

Journal Pre-proof

Graph Temporal Ensembling based Semi-supervised Convolutional Neural Network with Noisy Labels for Histopathology Image Analysis

Xiaoshuang Shi , Hai Su , Fuyong Xing , Yun Liang , Gang Qu , Lin Yang

PII: S1361-8415(19)30160-4
DOI: <https://doi.org/10.1016/j.media.2019.101624>
Reference: MEDIMA 101624



To appear in: *Medical Image Analysis*

Received date: 26 December 2018
Revised date: 22 November 2019
Accepted date: 25 November 2019

Please cite this article as: Xiaoshuang Shi , Hai Su , Fuyong Xing , Yun Liang , Gang Qu , Lin Yang , Graph Temporal Ensembling based Semi-supervised Convolutional Neural Network with Noisy Labels for Histopathology Image Analysis, *Medical Image Analysis* (2019), doi: <https://doi.org/10.1016/j.media.2019.101624>

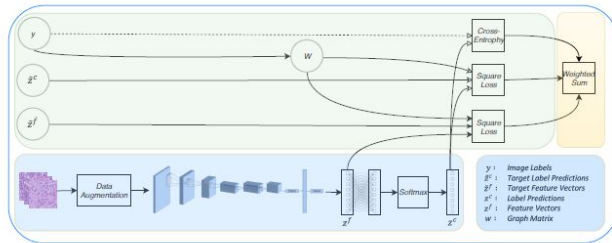
This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Published by Elsevier B.V.

Research highlights

- We utilize the graph to boost the strength of ensemble predictions of training samples.
- We propose a novel loss function to form consensus predicting features and labels.
- We propose a novel self-ensembling based robust semi-supervised deep architecture for histopathology image analysis.

Graphical abstract



Graph Temporal Ensembling based Semi-supervised Convolutional Neural Network with Noisy Labels for Histopathology Image Analysis

Xiaoshuang Shi^a, Hai Su^a, Fuyong Xing^b, Yun Liang^a, Gang Qu^a, Lin Yang^{a*}

^a*J. Crayton Pruitt Family Department of Biomedical Engineering, University of Florida*

^b*Department of Biostatistics and Informatics, University of Colorado Denver*

Abstract

Although convolutional neural networks have achieved tremendous success on histopathology image classification, they usually require large-scale clean annotated data and are sensitive to noisy labels. Unfortunately, labeling large-scale images is laborious, expensive and lowly reliable for pathologists. To address these problems, in this paper, we propose a novel self-ensembling based deep architecture to leverage the semantic information of annotated images and explore the information hidden in unlabeled data, and meanwhile being robust to noisy labels. Specifically, the proposed architecture first creates ensemble targets for feature and label predictions of training samples, by using exponential moving average (EMA) to aggregate feature and label predictions within multiple previous training epochs. Then, the ensemble targets within the same class are mapped into a cluster so that they are further enhanced. Next, a consistency cost is utilized to form consensus

*Corresponding author

Email address: lin.yang@bme.ufl.edu (Lin Yang^a)

predictions under different configurations. Finally, we validate the proposed method with extensive experiments on lung and breast cancer datasets that contain thousands of images. It can achieve 90.5% and 89.5% image classification accuracy using only 20% labeled patients on the two datasets, respectively. This performance is comparable to that of the baseline method with all labeled patients. Experiments also demonstrate its robustness to small percentage of noisy labels.

Keywords: Semi-supervised; Noisy labels; Convolutional neural network; Histopathology image classification

1. Introduction

Because of the advance in high-throughput tissue bank and archiving of digitized histological studies, histopathology images with computer aided diagnosis (CAD) using modern machine learning techniques have attracted considerable attention to facilitate disease grading and classification (Gurcan et al., 2009) (Shen et al., 2017) (Shi et al., 2017) (Xing et al., 2017). Recently, with an ever-increasing amount of images and the development of deep neural networks (Xu et al., 2018) (Xu et al., 2019), especially convolutional neural networks (CNNs), it is promising to bridge the semantic gap between images and diagnostic information (Zhang et al., 2015) (Litjens et al., 2017) (Shi et al., 2018) (Sapkota et al., 2018) (Chen et al., 2019). Most of current deep learning methods require a large amount of clean annotated data to achieve desired performance, and they are usually sensitive to noisy labels, i.e. a very small percentage of noisy labels might severely decrease the model performance (Patrini et al., 2017). Unfortunately, labeling large-scale

16 histopathology images is laborious, expensive and time-consuming. Addi-
17 tionally, the annotation process is conducted with low labeling reliability,
18 resulting in noisy labels due to the subjective assessment of pathologists.
19 Therefore, it is necessary to design effective and efficient deep neural networks
20 for histopathology image analysis, that can be trained on a small amount of
21 labeled data yet large-scale unlabeled data. Meanwhile, the network should
22 be robust to a small percentage of noisy labels.

23 To leverage the semantic information of labeled data and meanwhile ex-
24 plore the information hidden in unlabeled data, numerous semi-supervised
25 deep learning methods have been proposed and applied to various appli-
26 cations, such as image classification and retrieval (Zhang and Peng, 2017),
27 detection (Tang et al., 2018) and segmentation (Bai et al., 2017). Among
28 previous semi-supervised deep classification methods, self-ensembling based
29 methods have achieved state-of-the-art accuracy on multiple benchmark im-
30 age databases. This is because self-ensembling can successfully explore the
31 semantic information hidden in unlabeled data, by creating an ensemble tar-
32 get for each label prediction and forming consensus predictions under differ-
33 ent configurations, such as different epochs, dropout regularizations and in-
34 put augmentations. However, previous self-ensembling based semi-supervised
35 deep learning methods focus on natural images rather than histopathology
36 images, and these two types of images usually have different data distribu-
37 tions. Additionally, most of them fail to consider the relationship among
38 training samples (Luo et al., 2017), and this might result in suboptimal per-
39 formance on histopathology images.

40 Due to the subjective assessment of pathologists on disease grading or

41 classification, it is difficult and costly to obtain large-scale clean labels for
42 histopathology images. Recently, several methods have been developed to
43 enhance the robustness of deep neural networks on noisy labels. Most of ex-
44 isting methods adopt one of the following four strategies: (i) designing robust
45 loss functions (Ghosh et al., 2017); (ii) calculating a transformation matrix
46 (Mnih and Hinton, 2012); (iii) reweighting examples (Ren et al., 2018); (iv)
47 forming consensus predictions under different configurations (Reed et al.,
48 2014). Although these strategies can alleviate the effect of noisy labels, most
49 of them cannot simultaneously explore the semantic information in unlabeled
50 data except several self-ensembling based algorithms, e.g. Temporal Ensem-
51 bling (TE) (Laine and Aila, 2016). It creates an ensemble target for each
52 label prediction of training samples by applying exponential moving average
53 (EMA) to the predictions within multiple previous training epochs, and then
54 minimizes the difference between the label prediction and its ensemble target.
55 Therefore, in this paper, we focus on improving TE to exploit the semantic
56 information of a small amount of labeled histopathology images, explore the
57 information of unlabeled data, and meanwhile suppress the effect of noisy
58 labels.

59 Specifically, we propose a novel robust semi-supervised convolutional neu-
60 ral network, namely graph temporal ensembling (GTE). Inspired by TE, the
61 proposed method creates ensemble targets for feature and label predictions
62 of each training sample and forms consensus predictions under different con-
63 figurations, in order to take advantage of the semantic information from un-
64 labeled data and boost the model robustness to noisy labels (Laine and Aila,
65 2016). Because TE fails to take into account the relationship among labeled

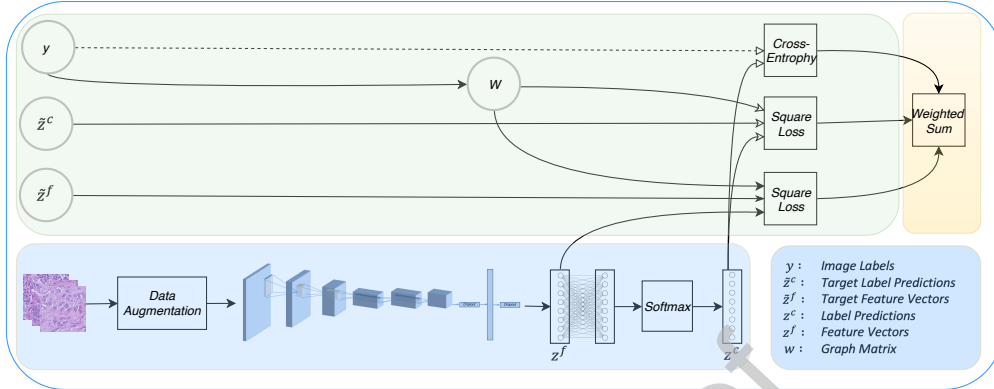


Fig 1: The proposed self-ensembling based deep architecture using AlexNet as the backbone network for histopathology image classification.

66 training samples, we exploit their connections by utilizing a graph to map all
 67 labeled samples from the same class into a single cluster so as to boost the
 68 strength of the ensemble targets. To the best of our knowledge, GTE is the
 69 first self-ensembling based semi-supervised deep method for histopathology
 70 image classification. The major contributions of this paper are summarized
 71 as follows:

- 72 • We present a novel self-ensembling based deep architecture to lever-
 73 age the semantic information of unlabeled histopathology images, and
 74 meanwhile being robust to noisy labels. For clarity, we show the pro-
 75 posed deep architecture in Fig 1.
- 76 • We propose a novel loss function to form consensus feature and label
 77 predictions, i.e. creating stronger ensemble targets for feature and label
 78 predictions via the graph that maps the targets of labeled training im-
 79 ages within the same class into a single cluster, and forming consensus
 80 predictions of all training images under different configurations.

- Extensive experiments on lung and breast cancer image datasets demonstrate that (i) The proposed method can achieve superior performance over recent state-of-the-art semi-supervised deep methods and deep neural networks with noisy labels; (ii) Using a graph-based approach (to map labeled samples of each class into a cluster) is more beneficial to semi-supervision compared to feature consistency (to form consensus predictions of feature representations), while feature consistency is more significant on the model robustness to noisy labels.

This paper is organized as follows. Section 2 briefly introduces the related work including semi-supervised deep methods and deep neural networks with noisy labels. Section 3 presents the proposed method GTE. Section 4 shows and analyzes experimental results on histopathology image classification. Section 5 concludes this paper and points out the future work.

