



UNIVERSIDAD
NACIONAL
DE COLOMBIA

1 **Medical image segmentation in a multiple**
2 **labelers context: Application to the study of**
3 **histopathology**

4 **Brandon Lotero Londoño**

5 Universidad Nacional de Colombia
6 Faculty of Engineering and Architecture
7 Department of Electric, Electronic and Computing Engineering
8 Manizales, Colombia
9 2023

10 **Medical image segmentation in a multiple**
11 **labelers context: Application to the study of**
12 **histopathology**

13 **Brandon Lotero Londoño**

14 Dissertation submitted as a partial requirement to receive the grade of:
15 **Master in Engineering - Industrial Automation**

16 Advisor:

17 Prof. Andrés Marino Álvarez-Meza, Ph.D.

18 Co-advisor:

19 Prof. Germán Castellanos-Domínguez, Ph.D.

20 Academic research group:

21 Signal Processing and Recognition Group - SPRG

22 Universidad Nacional de Colombia

23 Faculty of Engineering and Architecture

24 Department of Electric, Electronic and Computing Engineering

25 Manizales, Colombia

26 2025

27 **Segmentación de imágenes médicas en un**
28 **contexto de múltiples anotadores:**
29 **Aplicación al estudio de histopatologías**

30 **Brandon Lotero Londoño**

31 Disertación presentada como requisito parcial para recibir el título de:
32 **Magíster en Ingeniería - Automatización Industrial**

33 Director:

34 Prof. Andrés Marino Álvarez-Meza, Ph.D.

35 Codirector:

36 Prof. Germán Castellanos-Domínguez, Ph.D.

37 Grupo de investigación:

38 Grupo de Control y Procesamiento Digital de Señales - GCPDS

39 Universidad Nacional de Colombia

40 Facultad de Ingeniería y Arquitectura

41 Departamento de Ingeniería Eléctrica, Electrónica y Computación

42 Manizales, Colombia

43 2023

ACKNOWLEDGEMENTS

45 PENDING

ABSTRACT

49 PENDING

50 **Keywords:** PENDING

53 PENDIENTE

54 **Palabras clave:** PENDIENTE

56 **Contents**

57	Acknowledgements	vii
58	Abstract	ix
59	Resumen	xi
60	Contents	xiv
61	List of figures	xv
62	List of tables	xvii
63	Abbreviations	xix
64	1 Introduction	1
65	1.1 Motivation	1
66	1.2 Problem Statement	6
67	1.2.1 Variability in Expertise Levels	8
68	1.2.2 Technical Constraints and Image Quality	9
69	1.2.3 Research Question	9
70	1.3 Literature review	11
71	1.3.1 Facing annotation variability in medical images	12
72	1.3.2 Facing noisy annotations and low-quality data	19
73	1.4 Aims	21
74	1.4.1 General Aim	22
75	1.4.2 Specific Aims	22

76	1.5	Outline and Contributions	23
77	2	Conceptual preliminaries	25
78	2.1	Modern concept of digital image	25
79	2.1.1	Types of digital images	25
80	2.1.2	Mathematical representations	26
81	2.1.3	Digital histopathological images	28
82	2.2	Deep learning fundamentals	28
83	2.3	Datasets and data sources	28
84	3	Chained Gaussian Processes	29
85	3.1	Gaussian processes	29
86	3.2	Chained Gaussian processes	29
87	4	Truncated Generalized Cross Entropy for segmentation	31
88	4.1	Loss functions for multiple annotators	31
89	4.1.1	Generalized Cross Entropy	32
90	4.1.2	Extension to Multiple Annotators	34
91	4.1.3	Reliability Maps and Truncated GCE	34
92	4.2	Proposed Model	36
93	4.2.1	Backbone Architecture	36
94	4.2.2	UNET Architecture	36
95	4.2.3	Reliability Map Branch	37
96	4.2.4	Integration with $TGCE_{SS}$ Loss	37
97	4.2.5	Training Process	38
98	4.3	Experiments	38
99	4.3.1	Dataset	38
100	4.3.2	Metrics	38
101		Bibliography	40

