution. We
 478 compared the performance of our method against several baselines and the original STAPLE algorithm
 479 [9] and Spatial STAPLE [14]. The first baseline is the naive CNN trained on the mean labels and the
 480 majority vote labels across the 5 annotators. The second baseline is the separate CNNs trained on 5
 481 annotator labels and evaluate on their mean output. The “oracle” model is the idealistic scenario where
 482 CMs of the annotators are a priori known to the model while “annotators” indicate the average labeling
 483 accuracy of each annotator group. All the baselines and the annotator CNN, the segmentation CNN
 484 in our model are implemented with the NicMSLesions architecture described in [46]. We also evaluate
 485 on the validation set the effects of regularisation coefficient $\lambda \in \{0, 0.001, 0.01, 0.1, 0.4, 0.7, 0.9\}$ of the
 486 trace-norm in Eq. 4 on the accuracy of segmentation and CM estimation. Results are shown in Fig. 3.

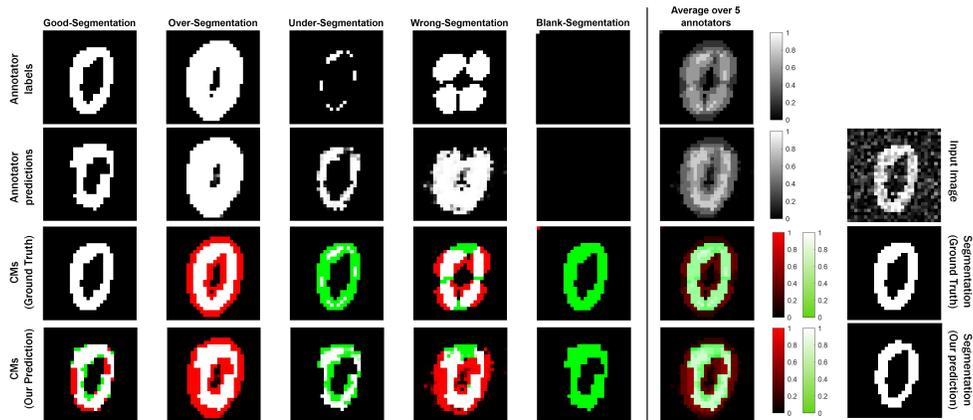


Figure 7: Visualisation of estimated true labels and confusion matrices on MNIST datasets (Best viewed in colour: white is the true positive, green is the false negative, red is the false positive and black is the true negative).

487 **A.3 BraTS and LIDC-IDRI**

488 We also evaluate our model on a multi-class segmentation task, using training data from 2019 Brain
 489 Tumour Segmentation Challenge (BraTS). In training data of BraTS 2019, there are 259 cases with
 490 high grade (HG) and 76 cases with low grade (LG) glioma. For each case, four MRI modalities are
 491 available, FLAIR, T1, T1-contrast and T2. The datasets are pre-processed by the organizers and
 492 co-registered to the same anatomical template, interpolated to the same resolution (1 mm^3) and
 493 skull-stripped. We centre cropped 2D images (192×192 pixels) and hold 1600 2D images for training,
 494 300 images for validation, 500 images for testing, we apply Gaussian normalization on each case of
 495 each modality, to have zero-mean and unit variance. Fig. 6 shows one such tumor case in four different
 496 modality. To create synthetic noisy labels in multi-class scenario, we first choose a target class and
 497 then apply morphological operations on the provided GT mask to create 4 synthetic noisy labels at
 498 different patterns, namely, over-segmentation, under-segmentation, wrong segmentation and good
 499 segmentation. Details of noisy label simulation are in Appendix A.3.