2. Related work

In this section, we briefly review the related work: semi-supervised deep learning and deep neural networks with noisy labels.

2.1. Semi-supervised deep learning

Semi-supervised deep learning algorithms can be roughly categorized into five groups: self-training, multi-view, generative adversarial networks (GAN), graph and self-ensembling. Self-training (Yarowsky, 1995; Lee, 2013) utilizes models' own predictions of unlabeled data in order to attain additional semantic information to improve the model generalization. Multi-view based methods (Zhou and Goldman, 2004; Ruder and Plank, 2018) train different

104 models with different views of data based on the assumption that these views
105 can complement each other and the models can collaborate to boost the per-
106 formance of each other. GAN based semi-supervised methods (Odena, 2016;
107 Kumar et al., 2017) usually utilize the discriminator of GAN as the classifier
108 to obtain the output containing $K + 1$ (K real classes and one fake class)
109 probabilities and adopt the generator to improve the performance of the dis-
110 criminator. Graph based methods (Weston et al., 2012; Kipf and Welling,
111 2016) utilize the transduction of graph to exploit the semantic information
112 of labeled data and meanwhile explore the underlying structure of unlabeled
113 data, i.e. samples that are close in feature spaces should be also close in
114 output spaces (local or global consistency) (Kamnitsas et al., 2018). This is
115 due to the smoothness and clustering assumption. Self-ensembling (Rasmus
116 et al., 2015; Laine and Aila, 2016) is to utilize a single model under different
117 configurations to create a stronger prediction for each training sample, in
118 order to form consensus predictions for boosting model robustness. Because
119 there are numerous semi-supervised deep algorithms, we mainly review the
120 most related ones: graph and self-ensembling based semi-supervised deep
121 methods, several semi-supervised and weakly supervised methods for pathol-
122 ogy image classification in the following.

123 Graph based methods usually construct graphs to preserve the relation-
124 ship of neighbors and then utilize the transduction of the graph to simultane-
125 ously exploit the semantic information of labeled data and explore the under-
126 lying structure of unlabeled data. (Weston et al., 2012) apply “shallow” semi-
127 supervised learning techniques to deep neural networks with adding a graph
128 Laplacian regularizer. Diffusion-convolutional neural networks (DCNN) (At-

129 wood and Towsley, 2016) and graph convolutional networks (GCN) (Kipf
130 and Welling, 2016) are proposed for graph-structured data, and they are
131 transductive and require a pre-constructed graph. (Luo et al., 2018) incor-
132 porate one self-ensembling method, e.g. mean teacher (MT) (Tarvainen and
133 Valpola, 2017), into GCN to further improve its performance. Instead of
134 relying on the pre-constructed graph, (Haeusser et al., 2017) seek the associ-
135 ation between labeled and unlabeled data using a two-step random walk in
136 a feature space, which starts and ends at labeled samples within the same
137 class through one intermediate unlabeled sample. Later, (Kamnitsas et al.,
138 2018) add a regularizer to capture the global structure hidden in the data.
139 Among graph based methods, the one most related to GTE is (Luo et al.,
140 2018). It is suitable for graph-structured data and requires a pre-constructed
141 graph, while GTE is applied to histopathology images and it constructs the
142 graph based on the labels of training samples in each batch.

143 Self-ensembling based methods aim to create strong ensemble predictions
144 of training samples, using the output of one single neural network under dif-
145 ferent training epochs, regularizations, input augmentation conditions, etc.
146 The Γ -model version of ladder network (Rasmus et al., 2015) contains clean
147 and corrupted branches, and the clean branch is to generate proxy labels of
148 corrupted unlabeled data produced by the corrupted branch. Π -model (Laine
149 and Aila, 2016) utilizes two corrupted branches to generate label predictions
150 of training samples, and then applies a consistency cost to minimize the dif-
151 ference between the predictions. Temporal ensembling (TE) (Laine and Aila,
152 2016) only contains one corrupted branch, and it generates an ensemble tar-
153 get for each label prediction by using EMA to aggregate the predictions of

154 multiple previous epochs. To further smooth the model, MT (Tarvainen and
155 Valpola, 2017) averages model weights instead of aggregating label predic-
156 tions. (Su et al., 2019) embed a label propagation step into MT in order to
157 maintain the local and global consistency of predictions. Virtual adversar-
158 ial training (VAT) (Miyato et al., 2018) proposes a regularization method
159 that utilizes a virtual adversarial loss to measure the local smoothness of the
160 conditional label distribution. Smooth neighbors on teacher graph (SNTG)
161 (Luo et al., 2017) constructs a graph using label predictions of the teacher
162 model to measure the similarity of neighbors, in order to learn the represen-
163 tation smoothly on a low-dimensional manifold. Among the aforementioned
164 graph and self-ensembling based methods, SNTG is the work most related to
165 the proposed method GTE. Their major differences are: (i) SNTG samples
166 several label predictions of the teacher model to construct a graph, while
167 GTE constructs the graph using all given labels; (ii) SNTG only learns fea-
168 ture representations on a low-dimensional manifold, but GTE aims to map
169 feature and label predictions of each class into a cluster, respectively, so that
170 it can create stronger ensemble targets for feature and label predictions.

171 Semi-supervised deep learning methods have been widely studied in the
172 fields of natural image recognition, structural data and natural language
173 processing (NLP). However, few efforts are devoted to digit pathology image
174 analysis. A cluster-then-label semi-supervised method (Peikari et al., 2018)
175 identifies high-density regions in the data space to help support vector ma-
176 chine (SVM) find the decision boundary. (Sparks and Madabhushi, 2016)
177 present an image query method based on semi-supervised manifold learning.
178 (Kapil et al., 2018) utilize class auxiliary generative adversarial networks

179 (AC-GANs) for programmed death ligand 1 (PD-L1) scoring on needle biop-
180 sies. In addition to semi-supervised algorithms, multiple instance learning
181 (MIL) and transfer learning (Quelleg et al., 2017; Cheplygina et al., 2019)
182 are also two popular strategies to reduce the workload of label annotations
183 for pathology images. (Xu et al., 2014) propose an effective MIL method,
184 multiple clustered instance learning, to simultaneously perform image-level
185 classification, medical image segmentation and patch-level clustering. (Kraus
186 et al., 2016) combine deep learning and MIL for medical image classification
187 using only whole image labels. (Chang et al., 2017) propose an unsuper-
188 vised transfer learning method using multi-scale convolutional sparse coding
189 to learn transferable base knowledge for medical image classification.

190 *2.2. Deep neural networks with noisy labels*

191 Although several methods have been proposed to improve the robustness
192 of deep neural networks to noisy labels, very limited efforts have been devoted
193 to boosting the model robustness to noisy labels for pathology image classi-
194 fication. (Mnih and Hinton, 2012) propose two robust loss functions based
195 on noise distribution to deal with omission and registration noise in aerial
196 images. (Ghosh et al., 2017) demonstrate that the ℓ_1 -norm based loss func-
197 tion is more robust than cross-entropy and least-squares loss functions. (Veit
198 et al., 2017) utilize the ℓ_1 -norm to clean labels. Unfortunately, the ℓ_1 -norm
199 based loss function is much more difficult to converge. (Natarajan et al.,
200 2013; Sukhbaatar and Fergus, 2014; Xiao et al., 2015; Patrini et al., 2017;
201 Vahdat, 2017) model the relationship between images, class labels and label
202 noise with probability graphs and then integrate them into end-to-end deep
203 neural networks. (Ren et al., 2018) assign weights to training samples based

204 on their gradient directions to decrease the effect of samples with noisy labels.
 205 (Reed et al., 2014; Laine and Aila, 2016; Li et al., 2017) aim to form consen-
 206 sus label predictions under different epochs to boost the model robustness.
 207 Compared to previous methods only forming consensus label predictions, the
 208 proposed GTE aims to form both consensus label and feature predictions to
 209 further improve the consistency and smoothness of predictions.

210 3. Methods

211 In this section, we first provide a brief introduction of the TE method
 212 (Laine and Aila, 2016) and then present the proposed method, GTE.

213 3.1. Temporal ensembling based convolutional neural network

214 Given a set of labeled and unlabeled images $\mathbf{X} = [\mathbf{X}_u; \mathbf{X}_l] = \{\mathbf{x}_i\}_{i=1}^N$,
 215 $\mathbf{y} = \{y_i\}_{i=1}^n$ ($y_i \in \{0, 1, \dots, K-1\}$) denotes the labels of labeled images,
 216 where \mathbf{X}_u and \mathbf{X}_l represent the labeled and unlabeled images, respectively,
 217 \mathbf{x}_i denotes the i -th training image and K is the number of classes. Let
 218 $\mathbf{z}_i^c \in \mathbb{R}^K$ be the label prediction (which is a predicted class probability vector)
 219 of \mathbf{x}_i and $\tilde{\mathbf{z}}_i^c \in \mathbb{R}^K$ represent the ensemble target for label prediction \mathbf{z}_i^c . $\tilde{\mathbf{z}}_i^c$
 220 is obtained by applying EMA to label predictions within multiple previous
 221 training epochs. Specifically, in each training epoch, \mathbf{z}_i^c is firstly accumulated
 222 into an ensemble vector $\tilde{\mathbf{z}}_i^{ce} \in \mathbb{R}^K$, i.e. $\tilde{\mathbf{z}}_i^{ce} = \alpha \tilde{\mathbf{z}}_i^{ce} + (1 - \alpha) \mathbf{z}_i^c$, and then $\tilde{\mathbf{z}}_i^c$ is
 223 computed by $\tilde{\mathbf{z}}_i^c = \tilde{\mathbf{z}}_i^{ce} / (1 - \alpha^t)$, where α is a momentum term to control how
 224 far the ensemble $\tilde{\mathbf{z}}_i^{ce}$ reaches into the training history, and t is the current
 225 number of training epochs.