LIST OF FIGURES

103	1-1	Estimation of the tasks and medical image types based on recent	
104		literature review (count of referenced terms).	3
105	1-2	AI and machine learning in medical imaging brief timeline.	4
106	1-3	Example of a histopathological image segmented by multiple	
107		annotators, illustrating variations in label assignment.	7
108	1-4	Summary diagram for problem Statement	10
109	1-5	Proposed framework for the approach in [López-Pérez et al., 2024]. .	17
110	4-1	Solution Architecture (mockup)	39

- 113 **CAD** Computer-Aided Diagnosis 2, 5, 6
114 **CCGP** Correlated Chained Gaussian Processes 18
115 **CCGPMA** Correlated Chained Gaussian Processes for Multiple Annotators 17, 19
116 **CE** Cross Entropy 32
117 **CGP** Chained Gaussian Processes 18
118 **CNN** Convolutional Neural Networks 3, 14, 22, 23
119 **CT** Computed Tomography 11

120 **ELBO** Evidence Lower Bound 18
121 **GCE** Generalized Cross Entropy 32
122 **GCECDL** Generalized Cross-Entropy-based Chained Deep Learning 20, 21

123 **ISS** Image Semantic segmentation 2, 3, 6, 11, 13, 21-23
124 **LF** Latent Function 18

125 **MAE** Mean Absolute Error 32, 33
126 **MITs** Medical Imaging Techniques 1
127 **ML** Machine Learning 11, 21
128 **MV** Majority Voting 11, 12

129 **OCR** Optical Character Recognition 11
130 **PET** Positron Emission Tomography 14

131 **ROI** Region of Interest 2, 6

132 **SLFM** Semi-Parametric Latent Factor Model 18
133 **SS** Semantic segmentation 3
134 **STAPLE** Simultaneous Truth and Performance Level Estimation 12-14
135 **WSI** Whole Slide Imaging 1, 5, 6, 8, 15

136

137

CHAPTER

138

ONE

139

140

INTRODUCTION

141

1.1 Motivation

142 Since Roentgen's discovery of X-rays in 1895, medical imaging has advanced
143 significantly, with modalities like radionuclide imaging, ultrasound, CT, MRI, and
144 digital radiography emerging over the past 50 years. Modern imaging extends
145 beyond image production to include processing, display, storage, transmission and
146 analysis. [Zhou et al., 2021]. Other Medical Imaging Techniques (MITs) have arose
147 during the last decades, some of them implying only the examination of certain
148 pieces or tissues instead of complete patients, like histopathological images, which
149 are images of tissue samples obtained from biopsies or surgical resections and are
150 widely used for the diagnosis of diseases like cancer through Whole Slide Imaging
151 (WSI) scanners [Rashmi et al., 2021].

152 Along with the advances in technologies for medical images acquisition,
153 computational technologies on pattern recognition and artificial intelligence have

also emerged, allowing the development of **Computer-Aided Diagnosis (CAD)** systems based on machine learning algorithms. These systems aim to assist physicians in the diagnosis and treatment of diseases, by providing a second opinion or by automating the analysis of medical images. [Panayides et al., 2020]. One of the most used tasks in which machine learning technologies is being used in the universe of medical images is **Image Semantic segmentation (ISS)**, which consists of assigning a label to each pixel in an image according to the object it belongs to. This task is crucial for the development of **CAD** systems, as it allows the identification of **Region of Interest (ROI)** in the images, which can be used to detect and classify diseases [Azad et al., 2024].