500 The LIDC-IDRI dataset contains 1018 lung CT scans from 1010 lung patients with manual lesion
 501 segmentations from four experts. For each scan, 4 radiologists provided annotation masks for lesions
 502 that they independently detected and considered to be abnormal. For our experiments, we use the same
 503 method in [24] to pre-process all scans. We split the dataset at case-wise into a training (722 patients),
 504 validation (144 patients) and testing (144 patients). We then resampled the CT scans to $1\text{mm} \times 1\text{mm}$
 505 in-plane resolution. We also centre cropped 2D images (180×180 pixels) around lesion positions, in
 506 order to focus on the annotated lesions. The lesion positions are those where at least one of the experts
 507 segmented a lesion. We hold 5000 images in the training set, 1000 images in the validation set and
 508 1000 images in the test set.

509 On both BraTS and LIDC-IDRI dataset, our proposed model achieves a higher dice similarity coefficient
 510 than STAPLE on both of the dense labels and single label scenarios (shown in Table. 4 and Table. 5
 511 in Appendix A.3). In addition, our model (with trace) outperforms STAPLE in terms of CM estimation
 512 by a large margin at 14.4% on BraTS. In Fig. 6, we visualized the segmentation results on BraTS and
 513 the corresponding annotators’ predictions. Fig. 7 presents three examples of the segmentation results
 514 and the corresponding four annotator contours, as well as the consensus. As shown in both figures, our
 515 model successfully predicts the both the segmentation of lesions and the variations of each annotator
 516 in different cases. Additionally, as shown in Table.3, our model consistently outperforms Probabilistic
 517 U-Net on generalized energy distance across the four test different datasets, which indicates that our
 518 method is better at capturing the inter-annotator variability than the baseline Probabilistic U-Net.

Models	BraTS DICE (%) (testing)	BraTS CM estimation (validation)	LIDC-IDRI DICE (%) (testing)	LIDC-IDRI CM estimation (validation)
Naive CNN on mean labels	29.42 ± 0.58	n/a	56.72 ± 0.61	n/a
Naive CNN on mode labels	34.12 ± 0.45	n/a	58.64 ± 0.47	n/a
Probabilistic U-net [24]	40.53 ± 0.75	n/a	61.26 ± 0.69	n/a
STAPLE [9]	46.73 ± 0.17	0.2147 ± 0.0103	69.34 ± 0.58	0.0832 ± 0.0043
Spatial STAPLE [14]	47.31 ± 0.21	0.1871 ± 0.0094	70.92 ± 0.18	0.0746 ± 0.0057
Ours without Trace	49.03 ± 0.34	0.1569 ± 0.0072	71.25 ± 0.12	0.0482 ± 0.0038
Ours	53.47 ± 0.24	0.1185 ± 0.0056	74.12 ± 0.19	0.0451 ± 0.0025
Oracle (Ours but with known CMs)	67.13 ± 0.14	0.0843 ± 0.0029	79.41 ± 0.17	0.0381 ± 0.0021

Table 4: Comparison of segmentation accuracy and error of CM estimation for different methods with dense labels (mean ± standard deviation). (ryu): On LIDC-IDRI, it is rather surprising that our method performs better than SpatialSTAPLE when the GT were created by SpatialSTAPLE! We should mention this clearly in the results section.