226 Let L represent the index set of labeled images in \mathbf{X} and B denote the
 227 index set of selected images from \mathbf{X} . The loss function of TE is (Laine and

228 Aila, 2016):

$$J_{TE} = -\frac{1}{|B|} \sum_{i \in (B \cap L)} \log \mathbf{z}_i^c[y_i] + \frac{\tau(t)}{K|B|} \sum_{i \in B} \|\mathbf{z}_i^c - \tilde{\mathbf{z}}_i^c\|_F^2, \quad (1)$$

229 where the first term is a cross-entropy loss function for labeled images, $|B|$
 230 indicates the number of selected images, and $\tau(t)$ indicates a time-dependent
 231 weighting function to gradually enhance the weight of the consistency cost,
 232 e.g. $\sum_{i \in B} \|\mathbf{z}_i^c - \tilde{\mathbf{z}}_i^c\|_F^2$. This term is to explore the semantic information
 233 hidden in unlabeled data and meanwhile smooth the model. Eq. (1) suggests
 234 that TE, which extrapolates the labels of unlabeled images by aggregating
 235 label predictions within multiple previous epochs, can effectively explore the
 236 semantic information in unlabeled data.

237 3.2. Graph temporal ensembling based convolutional neural network

238 TE has obtained promising classification performance on natural images
 239 and handwritten digits. However, it fails to take into account the connection
 240 among labeled images and the consistency of predicting features, which can
 241 further boost model smoothness. Considering these two factors, we construct
 242 a graph to connect feature and label predictions among labeled images, so
 243 that images of each class can be mapped into a cluster, which can create
 244 more robust and stronger ensemble feature and label predictions.

245 First, we define a graph to preserve the relationship among training sam-
 246 ples. Suppose that B represents the index set of images in each batch, and
 247 a symmetric matrix $\mathbf{S} \in \mathbb{R}^{|B| \times |B|}$ maintains their relations. For any two
 248 samples \mathbf{x}_i and \mathbf{x}_j , their relations are described by:

$$s_{ij} = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ belong to the same class} \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

249 where $s_{ij} \in \mathbf{S}$ ($i, j \in B$). Note that for unlabeled data \mathbf{x}_i , only $s_{ii} = 1$ and
 250 $s_{ij} = 0$ for $i \neq j$.

251 Because each class consists of different numbers of samples in one batch,
 252 we normalize the matrix \mathbf{S} as follows:

$$\mathbf{W} = \mathbf{D}^{-\frac{1}{2}} \mathbf{S} \mathbf{D}^{-\frac{1}{2}}, \quad (3)$$

253 where $\mathbf{D} \in \mathbb{R}^{|B| \times |B|}$ is a positive diagonal matrix with the i -th element being
 254 $d_{ii} = \sum_{j=1}^{|B|} s_{ij}$. Thus, the sum of each row of \mathbf{W} is one.

255 Next, we present the loss function of GTE by using the matrix \mathbf{W} . In
 256 addition to the label prediction \mathbf{z}_i^c of the image \mathbf{x}_i , we adopt its feature
 257 representation \mathbf{z}_i^f to enhance the relationship among training samples. For
 258 samples in each batch, to map feature and label predictions of labeled images
 259 within the same class into a cluster, we propose the following loss:

$$J_1 = \frac{1}{K|B|} \sum_{i \in B} \left\| \mathbf{z}_i^f - \sum_{j=1}^{|B|} w_{ij} \mathbf{z}_j^f \right\|_F^2 + \left\| \mathbf{z}_i^c - \sum_{j=1}^{|B|} w_{ij} \mathbf{z}_j^c \right\|_F^2. \quad (4)$$

260 For any unlabeled image \mathbf{x}_i , we have $s_{ii} = 1$ so that $w_{ii} = 1$ and $w_{ij} = 0$ for
 261 $i \neq j$, Eq. (4) will remove the effect of unlabeled samples and only map the
 262 labeled images of each class into a cluster.

263 Let $\tilde{\mathbf{z}}_i^f$ denote the ensemble target of the feature representation \mathbf{z}_i^f . Similar
 264 to $\tilde{\mathbf{z}}_i^c$, $\tilde{\mathbf{z}}_i^f$ is obtained by applying EMA to feature representations within
 265 multiple previous training epochs, i.e. $\mathbf{z}_i^{fe} = \alpha \mathbf{z}_i^{fe} + (1 - \alpha) \mathbf{z}_i^f$, and $\tilde{\mathbf{z}}_i^f =$
 266 $\mathbf{z}_i^{fe} / (1 - \alpha^t)$, where \mathbf{z}_i^{fe} is an ensemble feature vector over multiple previous
 267 feature representations. $\tilde{\mathbf{z}}_i^f$ and $\tilde{\mathbf{z}}_i^c$ usually have better prediction quality
 268 than \mathbf{z}_i^f and \mathbf{z}_i^c , respectively. Additionally, Eq. (4) does not take advantage
 269 of unlabeled data. Therefore, to map labeled images of each class into a

270 cluster and explore the semantic information in unlabeled images, we utilize
 271 $\sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^f$ and $\sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^c$ to replace $\sum_{j=1}^{|B|} w_{ij} \mathbf{z}_i^f$ and $\sum_{j=1}^{|B|} w_{ij} \mathbf{z}_i^c$ in Eq. (4).
 272 It becomes:

$$J_2 = \frac{1}{K|B|} \sum_{i \in B} \left\| \mathbf{z}_i^f - \sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^f \right\|_F^2 + \left\| \mathbf{z}_i^c - \sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^c \right\|_F^2. \quad (5)$$

273 It is worth noting when \mathbf{x}_i is one unlabeled sample, Eq. (5) explores its
 274 information via $\left\| \mathbf{z}_i^f - \tilde{\mathbf{z}}_i^f \right\|_F^2 + \left\| \mathbf{z}_i^c - \tilde{\mathbf{z}}_i^c \right\|_F^2$, because of $w_{ii} = 1$ and $w_{ij} = 0$ for
 275 $i \neq j$.

276 Each image should be more similar to itself than the others under different
 277 configurations, and thus we redefine the matrix \mathbf{S} by using $\mathbf{S} \leftarrow \mathbf{S} + \gamma \mathbf{I}_{n^b}$ to
 278 enhance the consistency of predictions, where $\mathbf{I}_{n^b} \in \mathbb{R}^{n^b \times n^b}$ is an identity
 279 matrix and γ is a non-negative constant to adjust the weight between cluster
 280 mapping and the consistency. The diagonal matrix \mathbf{D} is calculated by $d_{ii} =$
 281 $\sum_{j=1}^{n^b} s_{ij}$ and then \mathbf{W} is computed based on Eq. (3).

282 By introducing the time-dependent weighting function $\tau(t)$ to Eq. (5)
 283 and then combining it with the cross-entropy function on labeled images, we
 284 obtain the loss function of GTE as follows:

$$J = -\frac{1}{|B|} \sum_{i \in (B \cap L)} \log \mathbf{z}_i^c[y_i] + \frac{\tau(t)}{K|B|} \sum_{i,j \in (B)} \left(\lambda_1 \left\| \mathbf{z}_i^f - \sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^f \right\|_F^2 + \lambda_2 \left\| \mathbf{z}_i^c - \sum_{j=1}^{|B|} w_{ij} \tilde{\mathbf{z}}_j^c \right\|_F^2 \right), \quad (6)$$

285 where λ_1 and λ_2 are non-negative constants to balance the three terms. Let
 286 $f_\theta(\cdot)$ represent a stochastic neural network with parameters θ . Based on the
 287 loss function Eq. (6), the parameters θ can be updated by any optimizer, e.g.
 288 Adam (Kingma and Ba, 2014). For clarity, we present the detailed procedure
 289 using Eq. (6) to learn model parameters in Algorithm 1.

Algorithm 1: GTE

Input: Training images $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$, label index set L , labels $\mathbf{y} = \{y_i\}_{i=1}^n$, weighting function $\tau(t)$, parameters $\gamma, \lambda_1, \lambda_2$, ensembling momentum α , stochastic neural network with parameters θ : $f_\theta(x)$, stochastic input augmentation function: $g(x)$

Output: Parameters θ

1. Initialization:

$\mathbf{Z}^{fe} \leftarrow \mathbf{0}_{N \times K}$, \triangleright ensemble previous feature vectors

$\tilde{\mathbf{Z}}^f \leftarrow \mathbf{0}_{N \times K}$, \triangleright target feature vectors

$\mathbf{Z}^{ce} \leftarrow \mathbf{0}_{N \times K}$, \triangleright ensemble previous label predictions

$\tilde{\mathbf{Z}}^c \leftarrow \mathbf{0}_{N \times K}$, \triangleright target label predictions

2. for t in $[1, num_epochs]$ do
3. for each minibatch B do

4. Constructing \mathbf{S} via Eq. (2) and then calculating $\mathbf{S} \leftarrow \mathbf{S} + \gamma \mathbf{I}_{nb}$;

5. Calculating \mathbf{W} via Eq. (3);

6. $\mathbf{z}_{i \in B}^c, \mathbf{z}_{i \in B}^f \leftarrow f_\theta(g(\mathbf{x}_{i \in B}))$;

7. loss \leftarrow Eq. (6);

8. updating θ using optimizers, e.g. Adam (Kingma and Ba, 2014);

9. end for

10. $\mathbf{Z}^{fe} \leftarrow \alpha \mathbf{Z}^{fe} + (1 - \alpha) \mathbf{Z}^f$

11. $\tilde{\mathbf{Z}}^f \leftarrow \mathbf{Z}^{fe} / (1 - \alpha^t)$

12. $\mathbf{Z}^{ce} \leftarrow \alpha \mathbf{Z}^{ce} + (1 - \alpha) \mathbf{Z}^c$

13. $\tilde{\mathbf{Z}}^c \leftarrow \mathbf{Z}^c / (1 - \alpha^t)$

14. end for

290 3.3. Implementation details

291 We implement GTE with the PyTorch framework and adopt AlexNet
 292 (Krizhevsky et al., 2012) as our backbone network, which is pre-trained on
 293 the ImageNet database (Deng et al., 2009). The maximum learning rate η is
 294 0.00005 and the Adam momentum parameters are $\beta_1 = 0.9$ and $\beta_2 = 0.999$.
 295 For GTE, we empirically set the parameters $\alpha = 0.6$, $\lambda_1 = \lambda_2 = 0.1$ and
 296 $\gamma = 1$. We totally run 100 epochs for training with mini-batch size of 40. The
 297 time-dependent weighting function $\tau(t)$ ramps up from 0 to 1 during the first
 298 40 epochs. $\tau(t)$ is a Gaussian ramp-up curve $e^{-5\|1-T\|^2}$, where T advances
 299 linearly from 0 to 1 during the first 40 epochs. The learning rate η also
 300 gradually increases to the maximum using $\tau(t)$ during the first 40 epochs and
 301 then keeps unchanged in the following 30 epochs, but decreases to 0 by using
 302 a Gaussian ramp-down curve during the last 30 epochs. The ramp-down
 303 curve is similar to the ramp-up curve while using a scaling constant 12.5, and
 304 linearly decreasing T from 1 to 0 during the last 30 epochs. Additionally,
 305 the Adam momentum parameter β_1 ramps down to 0.5 during the last 30
 306 epochs.