The application of Machine Learning in medical imaging has grown significantly, with key tasks including classification, segmentation, anomaly detection, super-resolution, image registration, and synthetic image generation [Brito-Pacheco et al., 2025]. Among imaging modalities, X-rays and CT scans are widely used for classification and anomaly detection, especially in pulmonary and oncological applications. MRI and ultrasound play a crucial role in segmentation and resolution enhancement, while PET/SPECT imaging is essential for anomaly detection in oncology and neurodegenerative diseases [Brito-Pacheco et al., 2025]. Histopathology is rapidly gaining prominence, particularly in segmentation and feature extraction, where AI-driven techniques aid in automated cancer diagnosis and tissue structure analysis. The integration of Deep Learning in histological image processing is revolutionizing pathology, enabling more precise and efficient diagnostics. A brief comparison of the tasks and medical image types based on recent literature review, can be seen in Figure 1-1. [Yu et al., 2025], [Brito-Pacheco et al., 2025], [Ryou et al., 2025], [Hu et al., 2025], [Elhaminia et al., 2025]

For solving the different requirements of tasks in medical images, a variety of computational techniques have been developed [Zhou et al., 2021]. Initially, these needs were covered with simple morphological filters, which implied no training process or elaborated optimization. However, as the complexity of the tasks increased, the need for more sophisticated techniques arose, leading to the

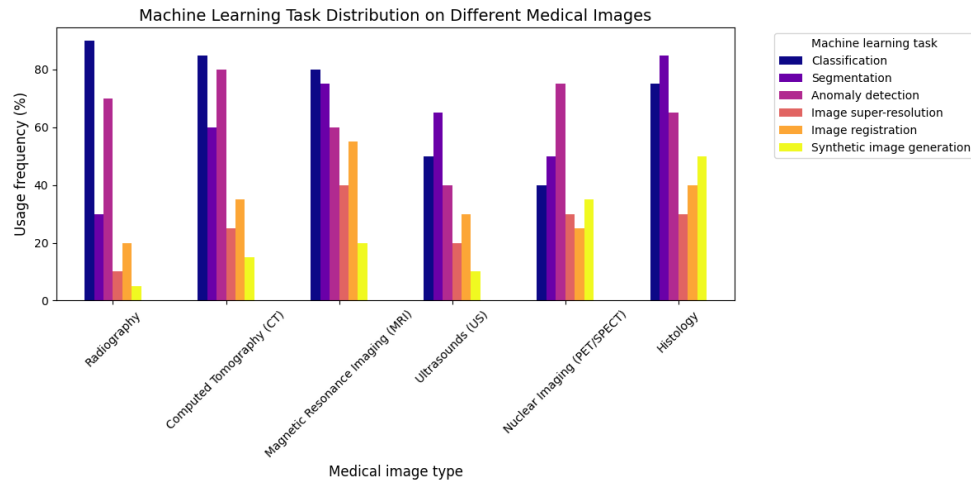


Figure 1-1 Estimation of the tasks and medical image types based on recent literature review (count of referenced terms).

184 application of advanced statistical tools and machine learning algorithms like
 185 Support Vector Machines, Decision Trees, and SGD Neural Networks [Avanzo
 186 et al., 2024]. The coevolution of advances in medical image acquisition,
 187 computational power (i.e. Moore's law) and statistical/mathematical techniques
 188 have led to a convergence for merging state of the art algorithms with medical
 189 imaging [Shalf, 2020]. Figure 1-2 shows a brief timeline of coevolution between
 190 some conspicuous advances in computational pattern recognition and its medical
 191 applications in different scopes (besides medical imaging) [Avanzo et al., 2024].