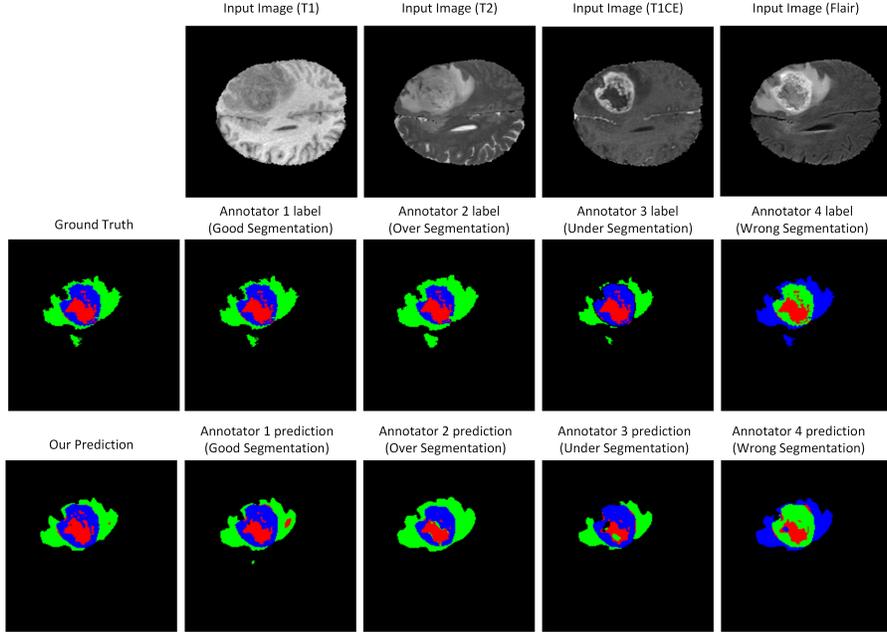


Figure 8: The final segmentation of our model on BraTS and each annotator network predictions visualization. (Best viewed in colour: the target label is red.)

Models	BraTS DICE (%) (testing)	BraTS CM estimation (validation)	LIDC-IDRI DICE (%) (testing)	LIDC-IDRI CM estimation (validation)
Naive CNN on mean & mode labels	36.12 ± 0.93	n/a	48.36 ± 0.79	n/a
STAPLE [9]	38.74 ± 0.85	0.2956 ± 0.1047	57.32 ± 0.87	0.1715 ± 0.0134
Spatial STAPLE [14]	41.59 ± 0.74	0.2543 ± 0.0867	62.35 ± 0.64	0.1419 ± 0.0207
Ours without Trace	43.74 ± 0.49	0.1825 ± 0.0724	66.95 ± 0.51	0.0921 ± 0.0167
Ours	46.21 ± 0.28	0.1576 ± 0.0487	68.12 ± 0.48	0.0587 ± 0.0098

Table 5: Comparison of segmentation accuracy and error of CM estimation for different methods with one label per image (mean \pm standard deviation).

519 A.4 Low-rank approximation

520 In particular, we parametrise the spatial CM $\hat{\mathbf{A}}_{\phi}^{(r)}(\mathbf{x}) = \mathbf{B}_{1,\phi}^{(r)}(\mathbf{x}) \cdot \mathbf{B}_{2,\phi}^{T,(r)}(\mathbf{x})$ where both $\mathbf{B}_{1,\phi}^{(r)}$ and $\mathbf{B}_{2,\phi}^{(r)}$
521 are smaller matrices of size $W \times H \times L \times l$ where $l \ll L$. Two separate rectangular matrices are used
522 since the confusion matrices are not necessarily symmetric. Such low-rank approximation reduces
523 the total number of variables to $2WHLl$ and the FLOPs to $WH(4L(l-0.25)-l)$. **Still need to decide**
524 **whether to include this paragraph depending on the results on DICE.**

Rank	Dice	CM estimation	GPU Memory	No. Parameters & FLOPS
Default	53.47 ± 0.24	0.1185 ± 0.0056	2.68GB	192×192
rank 1	50.56 ± 2.00	-	2.57GB	

Table 6: Segmentation performance of low-rank approximation on BraTS. GPU memory is when batch size 1 is used (mean \pm standard deviation).

525 A.5 Algorithm

526 (Copy from CVPR): Here we provide pseudo-codes of our method (Algorithm 1), generalized EM [34]
527 (Algorithm 2) and model-bootstrapped EM [27] (Algorithm 3) to clarify the differences between differ-
528 ent methods for jointly learning the true label distribution and confusion matrices of annotators in eq. 2
529 in the main text. Given the training set $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N$, each example may not be labelled by

530 all the annotators. In such cases, for ease of notation, we assign pseudo class $\tilde{y}_n^{(r)} = -1$ to fill the missing
 531 labels. The comparison between these three algorithms illustrates the implementational simplicity of
 our method, despite the comparable or superior performance demonstrated on all three datasets.