307 4. Experimental Results and Analysis

308 To evaluate the proposed method, GTE, we conduct experiments on two
 309 histopathology image datasets consisting of lung cancer and breast cancer
 310 images, respectively. All images are selected from The Cancer Genome Atlas
 311 (TCGA) and each of them is stained with Hematoxylin and eosin (H&E).
 312 Both lung and breast cancer datasets contain images with both 10x and 20x
 313 magnifications.

314 **The lung cancer dataset** contains two types of diseases: adenocarci-
315 noma (AC) and squamous cell carcinoma (SC). We crop and select 2,904
316 (1,456 AC and 1,448 SC) lung cancer image patches from whole slide images
317 of 42 patients. The size of each cropped image patch is 500×500 . Each
318 patient is composed of a set of image patches, and the number of patches
319 varies from 24 to 196. We randomly select 30 patients containing 2104 im-
320 age patches (1,008 AC and 1,096 SC) to construct a training set, and the
321 remaining 12 patients with 800 image patches (448 AC and 352 SC) are used
322 for testing.

323 **The breast cancer dataset** consists of 1,763 image patches selected
324 from whole slide images of 193 patients in TCGA. The nuclear pleomorphism
325 score of each whole slide image is annotated by pathologists. Then we select
326 and crop image patches corresponding to the nuclear pleomorphism score of
327 whole slide images. Each image patch is cropped with the size of 1000×1000 .
328 When one slide includes examples of several grades, we only select the patches
329 corresponding to the highest grade. The number of image patches for each
330 patient ranges from 1 to 15. These patches are classified into three groups
331 by grading the nuclear pleomorphism based on the nucleus, shape and size
332 of cells. The grading nuclear pleomorphism score is judged by the following
333 rules (Galea et al., 1992):

334 **Score 1:** Cells are uniform in size and shape compared to healthy cells
335 (breast epithelial cells);

336 **Score 2:** Cells might be more extensive and slightly variational in size and
337 shape compared to healthy cells;

338 **Score 3:** Cells have remarkable variation in size and shape compared to

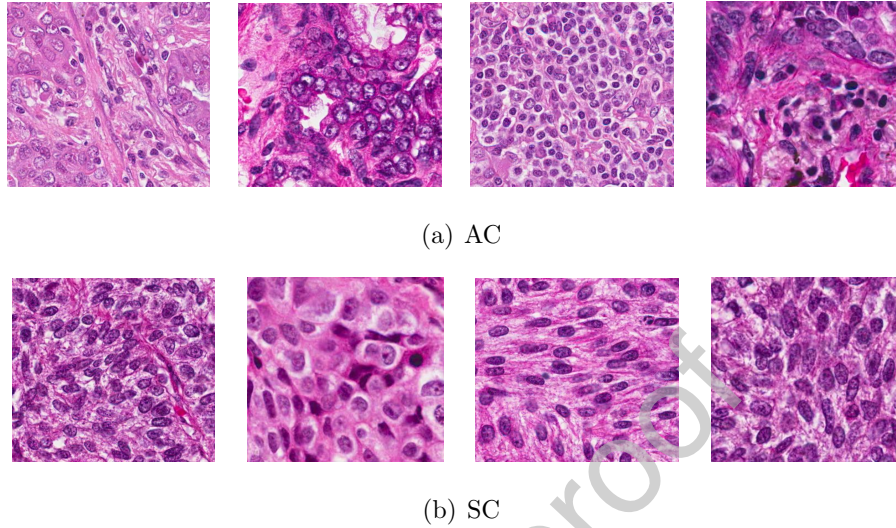


Fig 2: Examples of lung cancer images.

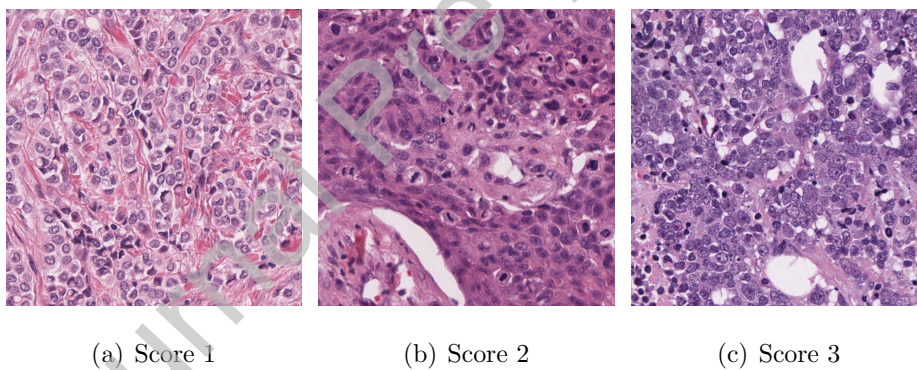


Fig 3: Examples of breast cancer images with different nuclear pleomorphism scores.

339 healthy cells;

340 Fig 2 and Fig. 3 show several example lung cancer images with two types
 341 of diseases and breast cancer images with the nuclear pleomorphism graded
 342 by scores 1-3, respectively. Because the images with nuclear pleomorphism
 343 scores 2 and 3 are very difficult to distinguish even for pathologists, we com-
 344 bine the images with nuclear pleomorphism scores 2 and 3 into one class. As

345 a result, the selected breast cancer images contain two classes. We name nu-
 346 clear pleomorphism score 1 as class 1 and label nuclear pleomorphism scores
 347 2-3 as class 2. We randomly select 1,421 images (476 class 1 and 945 class
 348 2) of 152 patients for training, and the remaining 342 images (88 class 1 and
 349 254 class 2) of 41 patients are utilized for testing. All labels and nuclear
 350 pleomorphism scores of image patches are confirmed by pathologists.

351 10x and 20x image patches are fed into a single network together, in
 352 order to enhance the model robustness to the scale of cellular information.
 353 All of them are resized to 224×224 as the given input and training images
 354 are augmented with random translations ($\{\Delta x, \Delta y\} \sim [-32, 32]$), horizontal
 355 flips ($p = 0.5$), Gaussian noise ($\sigma = 0.15$) and color normalization.

356 4.1. Experiments for semi-supervision

357 We compare GTE against six popular methods, including one supervised
 358 method: AlexNet with the cross-entropy loss function (Baseline), and five
 359 semi-supervised methods: Π -model and TE (Laine and Aila, 2016), MT (Tar-
 360 vainen and Valpola, 2017), VAT (Miyato et al., 2018) and SNTG (Luo et al.,
 361 2017). Here, SNTG utilizes TE as the base. For fairness, all the five semi-
 362 supervised methods adopt the pre-trained AlexNet as the backbone network.
 363 Additionally, we set their maximum learning rate to be 0.00005, upon which
 364 they can obtain the best or suboptimal performance on the pre-trained net-
 365 work. We randomly select 20%, 35% and 50% patients of each class from
 366 the lung or breast cancer training set to construct a labeled set, and the
 367 remaining patients are used as unlabeled ones. We repeat this process ten
 368 times and report the average image classification accuracy, recall, F-score of
 369 various methods. Additionally, we also present their performance when all

Table 1: Results for semi-supervision on the lung cancer dataset. We bold the best results and highlight the second best results via underlines.

Method	Accuracy			Recall	F-Score
	AC	SC	Avg		
20% labeled patients					
Baseline	0.875 ± 0.099	0.591 ± 0.300	0.750 ± 0.106	0.591 ± 0.230	0.634 ± 0.248
II-model	0.832 ± 0.201	0.676 ± 0.166	0.763 ± 0.130	0.636 ± 0.188	0.672 ± 0.150
TE	0.855 ± 0.121	0.620 ± 0.284	0.752 ± 0.119	0.620 ± 0.284	0.657 ± 0.234
MT	0.855 ± 0.099	0.567 ± 0.276	0.728 ± 0.091	0.543 ± 0.248	0.590 ± 0.211
VAT	0.827 ± 0.224	0.696 ± 0.259	<u>0.769 ± 0.115</u>	<u>0.692 ± 0.259</u>	<u>0.700 ± 0.197</u>
SNTG	0.838 ± 0.135	0.491 ± 0.268	0.685 ± 0.089	0.460 ± 0.242	0.515 ± 0.217
GTE	0.921 ± 0.046	0.886 ± 0.101	0.905 ± 0.041	0.886 ± 0.101	0.894 ± 0.071
35% labeled patients					
Baseline	0.842 ± 0.149	0.767 ± 0.230	0.808 ± 0.147	0.761 ± 0.231	<u>0.774 ± 0.204</u>
II-model	0.893 ± 0.061	0.765 ± 0.124	<u>0.837 ± 0.060</u>	0.730 ± 0.155	0.766 ± 0.119
TE	0.855 ± 0.114	0.772 ± 0.197	0.818 ± 0.081	<u>0.772 ± 0.197</u>	0.785 ± 0.147
MT	0.862 ± 0.098	0.750 ± 0.196	0.813 ± 0.100	0.738 ± 0.193	0.759 ± 0.158
VAT	0.902 ± 0.065	0.581 ± 0.310	0.761 ± 0.122	0.571 ± 0.311	0.622 ± 0.258
SNTG	0.782 ± 0.210	0.751 ± 0.210	0.768 ± 0.152	0.716 ± 0.207	0.722 ± 0.181
GTE	0.936 ± 0.026	0.878 ± 0.059	0.910 ± 0.026	0.878 ± 0.059	0.893 ± 0.042
50% labeled patients					
Baseline	0.897 ± 0.087	0.584 ± 0.225	0.759 ± 0.090	0.584 ± 0.225	0.645 ± 0.188
II-model	0.925 ± 0.044	0.696 ± 0.239	0.824 ± 0.104	0.688 ± 0.232	0.735 ± 0.191
TE	0.891 ± 0.060	0.821 ± 0.067	0.860 ± 0.045	<u>0.821 ± 0.067</u>	<u>0.839 ± 0.052</u>
MT	0.905 ± 0.060	0.733 ± 0.267	0.829 ± 0.129	0.733 ± 0.264	0.759 ± 0.241
VAT	0.921 ± 0.076	0.654 ± 0.266	0.804 ± 0.125	0.644 ± 0.262	0.699 ± 0.222
SNTG	0.933 ± 0.032	0.820 ± 0.169	<u>0.883 ± 0.081</u>	0.802 ± 0.183	0.833 ± 0.145
GTE	0.933 ± 0.034	0.884 ± 0.063	0.911 ± 0.019	0.884 ± 0.063	0.898 ± 0.039
All labeled patients					
Baseline	0.953 ± 0.028	0.922 ± 0.031	0.939 ± 0.011	0.917 ± 0.025	<u>0.928 ± 0.017</u>
II-model	0.965 ± 0.013	0.897 ± 0.036	0.935 ± 0.010	0.884 ± 0.048	0.904 ± 0.032
TE	0.951 ± 0.024	0.897 ± 0.055	0.927 ± 0.024	0.871 ± 0.055	0.897 ± 0.039
MT	0.935 ± 0.030	0.921 ± 0.019	0.928 ± 0.009	0.896 ± 0.027	0.906 ± 0.017
VAT	0.950 ± 0.023	0.942 ± 0.026	0.947 ± 0.006	0.939 ± 0.026	0.941 ± 0.014
SNTG	0.951 ± 0.019	0.916 ± 0.020	0.936 ± 0.006	0.910 ± 0.019	0.921 ± 0.011
GTE	0.954 ± 0.013	0.923 ± 0.035	<u>0.940 ± 0.012</u>	<u>0.919 ± 0.035</u>	<u>0.928 ± 0.025</u>