192 Convolutional Neural Networks (CNN) have been widely used in Semantic
 193 segmentation (SS) tasks, as they have outperformed traditional machine learning
 194 algorithms in this task for both medical and non medical images [Xu et al., 2024]
 195 [Sarvamangala and Kulkarni, 2022]. However, most CNN architectures are deep,
 196 which imply a necessity of a large amount of data to train them. This introduces a
 197 problem since both the acquisition and annotation of medical images are
 198 expensive and time-consuming processes. This is especially true for ISS tasks, as
 199 they require pixel-level annotations, which is taxing in terms of cost, time and
 200 logistics involved [Bhalgat et al., 2018]. Other fashions face this problem through
 201 less expensive annotation strategies like bounding boxes or anatomical landmarks

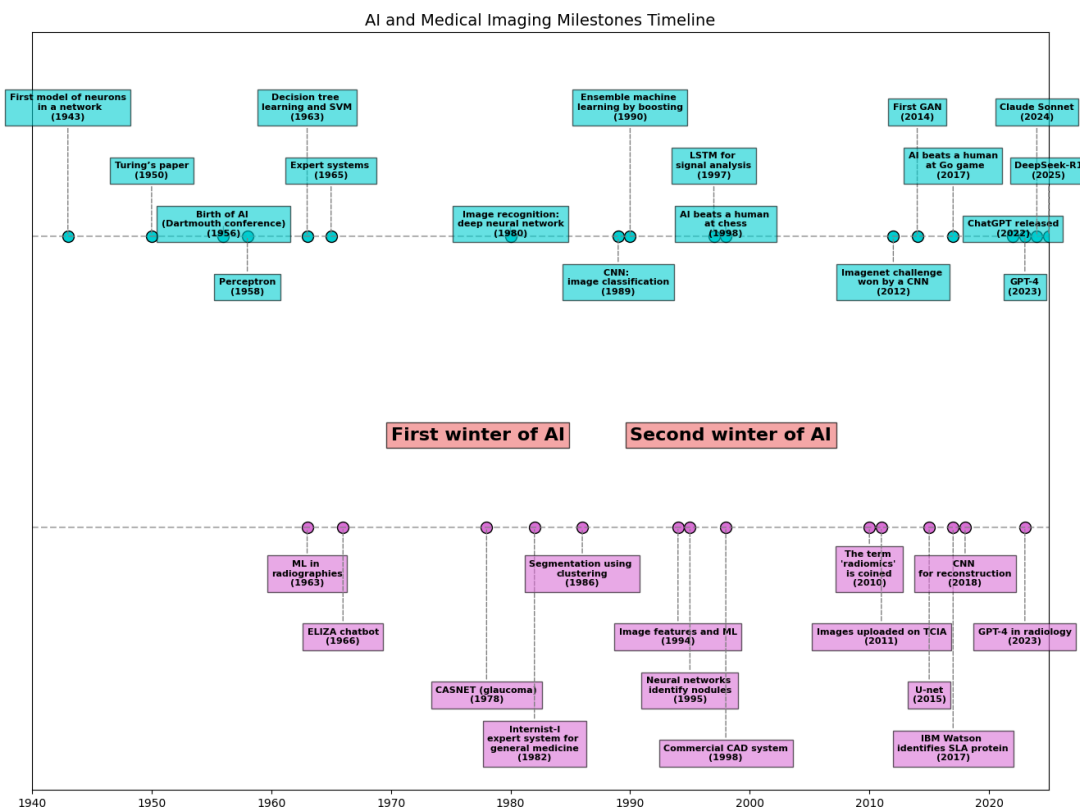


Figure 1-2 AI and machine learning in medical imaging brief timeline.

202 for being used in a semi-supervised strategy [Shah et al., 2018].

203 Many medical images datasets however, contain a high variability in class sizes
204 and variations in colors, which is specially noticeable in histopathological images
205 because of the usage of different staining and other factors which can affect the
206 color of the images. This variability can lead to a significant loss of efficiency of
207 machine learning models when using a mixed supervision strategy, as the model
208 can be biased towards the most common classes or colors in the dataset [Shah
209 et al., 2018].