Algorithm 1 Our method

Inputs: $\mathcal{D} = \{\mathbf{x}_n, \tilde{y}_n^{(1)}, \dots, \tilde{y}_n^{(R)}\}_{n=1}^N$, λ : scale of trace regularizer

Initialize the confusion matrices $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ **to identity matrices**

Initialize the parameters of the base classifier θ

Learn θ **and** $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$ **by performing minibatch SGD on the combined loss:**

$$\theta, \{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R \leftarrow \operatorname{argmin}_{\theta, \{\hat{\mathbf{A}}^{(r)}\}} \left[\sum_{i=1}^N \sum_{r=1}^R \mathbb{1}(\tilde{y}_i^{(r)} \neq -1) \cdot \operatorname{CE}(\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_{\theta}(\mathbf{x}_i), \tilde{y}_i^{(r)}) + \lambda \sum_{r=1}^R \operatorname{tr}(\hat{\mathbf{A}}^{(r)}) \right]$$

Return: $\hat{\mathbf{p}}_{\theta}$ and $\{\hat{\mathbf{A}}^{(r)}\}_{r=1}^R$

532

533 **B Proof of Theorem 2**

534 (Ryu): still need to change the statement of the lemma and the proof.

535 Here we intend to motivate the addition of the trace regularizer in eq. (4). In the last section, we saw
 536 that minimizing cross-entropy loss alone encourages $\hat{\mathbf{A}}^{(r)} \mathbf{P} \rightarrow \mathbf{A}^{(r)}$. Therefore, if we could devise
 537 a regularizer which, when minimized, uniquely ensures the convergence $\hat{\mathbf{A}}^{(r)} \rightarrow \mathbf{A}^{(r)}$, then this
 538 would make \mathbf{P} tend to the identity matrix, implying that the base model fully captures the true label
 539 distribution i.e. $\operatorname{argmax}_k [\mathbf{p}(\mathbf{x})]_k = y \forall \mathbf{x}$. We describe below the trace regularizer is indeed a such
 540 regularizer when both $\hat{\mathbf{A}}^{(r)}$ and $\mathbf{A}^{(r)}$ satisfy some conditions. We first show this result assuming that
 541 there is a single annotator, and then extend to the scenario with multiple annotators.

542 **Lemma 1** (Single Annotator). *Let \mathbf{P} be the CM of the estimated true labels $\hat{\mathbf{p}}_\theta$ and $\hat{\mathbf{A}}$ be the estimated*
 543 *CM of the annotator. If the model matches the noisy label distribution of the annotator i.e. $\hat{\mathbf{A}}\mathbf{P} = \mathbf{A}$,*
 544 *and both $\hat{\mathbf{A}}$ and \mathbf{A} are diagonally dominant ($a_{ii} > a_{ij}$, $\hat{a}_{ii} > \hat{a}_{ij}$) for all $i \neq j$, then $\hat{\mathbf{A}}$ with the minimal*
 545 *trace uniquely coincides with the true \mathbf{A} .*

546 *Proof.* We show that each diagonal element in the true CM \mathbf{A} forms a lower bound to the corresponding
 547 element in its estimation.

$$a_{ii} = \sum_j \hat{a}_{ij} p_{ji} \leq \sum_j \hat{a}_{ii} p_{ji} = \hat{a}_{ii} \left(\sum_j p_{ji} \right) = \hat{a}_{ii} \quad (5)$$

548 for all $i \in \{1, \dots, L\}$. It therefore follows that $\operatorname{tr}(\mathbf{A}) \leq \operatorname{tr}(\hat{\mathbf{A}})$. We now show that the equality $\hat{\mathbf{A}} = \mathbf{A}$ is
 549 uniquely achieved when the trace is the smallest i.e. $\operatorname{tr}(\mathbf{A}) = \operatorname{tr}(\hat{\mathbf{A}}) \Rightarrow \mathbf{A} = \hat{\mathbf{A}}$. From (5), if the trace
 550 of \mathbf{A} and $\hat{\mathbf{A}}$ are the same, we see that their diagonal elements also match i.e. $a_{ii} = \hat{a}_{ii} \forall i \in \{1, \dots, L\}$.
 551 Now, the non-negativity of all elements in CMs \mathbf{P} and $\hat{\mathbf{A}}$, and the equality $a_{ii} = \sum_j \hat{a}_{ij} p_{ji}$ imply that
 552 $p_{ji} = \mathbb{1}[i = j]$ i.e. \mathbf{P} is the identity matrix.