Table 2: Results for semi-supervision on the breast cancer dataset. We bold the best results and highlight the second best results via underlines.

Method	Accuracy			Recall	F-Score
	Class 1	Class 2	Avg		
20% labeled patients					
Baseline	0.523 ± 0.289	0.896 ± 0.051	0.800 ± 0.058	0.896 ± 0.051	0.844 ± 0.037
II-model	0.375 ± 0.367	0.963 ± 0.040	0.812 ± 0.063	0.958 ± 0.036	<u>0.875</u> ± 0.043
TE	0.572 ± 0.246	0.884 ± 0.092	0.804 ± 0.057	0.884 ± 0.092	0.841 ± 0.064
MT	0.401 ± 0.273	0.938 ± 0.060	0.807 ± 0.068	0.938 ± 0.111	0.851 ± 0.100
VAT	0.480 ± 0.196	0.937 ± 0.068	<u>0.820</u> ± 0.064	0.922 ± 0.076	0.860 ± 0.066
SNTG	0.488 ± 0.285	0.935 ± 0.065	<u>0.820</u> ± 0.060	0.922 ± 0.068	0.861 ± 0.059
GTE	0.738 ± 0.111	0.949 ± 0.016	0.895 ± 0.023	<u>0.951</u> ± 0.016	0.922 ± 0.014
35% labeled patients					
Baseline	0.489 ± 0.235	0.939 ± 0.049	0.823 ± 0.040	0.939 ± 0.047	0.876 ± 0.027
II-model	0.600 ± 0.192	0.898 ± 0.053	0.821 ± 0.054	0.899 ± 0.053	0.853 ± 0.044
TE	0.564 ± 0.199	0.935 ± 0.045	<u>0.840</u> ± 0.035	0.935 ± 0.049	<u>0.884</u> ± 0.029
MT	0.393 ± 0.229	0.956 ± 0.033	0.811 ± 0.056	0.970 ± 0.027	0.872 ± 0.023
VAT	0.696 ± 0.258	0.873 ± 0.077	0.827 ± 0.093	0.845 ± 0.098	0.824 ± 0.096
SNTG	0.581 ± 0.279	0.913 ± 0.076	0.827 ± 0.066	0.893 ± 0.075	0.849 ± 0.059
GTE	0.751 ± 0.119	0.957 ± 0.025	0.904 ± 0.026	<u>0.959</u> ± 0.025	0.930 ± 0.019
50% labeled patients					
Baseline	0.642 ± 0.183	0.905 ± 0.085	0.837 ± 0.060	0.905 ± 0.085	0.868 ± 0.008
II-model	0.599 ± 0.254	0.933 ± 0.070	0.847 ± 0.060	0.932 ± 0.069	0.884 ± 0.055
TE	0.638 ± 0.186	0.901 ± 0.055	0.833 ± 0.050	0.901 ± 0.055	0.865 ± 0.044
MT	0.581 ± 0.224	0.955 ± 0.042	0.835 ± 0.040	0.958 ± 0.040	0.884 ± 0.021
VAT	0.758 ± 0.143	0.956 ± 0.053	<u>0.905</u> ± 0.045	0.945 ± 0.055	<u>0.916</u> ± 0.043
SNTG	0.514 ± 0.273	0.954 ± 0.024	0.841 ± 0.073	0.937 ± 0.049	0.878 ± 0.067
GTE	0.791 ± 0.053	0.950 ± 0.024	0.909 ± 0.023	<u>0.950</u> ± 0.024	0.929 ± 0.023
All labeled patients					
Baseline	0.825 ± 0.026	0.965 ± 0.012	<u>0.919</u> ± 0.007	0.965 ± 0.012	<u>0.941</u> ± 0.008
II-model	0.768 ± 0.083	0.947 ± 0.029	0.901 ± 0.019	0.954 ± 0.023	0.923 ± 0.014
TE	0.779 ± 0.044	0.963 ± 0.016	0.916 ± 0.011	0.960 ± 0.015	0.938 ± 0.011
MT	0.719 ± 0.080	0.911 ± 0.033	0.862 ± 0.021	0.931 ± 0.032	0.891 ± 0.023
VAT	0.848 ± 0.013	0.961 ± 0.015	0.932 ± 0.010	<u>0.963</u> ± 0.018	0.942 ± 0.015
SNTG	0.753 ± 0.059	0.928 ± 0.019	0.883 ± 0.017	0.914 ± 0.024	0.889 ± 0.020
GTE	0.798 ± 0.039	0.961 ± 0.018	<u>0.919</u> ± 0.008	<u>0.963</u> ± 0.016	0.940 ± 0.011

370 training images are labeled.

371 Table 1 presents the image classification accuracy, recall and F-score of
372 various methods on the lung cancer image dataset. GTE can achieve superior
373 performance over the others when 20%, 35% and 50% patients are labeled.
374 For example, the gain of GTE in terms of the average accuracy ranges from
375 3.2% to 17.7% over the best competitors on 20%, 35% and 50% labeled
376 patients. Additionally, GTE obtains the best recall and F-score among all
377 methods on the three cases. When all patients are labeled, GTE attains
378 inferior performance to VAT, while it achieves the similar performance to
379 Baseline and outperforms the other four methods. Table 2 shows the image
380 classification results on the breast cancer image dataset. GTE can also obtain
381 superior performance over the other six methods when 20%, 35% and 50%
382 patients are labeled. For instance, the gain of GTE in accuracy is from 0.4%
383 to 9.2% over the best competitors on 20%, 35% and 50% labeled patients.
384 When using all labeled patients, GTE exhibits similar performance to Base-
385 line and outperforms the other methods except VAT. As shown in Tables 1
386 and 2, some methods with less labeled patients obtain better performance
387 than that with more labeled patients, like VAT using 20% and 35% labeled
388 patients in Table 1, and TE using 35% and 50% labeled patients in Table 2.
389 *This might be because each patient in training and test sets contains different*
390 *number of images, leading to distinct significance of patients for model train-*
391 *ing.* To more clearly compare their classification performance, we show the
392 average accuracy and F-score of seven methods on lung and breast cancer
393 databases in Fig 4.

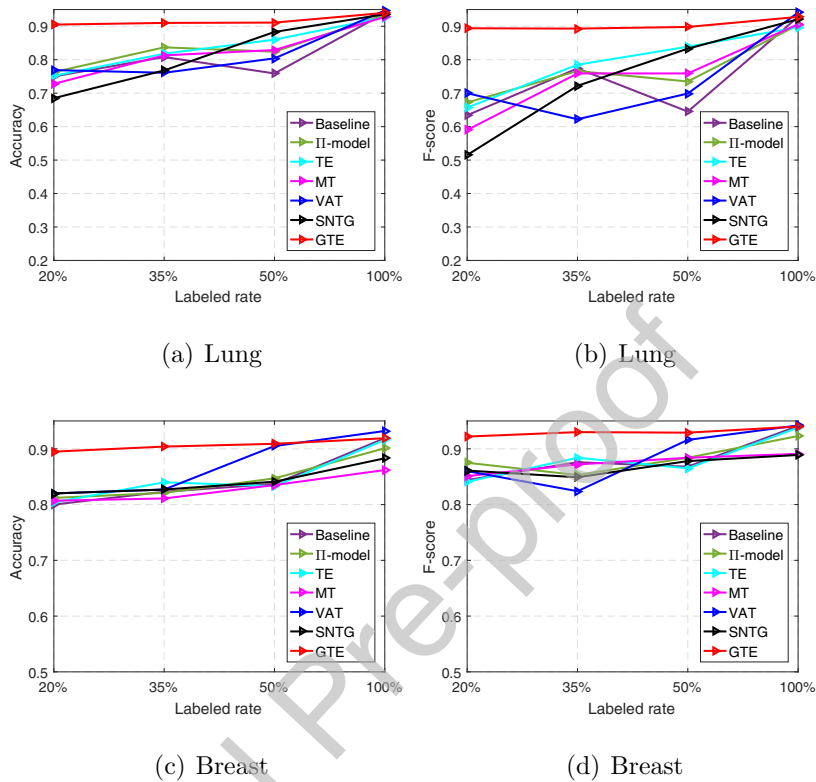


Fig 4: The average accuracy and F-score of seven methods with different rates of labeled patients for semi-supervised learning on lung and breast cancer databases.

394 4.2. Experiments for noisy labels

395 We compare GTE against six popular methods, including one baseline,
 396 AlexNet with the cross-entropy loss function, and five state-of-the-art ap-
 397 proaches, such as Bootstrap (Reed et al., 2014), Forward and Backward (Pa-
 398 trini et al., 2017), TE (Laine and Aila, 2016) and SNTG (Luo et al., 2017).
 399 These approaches also utilize the pre-trained AlexNet as the backbone net-
 400 work and their maximum learning rate is 0.00005. We generate symmetric
 401 label noise by randomly flipping 10% and 20% training patients of each class

402 into the other class. We repeat this process ten times and present their
403 average image classification accuracy, recall and F-score.