210 This is where other solutions arise to tackle the problem of the weak image
211 annotation while maintaining low costs. One of these solutions is crowdsourcing
212 strategy, which consists of having multiple annotators labeling the same image,
213 and then combining the labels to obtain a consensus label [Lu et al., 2023]. This
214 strategy can lead to a labeling cost reduction when different levels of expertise are
215 combined, since the crowd may be composed of both experts and laymen, being
216 the latter less expensive to hire [López-Pérez et al., 2023].

217 Recently, diagnosis, prognosis and treatment of cancer have heavily relied on
218 histopathology, where tissue samples are obtained through biopsies or surgical
219 resections and critical information that helps pathologists determine the presence
220 and severity of malignancies [López-Pérez et al., 2024]. The segmentation of
221 histopathological images enables precise identification of structures such as
222 nuclei, glands, and tumors, which are essential for assessing disease progression
223 and treatment response [Rashmi et al., 2021]. Accurate segmentation is
224 particularly crucial in digital pathology, where whole-slide images (WSI) are
225 analyzed using AI-powered CAD systems to support clinical decision-making
226 [López-Pérez et al., 2024].

227 A major challenge in histopathological image segmentation arises from the
228 variability in annotations provided by different pathologists. Unlike natural
229 images, where object boundaries are often well-defined, histological structures
230 may have ambiguous borders, leading to inconsistencies among annotators

[López-Pérez et al., 2023]. Because of this, crowdsourcing labeling is one of the most popular approaches, as illustrated in Figure 1-3, an example of how histopathological images are segmented by multiple experts, showing some variations in label assignment ¹. These discrepancies highlight the need for models that can handle annotation uncertainty effectively. Leveraging crowdsourcing strategies and machine learning techniques that infer annotator reliability can enhance segmentation performance while reducing costs.

1.2 Problem Statement

Throughout the development of medical technology and CAD, the task of ISS has become a crucial step in delivering precise diagnosis and treatment planning [Giri and Bhatia, 2024]. Particularly, in the area of histopathological studies, the usage of Whole Slide Images (WSI) is rather common since this method delivers high quality imaging and allows for the diagnosis of diseases like cancer [Lin et al., 2024].

ISS task consists of assigning a label to each pixel in an image according to the object it belongs to. Accurate segmentation is essential for the development of CAD systems, as it allows the identification of regions of interest (ROI) in the images, which can be used to detect and classify diseases and hence, treatment planning [Sarvamangala and Kulkarni, 2022]. However, modern computational solutions for ISS tasks involve the use of deep learning, which mostly rely large amounts of labeled data to train the models on supervised learning techniques. This means that the model is trained on a dataset with ground-truth labels, which are assumed to be correct and consistent across all samples. In practice, this assumption is often violated due to the high technical complexity of labeling these segments ².

¹obtained from a real world Triple Negative Breast Cancer (TNBC) dataset published in [López-Pérez et al., 2023]

²compared to a more trivial task like image classification on ordinary an well known classes like MNIST