553 First, let us set up the notations. For brevity, for a given input image $\mathbf{x} \in \mathbb{R}^{W \times H \times C}$, we denote the
 554 estimated CM of annotator r at (i, j) th pixel by $\hat{\mathbf{A}}^{(r)} := [\mathbf{A}^{(r)}(\mathbf{x})]_{ij} \in [0, 1]^{L \times L}$. We also define the
 555 mean CM $\mathbf{A}^* := \sum_{r=1}^R \pi_r \hat{\mathbf{A}}^{(r)}$ and its estimate $\hat{\mathbf{A}}^* := \sum_{r=1}^R \pi_r \hat{\mathbf{A}}^{(r)}$ where $\pi_r \in [0, 1]$ is the probability
 556 that the annotator r labels image \mathbf{x} . Lastly, as we stated earlier, we assume there is a single GT
 557 segmentation label per image — thus the true L -dimensional probability vector at pixel (i, j) takes
 558 the form of a one-hot vector i.e., $\mathbf{p}(\mathbf{x}) = \mathbf{e}_k$ for, say, class $k \in [1, \dots, L]$. Then, the followings result
 559 motivates the use of the trace regularisation:

560 **Theorem 2.** *If the annotator’s segmentation probabilities are perfectly modelled by the model*
 561 *for the given image i.e., $\hat{\mathbf{A}}^{(r)} \hat{\mathbf{p}}_\theta(\mathbf{x}) = \mathbf{A}^{(r)} \mathbf{p}(\mathbf{x}) \forall r = 1, \dots, R$, and the average true CM \mathbf{A}^* at a*
 562 *given pixel and its estimate $\hat{\mathbf{A}}^*$ are diagonally dominant ($a_{ii}^* > a_{ij}^*$, $\hat{a}_{ii}^* > \hat{a}_{ij}^*$ for all $i \neq j$), then*
 563 $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(R)} = \operatorname{argmin}_{\hat{\mathbf{A}}^{(1)}, \dots, \hat{\mathbf{A}}^{(R)}} \left[\operatorname{tr}(\hat{\mathbf{A}}^*) \right]$ *and such solutions are **unique** up to the k^{th} column where*
 564 *k is the correct pixel class.*

565 □

566 We note that the above result was also mentioned in [47] in a more general context of label noise
 567 modelling (that neglects annotator information). Here we further augment their proof by showing
 568 the uniqueness of solutions (i.e. $\operatorname{tr}(\mathbf{A}) = \operatorname{tr}(\hat{\mathbf{A}}) \Rightarrow \mathbf{A} = \hat{\mathbf{A}}$). In addition, the trace regularization was
 569 never used in practice in [47] — for implementation reason, the Frobenius norm was used in all their
 570 experiments. We now extend this to the multiple annotator regime.

571 (Ryu): need to adapt this result too to the new setting.

572 *Proof.* As the average CMs \mathbf{A}^* and $\hat{\mathbf{A}}^*$ are diagonally dominant and we have $\mathbf{A}^* = \hat{\mathbf{A}}^* \mathbf{P}$, Lemma 1
 573 yields that $\operatorname{tr}(\mathbf{A}^*) \leq \operatorname{tr}(\hat{\mathbf{A}}^*)$ with equality if and only if $\mathbf{A}^* = \hat{\mathbf{A}}^*$. Therefore, when the trace of the

574 average CM of annotators is minimized i.e. $\text{tr}(\hat{\mathbf{A}}^*) = \text{tr}(\mathbf{A}^*)$, the estimated CM of the true label
575 distribution \mathbf{P} reduces to identity, giving $\hat{\mathbf{A}}^{(r)} = \mathbf{A}^{(r)}$ for all $r \in \{1, \dots, R\}$. \square