404 Tables 3 illustrates that GTE can obtain competitive or even better av-
405 erage image classification accuracy and F-score than the others on lung and
406 breast cancer datasets. For example, GTE only achieves 0.2% higher average
407 accuracy than Backward when using 10% patients with noisy labels on the
408 lung caner dataset, but it obtains 4.2% higher accuracy than Backward when
409 20% patients are with noisy labels. Therefore, GTE is more robust to noisy
410 labels than the other methods.

411 4.3. Experiments for semi-supervision and noisy labels

412 To better demonstrate the strength of GTE, we conduct experiments on
413 the training set with noisy labeled images and unlabeled ones. Specifically,
414 we uniformly select 50% patients from the lung or breast cancer training set
415 to construct a labeled set and utilize the remaining patients as unlabeled
416 ones. Then we randomly flip 20% labeled patients of each class into the
417 other class. We run the experiment ten times and present the average image
418 classification accuracy, recall and F-score of four methods: Baseline, TE,
419 SNTG and GTE.

420 Fig 5 shows their results on lung and breast cancer datasets. It demon-
421 strates that GTE outperforms the other three methods when using noisy
422 labeled and unlabeled data for lung and breast cancer image classification.
423 Specifically, the gain of GTE in terms of average accuracy is 3.1% and 14.5%
424 over the best competitors on lung and breast cancer datasets, respectively.
425 GTE also achieves better F-score than the other three methods.

426 We further utilize a popular dimensionality reduction method, t-SNE

Table 3: Results for noisy labels on the lung cancer dataset. We bold the best results and highlight the second best results via underlines.

Method	Accuracy			Recall	F-Score
	AC	SC	Avg		
Lung cancer					
10% patients with noisy labels					
Baseline	0.905 ± 0.033	0.908 ± 0.035	0.906 ± 0.029	0.898 ± 0.035	0.896 ± 0.032
Bootstrap	0.916 ± 0.076	0.857 ± 0.135	0.890 ± 0.101	0.857 ± 0.135	0.871 ± 0.124
Forward	0.902 ± 0.071	0.834 ± 0.189	0.872 ± 0.121	0.832 ± 0.188	0.844 ± 0.160
Backward	0.938 ± 0.021	0.889 ± 0.049	<u>0.916 ± 0.019</u>	0.884 ± 0.049	0.899 ± 0.024
TE	0.928 ± 0.030	0.887 ± 0.033	0.910 ± 0.020	<u>0.887 ± 0.033</u>	<u>0.899 ± 0.022</u>
SNTG	0.941 ± 0.038	0.828 ± 0.161	0.891 ± 0.084	0.824 ± 0.154	0.860 ± 0.124
GTE	0.946 ± 0.025	0.883 ± 0.035	0.918 ± 0.022	0.883 ± 0.035	0.903 ± 0.036
20% patients with noisy labels					
Baseline	0.860 ± 0.071	0.778 ± 0.207	0.824 ± 0.100	0.778 ± 0.207	0.783 ± 0.148
Bootstrap	0.869 ± 0.090	0.756 ± 0.175	0.819 ± 0.074	0.756 ± 0.175	0.777 ± 0.122
Forward	0.884 ± 0.074	0.845 ± 0.076	<u>0.867 ± 0.065</u>	0.845 ± 0.076	<u>0.848 ± 0.072</u>
Backward	0.822 ± 0.132	0.847 ± 0.107	0.833 ± 0.105	0.847 ± 0.107	0.819 ± 0.108
TE	0.864 ± 0.055	0.861 ± 0.089	0.863 ± 0.064	0.861 ± 0.089	0.847 ± 0.073
SNTG	0.891 ± 0.068	0.697 ± 0.201	0.806 ± 0.125	0.695 ± 0.184	0.735 ± 0.172
GTE	0.895 ± 0.059	0.850 ± 0.097	0.875 ± 0.060	<u>0.850 ± 0.097</u>	0.855 ± 0.075
Breast cancer					
10% patients with noisy labels					
Baseline	0.749 ± 0.098	0.886 ± 0.046	0.851 ± 0.046	0.886 ± 0.046	0.892 ± 0.032
Bootstrap	0.708 ± 0.096	0.919 ± 0.046	<u>0.865 ± 0.035</u>	0.919 ± 0.046	<u>0.910 ± 0.025</u>
Forward	0.753 ± 0.063	0.876 ± 0.066	0.844 ± 0.043	0.876 ± 0.066	0.892 ± 0.034
Backward	0.780 ± 0.082	0.878 ± 0.065	0.853 ± 0.054	0.878 ± 0.065	0.898 ± 0.041
TE	0.709 ± 0.252	0.897 ± 0.071	0.849 ± 0.067	0.897 ± 0.071	0.898 ± 0.044
SNTG	0.765 ± 0.095	0.700 ± 0.830	0.854 ± 0.053	0.870 ± 0.041	0.886 ± 0.033
GTE	0.748 ± 0.119	0.911 ± 0.063	0.869 ± 0.058	<u>0.911 ± 0.063</u>	0.913 ± 0.041
20% patients with noisy labels					
Baseline	0.641 ± 0.259	0.849 ± 0.102	0.796 ± 0.076	0.849 ± 0.102	0.854 ± 0.058
Bootstrap	0.677 ± 0.254	0.837 ± 0.071	0.796 ± 0.050	0.837 ± 0.071	0.859 ± 0.062
Forward	0.656 ± 0.130	0.865 ± 0.060	0.811 ± 0.060	0.863 ± 0.069	0.871 ± 0.043
Backward	0.659 ± 0.133	0.857 ± 0.067	0.806 ± 0.053	0.856 ± 0.069	0.867 ± 0.040
TE	0.644 ± 0.209	0.870 ± 0.064	<u>0.812 ± 0.056</u>	0.874 ± 0.064	<u>0.883 ± 0.044</u>
SNTG	0.700 ± 0.147	0.830 ± 0.065	0.797 ± 0.054	0.826 ± 0.063	0.854 ± 0.036
GTE	0.703 ± 0.140	0.855 ± 0.096	0.815 ± 0.058	<u>0.872 ± 0.086</u>	0.889 ± 0.048

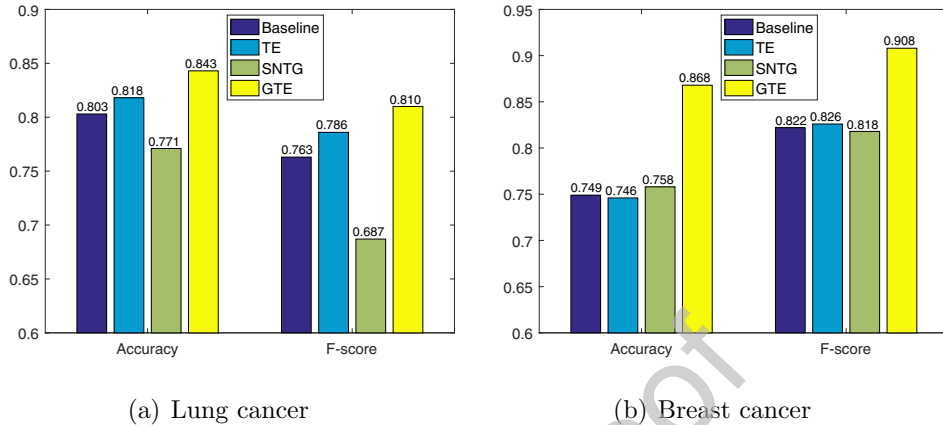


Fig 5: Results including average accuracy and F-score for semi-supervision and noisy labels on the lung and breast cancer datasets. Each experiment is repeated 10 times and an average result is reported.

427 (Maaten and Hinton, 2008), to project the feature vectors of lung cancer
 428 test images obtained by different methods (Baseline, TE, SNTG and GTE)
 429 onto a two-dimensional plane. As shown in Fig. 6, GTE can make features
 430 between classes more separable, because GTE aims to map images of each
 431 class into a cluster, leading to more compact features in each class.

432 4.4. Effects of Parameter Selection

433 First, we show the influence of the feature and graph in GTE on semi-
 434 supervision and noisy labels. Specifically, we set $\lambda_1 = 0$ to remove the effect
 435 of features, and utilize $\mathbf{S} = \mathbf{I}_{n_b}$ to eliminate the influence of graph. We
 436 conduct semi-supervised experiments on two databases, and we randomly
 437 select 20% patients from the training set of each database to construct a
 438 labeled set, and the remaining images are used as unlabeled ones. For noisy
 439 label experiments, we also utilize the two training sets and randomly flip 20%

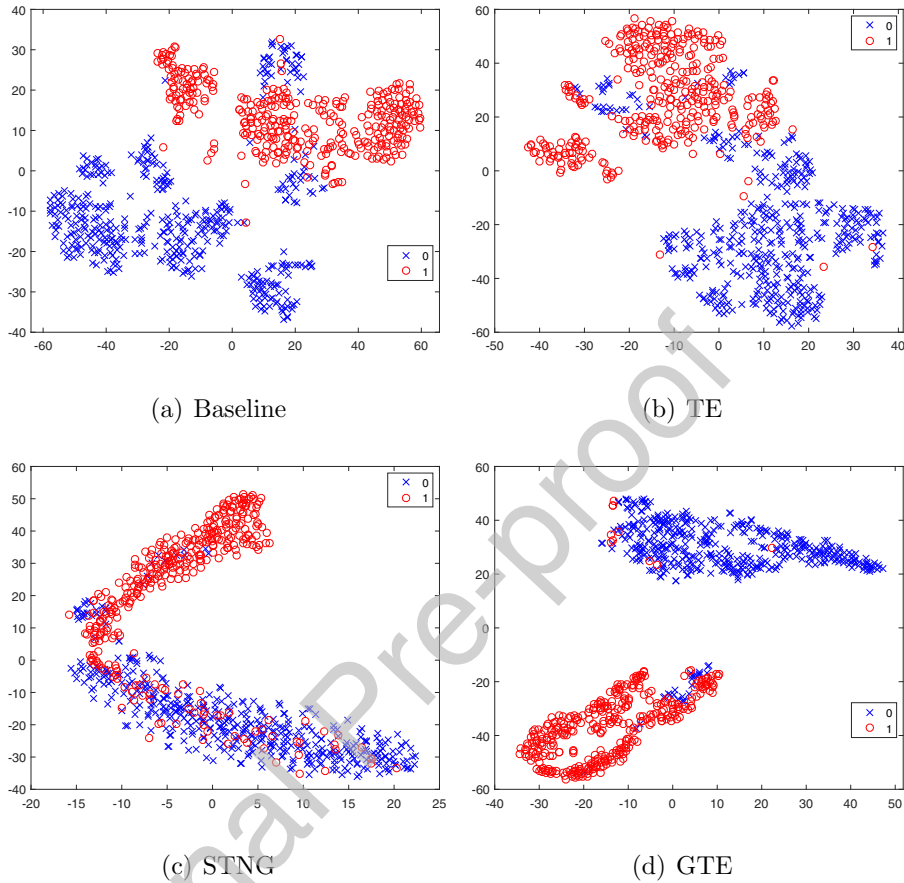


Fig 6: Feature projection of lung cancer test images onto a two-dimensional plane by using Baseline, TE, STNG and GTE. 0 and 1 represent AC and SC, respectively.

440 patients of each class into the other class. We repeat this process ten times
 441 and report the average image accuracy, recall and F-score in Tables 4-5, which
 442 illustrate that GTE has superior performance over $GTE(\lambda_1 = 0)$ and GTE
 443 ($\mathbf{S} = \mathbf{I}_{n^b}$). Additionally, $GTE(\lambda_1 = 0)$ achieves better performance than
 444 $GTE(\mathbf{S} = \mathbf{I}_{n^b})$ on semi-supervised experiments, but $GTE(\mathbf{S} = \mathbf{I}_{n^b})$ performs
 445 better than $GTE(\lambda_1 = 0)$ on experiments for noisy labels. They suggest that
 446 graph (mapping samples of each class into a cluster) is more important than

Table 4: Results on the lung cancer image dataset. We bold the best results and highlight the second best results via underlines.

Method	Accuracy			Recall	F-Score
	Class 1	Class 2	Avg		
Semi-supervised					
GTE	0.921 ± 0.046	0.886 ± 0.101	0.905 ± 0.041	0.886 ± 0.101	0.894 ± 0.071
GTE ($\lambda_1 = 0$)	0.935 ± 0.017	0.831 ± 0.150	<u>0.890</u> ± 0.069	<u>0.831</u> ± 0.150	<u>0.862</u> ± 0.110
GTE ($\mathbf{S} = \mathbf{I}_{n^b}$)	0.875 ± 0.086	0.556 ± 0.220	0.735 ± 0.091	0.556 ± 0.220	0.628 ± 0.177
Noisy labels					
GTE	0.895 ± 0.059	0.850 ± 0.097	0.875 ± 0.060	0.850 ± 0.097	0.855 ± 0.075
GTE ($\lambda_1 = 0$)	0.884 ± 0.059	0.756 ± 0.203	0.828 ± 0.106	0.756 ± 0.023	0.783 ± 0.164
GTE ($\mathbf{S} = \mathbf{I}_{n^b}$)	0.936 ± 0.041	0.783 ± 0.147	<u>0.869</u> ± 0.048	<u>0.790</u> ± 0.142	<u>0.822</u> ± 0.123

Table 5: Results on the breast cancer image dataset. We bold the best results and highlight the second best results via underlines.

Method	Accuracy			Recall	F-Score
	Class 1	Class 2	Avg		
Semi-supervised					
GTE	0.738 ± 0.111	0.949 ± 0.016	0.895 ± 0.023	<u>0.951</u> ± 0.016	0.922 ± 0.014
GTE ($\lambda_1 = 0$)	0.556 ± 0.223	0.967 ± 0.029	<u>0.861</u> ± 0.040	0.967 ± 0.029	<u>0.913</u> ± 0.021
GTE ($\mathbf{S} = \mathbf{I}_{n^b}$)	0.485 ± 0.320	0.924 ± 0.045	0.811 ± 0.092	0.924 ± 0.045	0.881 ± 0.055
Noisy labels					
GTE	0.703 ± 0.140	0.855 ± 0.096	0.815 ± 0.058	0.872 ± 0.086	0.889 ± 0.048
GTE ($\lambda_1 = 0$)	0.575 ± 0.255	0.813 ± 0.080	0.752 ± 0.054	0.813 ± 0.080	0.829 ± 0.038
GTE ($\mathbf{S} = \mathbf{I}_{n^b}$)	0.694 ± 0.147	0.851 ± 0.060	<u>0.810</u> ± 0.054	<u>0.863</u> ± 0.060	<u>0.877</u> ± 0.038

447 features (forming consensus feature predictions) on semi-supervision, while
 448 feature consistency is more significant on the robustness to noisy labels.

449 Then we present the effect of three essential parameters γ , λ_1 and λ_2
 450 in GTE for semi-supervision and noisy labels with lung cancer images. For
 451 semi-supervised experiments, we randomly select 50% patients from the lung
 452 cancer training set to construct a labeled set, and utilize the remaining
 453 images as unlabeled ones; for noisy label experiments, we also utilize the

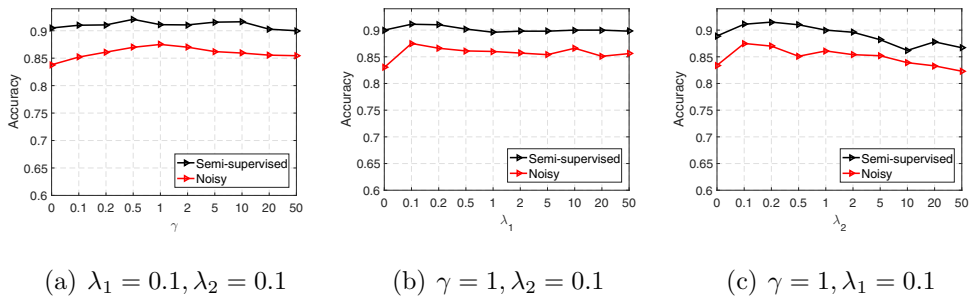


Fig 7: The effect of different values of γ , λ_1 and λ_2 on the proposed GTE for semi-supervision and noisy labels with lung cancer images. Each experiment is repeated 10 times and an average result is reported.

454 lung cancer training set and randomly flip 20% patients of each class into
 455 the other class. We run the experiments ten times and report the aver-
 456 age image accuracy of different values of γ , λ_1 and λ_2 during the range
 457 of $[0, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20, 50]$ in Fig 7, upon which we can see that
 458 $\gamma = 1, \lambda_1 = \lambda_2 = 0.1$ can achieve the best or suboptimal accuracy for GTE
 459 on both semi-supervision and noisy labels. Similar findings can be observed
 460 on other percentages and breast cancer images.

461 5. Conclusion and Future Work

462 In this paper, we propose a novel robust deep architecture to leverage
 463 the semantic information of both labeled and unlabeled data. The proposed
 464 architecture utilizes EMA to create ensemble targets for feature and label
 465 predictions, and then adopts a graph to map labeled images of each class
 466 into a cluster to boost the strength of ensemble predictions. Meanwhile, it
 467 applies the consistency cost on feature representations and label predictions
 468 of all training images, in order to form consensus predictions under different

469 configurations. Experiments on lung and breast cancer datasets demonstrate
470 the effectiveness and efficiency of the proposed method. Additionally, exper-
471 iments illustrate that mapping labeled images of each class into a cluster
472 is more beneficial to semi-supervised classification and forming consensus
473 feature predictions is more helpful to the model robustness.

474 Although the proposed robust semi-supervised deep method has achieved
475 promising performance on the two datasets, there is still much work to ex-
476 plore: (i) deep neural networks on semi-supervised learning and with noisy
477 labels for multi-class histopathology image classification; (ii) training deep
478 neural networks with extremely noisy labels; (iii) robust semi-supervised deep
479 hashing for histopathology image retrieval (because image retrieval can not
480 only provide the class information of images, but also search the most simi-
481 lar images); (iv) robust semi-supervised deep neural networks for whole-slide
482 image classification.

483 References

- 484 Atwood, J., Towsley, D., 2016. Diffusion-convolutional neural networks. In: Pro-
485 ceedings of Advances in Neural Information Processing Systems. pp. 1993–2001.
- 486 Bai, W., Oktay, O., Sinclair, M., Suzuki, H., Rajchl, M., Tarroni, G., Glocker, B.,
487 King, A., Matthews, P. M., Rueckert, D., 2017. Semi-supervised learning for
488 network-based cardiac mr image segmentation. In: International Conference on
489 Medical Image Computing and Computer-Assisted Intervention. pp. 253–260.
- 490 Chang, H., Han, J., Zhong, C., Snijders, A. M., Mao, J.-H., 2017. Unsuper-
491 vised transfer learning via multi-scale convolutional sparse coding for biomed-

- 492 cal applications. *IEEE transactions on pattern analysis and machine intelligence*
493 40 (5), 1182–1194.
- 494 Chen, P., Gao, L., Shi, X., Allen, K., Yang, L., 2019. Fully automatic knee os-
495 teoarthritis severity grading using deep neural networks with a novel ordinal
496 loss. *Computerized Medical Imaging and Graphics*.
- 497 Cheplygina, V., de Bruijne, M., Pluim, J. P., 2019. Not-so-supervised: a survey of
498 semi-supervised, multi-instance, and transfer learning in medical image analysis.
499 *Medical Image Analysis*.
- 500 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. Imagenet: A
501 large-scale hierarchical image database. In: *Proceedings of the IEEE Conference*
502 *on Computer Vision and Pattern Recognition*. pp. 248–255.
- 503 Galea, M. H., Blamey, R. W., Elston, C. E., Ellis, I. O., 1992. The nottingham
504 prognostic index in primary breast cancer. *Breast cancer research and treatment*
505 22 (3), 207–219.
- 506 Ghosh, A., Kumar, H., Sastry, P., 2017. Robust loss functions under label noise
507 for deep neural networks. In: *Proceedings of the AAAI Conference on Artificial*
508 *Intelligence*. pp. 1919–1925.
- 509 Gurcan, M. N., Boucheron, L., Can, A., Madabhushi, A., Rajpoot, N., Yener, B.,
510 2009. Histopathological image analysis: A review. *IEEE Reviews in Biomedical*
511 *Engineering* 2, 147.
- 512 Haeusser, P., Mordvintsev, A., Cremers, D., 2017. Learning by association—a versa-
513 tile semi-supervised training method for neural networks. In: *IEEE Conference*
514 *on Computer Vision and Pattern Recognition*. Vol. 3. p. 6.

- 515 Kamnitsas, K., Castro, D. C., Folgoc, L. L., Walker, I., Tanno, R., Rueckert, D.,
516 Glocker, B., Criminisi, A., Nori, A., 2018. Semi-supervised learning via compact
517 latent space clustering. arXiv preprint arXiv:1806.02679.
- 518 Kapil, A., Meier, A., Zuraw, A., Steele, K., Rebelatto, M., Schmidt, G., Brieu, N.,
519 2018. Deep semi supervised generative learning for automated pd-l1 tumor cell
520 scoring on nslc tissue needle biopsies. arXiv preprint arXiv:1806.11036.
- 521 Kingma, D. P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv
522 preprint arXiv:1412.6980.
- 523 Kipf, T. N., Welling, M., 2016. Semi-supervised classification with graph convolu-
524 tional networks. arXiv preprint arXiv:1609.02907.
- 525 Kraus, O. Z., Ba, J. L., Frey, B. J., 2016. Classifying and segmenting microscopy
526 images with deep multiple instance learning. *Bioinformatics* 32 (12), i52–i59.
- 527 Krizhevsky, A., Sutskever, I., Hinton, G. E., 2012. Imagenet classification with
528 deep convolutional neural networks. In: *Proceedings of Advances in Neural
529 Information Processing Systems*. pp. 1097–1105.
- 530 Kumar, A., Sattigeri, P., Fletcher, T., 2017. Semi-supervised learning with gans:
531 manifold invariance with improved inference. In: *Proceedings of Advances in
532 Neural Information Processing Systems*. pp. 5534–5544.
- 533 Laine, S., Aila, T., 2016. Temporal ensembling for semi-supervised learning. arXiv
534 preprint arXiv:1610.02242.
- 535 Lee, D.-H., 2013. Pseudo-label: The simple and efficient semi-supervised learn-
536 ing method for deep neural networks. In: *ICML Workshop on Challenges in
537 Representation Learning*. Vol. 3. p. 2.

- 538 Li, Y., Yang, J., Song, Y., Cao, L., Luo, J., Li, L.-J., 2017. Learning from noisy labels with distillation. In: Proceedings of International Conference on Computer
539 Vision. pp. 1928–1936.
- 541 Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian,
542 M., Van Der Laak, J. A., Van Ginneken, B., Sánchez, C. I., 2017. A survey on
543 deep learning in medical image analysis. *Medical Image Analysis* 42, 60–88.
- 544 Luo, Y., Guan, T., Yu, J., Liu, P., Yang, Y., 2018. Every node counts: Self-
545 ensembling graph convolutional networks for semi-supervised learning. arXiv
546 preprint arXiv:1809.09925.
- 547 Luo, Y., Zhu, J., Li, M., Ren, Y., Zhang, B., 2017. Smooth neighbors on teacher
548 graphs for semi-supervised learning. arXiv preprint arXiv:1711.00258.
- 549 Maaten, L. v. d., Hinton, G., 2008. Visualizing data using t-sne. *Journal of machine
550 learning research* 9 (Nov), 2579–2605.
- 551 Miyato, T., Maeda, S.-i., Ishii, S., Koyama, M., 2018. Virtual adversarial train-
552 ing: a regularization method for supervised and semi-supervised learning. *IEEE
553 Transactions on Pattern Analysis and Machine Intelligence*.
- 554 Mnih, V., Hinton, G. E., 2012. Learning to label aerial images from noisy data. In:
555 Proceedings of International Conference on Machine Learning). pp. 567–574.
- 556 Natarajan, N., Dhillon, I. S., Ravikumar, P. K., Tewari, A., 2013. Learning with
557 noisy labels. In: Proceedings of Advances in Neural Information Processing
558 Systems. pp. 1196–1204.
- 559 Odena, A., 2016. Semi-supervised learning with generative adversarial networks.
560 arXiv preprint arXiv:1606.01583.

- 561 Patrini, G., Rozza, A., Menon, A. K., Nock, R., Qu, L., 2017. Making deep neural
562 networks robust to label noise: A loss correction approach. In: Proceedings of
563 the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2233–
564 2241.
- 565 Peikari, M., Salama, S., Nofech-Mozes, S., Martel, A. L., 2018. A cluster-then-label
566 semi-supervised learning approach for pathology image classification. Scientific
567 Reports 8 (1), 7193.
- 568 Quellec, G., Cazuguel, G., Cochener, B., Lamard, M., 2017. Multiple-instance
569 learning for medical image and video analysis. IEEE reviews in biomedical en-
570 gineering 10, 213–234.
- 571 Rasmus, A., Berglund, M., Honkala, M., Valpola, H., Raiko, T., 2015. Semi-
572 supervised learning with ladder networks. In: Proceedings of Advances in Neural
573 Information Processing Systems. pp. 3546–3554.
- 574 Reed, S., Lee, H., Anguelov, D., Szegedy, C., Erhan, D., Rabinovich, A.,
575 2014. Training deep neural networks on noisy labels with bootstrapping. arXiv
576 preprint arXiv:1412.6596.
- 577 Ren, M., Zeng, W., Yang, B., Urtasun, R., 2018. Learning to reweight examples
578 for robust deep learning. arXiv preprint arXiv:1803.09050.
- 579 Ruder, S., Plank, B., 2018. Strong baselines for neural semi-supervised learning
580 under domain shift. arXiv preprint arXiv:1804.09530.
- 581 Sapkota, M., Shi, X., Xing, F., Yang, L., 2018. Deep convolutional hashing for
582 low-dimensional binary embedding of histopathological images. IEEE journal of
583 biomedical and health informatics 23 (2), 805–816.

- 584 Shen, D., Wu, G., Suk, H.-I., 2017. Deep learning in medical image analysis.
585 Annual Review of Biomedical Engineering 19, 221–248.
- 586 Shi, X., Sapkota, M., Xing, F., Liu, F., Cui, L., Yang, L., 2018. Pairwise based deep
587 ranking hashing for histopathology image classification and retrieval. Pattern
588 Recognition 81, 14–22.
- 589 Shi, X., Xing, F., Xu, K., Xie, Y., Su, H., Yang, L., 2017. Supervised graph hashing
590 for histopathology image retrieval and classification. Medical image analysis 42,
591 117–128.
- 592 Sparks, R., Madabhushi, A., 2016. Out-of-sample extrapolation utilizing semi-
593 supervised manifold learning (ose-ssl): Content based image retrieval for
594 histopathology images. Scientific Reports 6, 27306.
- 595 Su, H., Shi, X., Cai, J., Yang, L., 2019. Local and global consistency regular-
596 ized mean teacher for semi-supervised nuclei classification. In: International
597 Conference on Medical Image Computing and Computer-Assisted Intervention.
598 Springer, pp. 559–567.
- 599 Sukhbaatar, S., Fergus, R., 2014. Learning from noisy labels with deep neural
600 networks. arXiv preprint arXiv:1406.2080 2 (3), 4.
- 601 Tang, X., Guo, F., Shen, J., Du, T., 2018. Facial landmark detection by semi-
602 supervised deep learning. Neurocomputing 297, 22–32.
- 603 Tarvainen, A., Valpola, H., 2017. Mean teachers are better role models: Weight-
604 averaged consistency targets improve semi-supervised deep learning results. In:
605 Proceedings of Advances in Neural Information Processing Systems. pp. 1195–
606 1204.

- 607 Vahdat, A., 2017. Toward robustness against label noise in training deep discrim-
608 inative neural networks. In: Proceedings of Advances in Neural Information
609 Processing Systems. pp. 5596–5605.
- 610 Veit, A., Alldrin, N., Chechik, G., Krasin, I., Gupta, A., Belongie, S. J., 2017.
611 Learning from noisy large-scale datasets with minimal supervision. In: Pro-
612 ceedings of the IEEE Conference on Computer Vision and Pattern Recognition.
613 pp. 6575–6583.
- 614 Weston, J., Ratle, F., Mobahi, H., Collobert, R., 2012. Deep learning via semi-
615 supervised embedding. In: Neural Networks: Tricks of the Trade. Springer, pp.
616 639–655.
- 617 Xiao, T., Xia, T., Yang, Y., Huang, C., Wang, X., 2015. Learning from massive
618 noisy labeled data for image classification. In: Proceedings of the IEEE Confer-
619 ence on Computer Vision and Pattern Recognition. pp. 2691–2699.
- 620 Xing, F., Xie, Y., Su, H., Liu, F., Yang, L., 2017. Deep learning in microscopy
621 image analysis: A survey. IEEE transactions on neural networks and learning
622 systems 29 (10), 4550–4568.
- 623 Xu, K., Chen, H., Liu, S., Chen, P.-Y., Weng, T.-W., Hong, M., Lin, X., 2019.
624 Topology attack and defense for graph neural networks: An optimization per-
625 spective. arXiv preprint arXiv:1906.04214.
- 626 Xu, K., Liu, S., Zhao, P., Chen, P.-Y., Zhang, H., Fan, Q., Erdogmus, D., Wang,
627 Y., Lin, X., 2018. Structured adversarial attack: Towards general implementa-
628 tion and better interpretability. arXiv preprint arXiv:1808.01664.
- 629 Xu, Y., Zhu, J.-Y., Eric, I., Chang, C., Lai, M., Tu, Z., 2014. Weakly super-

- 630 vised histopathology cancer image segmentation and classification. *Medical im-*
631 *age analysis* 18 (3), 591–604.
- 632 Yarowsky, D., 1995. Unsupervised word sense disambiguation rivaling supervised
633 methods. In: *Proceedings of annual meeting on Association for Computational*
634 *Linguistics*. Association for Computational Linguistics, pp. 189–196.
- 635 Zhang, J., Peng, Y., 2017. Ssdh: semi-supervised deep hashing for large scale image
636 retrieval. *IEEE Transactions on Circuits and Systems for Video Technology*.
- 637 Zhang, X., Liu, W., Dundar, M., Badve, S., Zhang, S., 2015. Towards large-
638 scale histopathological image analysis: Hashing-based image retrieval. *IEEE*
639 *Transactions on Medical Imaging* 34 (2), 496–506.
- 640 Zhou, Y., Goldman, S., 2004. Democratic co-learning. In: *Proceedings of Interna-*
641 *tional Conference on Tools with Artificial Intelligence*. pp. 594–602.