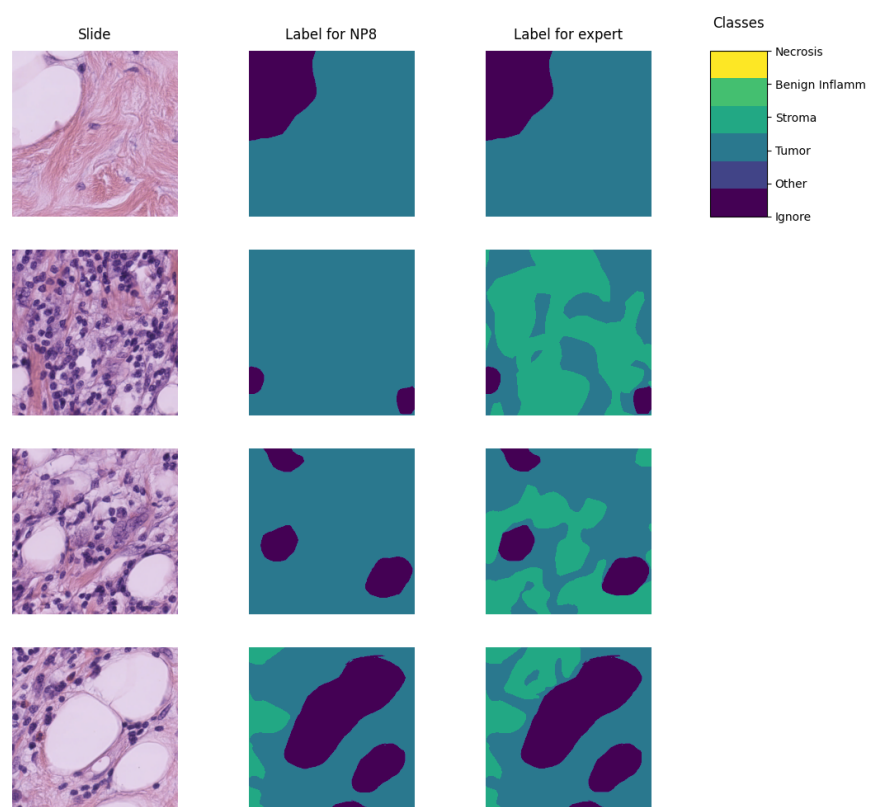


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

The process of labeling medical images is often managed with the help of specialized software tools that allow the annotators to draw the regions, delivering an standard format for the labeled masks [Habis, 2024]. Despite the help of these tools, the labeling process in WSI can have high costs, as it requires long hours of work from specialized personnel. Because of cost constraints in many medical institutions, the labeling processes is often done by multiple labelers with varying levels of expertise, equalizing the cost of the labeling process. However, this strategy can lead to inconsistent labels, as the consensus between the labelers may not be exact due to the diversity in depth of knowledge and experience of the labelers [Xu et al., 2024]. These inconsistencies are mostly represented in the subsections 1.2.1 and 1.2.2.

1.2.1 Variability in Expertise Levels

One of the primary sources of inter-observer variability in medical image segmentation is the difference in expertise levels among annotators [López-Pérez et al., 2023]. Experienced radiologists and pathologists tend to produce highly precise annotations, whereas novice labelers may introduce systematic biases due to their limited familiarity with subtle image features. Studies have demonstrated that annotation accuracy tends to improve with experience, yet medical institutions often rely on a mix of annotators to manage costs and workload distribution [Lu et al., 2023].

The training background of annotators and institutional guidelines play a crucial role in shaping labeling practices. Different medical schools and hospitals may adopt distinct segmentation protocols, leading to inconsistencies when datasets are combined from multiple sources [López-Pérez et al., 2023]. For example, some institutions may emphasize conservative delineation of tumor boundaries, while others adopt a more inclusive approach. Such variations contribute to systematic biases in medical image datasets [Banerjee et al., 2025].

282 Medical images frequently contain structures with ambiguous boundaries, making
283 segmentation inherently subjective. For instance, tumor margins in
284 histopathological slides may not have well-defined edges, leading to variations in
285 how different annotators delineate the regions of interest [Carmo et al., 2025].
286 These discrepancies arise not only from technical expertise but also from
287 differences in perception and interpretation.

288 1.2.2 Technical Constraints and Image Quality

289 Technical constraints in medical imaging, such as resolution differences, noise
290 levels, and contrast variations, can significantly impact segmentation accuracy.
291 Lower-resolution images may obscure fine structures, leading to inconsistencies in
292 boundary delineation [Zhou et al., 2024].

293 When combined with long sessions, bad images might also increase the cognitive
294 load of the annotators, leading to fatigue and reduced precision in labeling [Kim
295 et al., 2024]. This is particularly relevant in histopathological studies, where the
296 staining process and tissue preparation can introduce color variations and artifacts
297 that affect image quality, even if the same scanning equipment is used [Karthikeyan
298 et al., 2023].

299 1.2.3 Research Question

300 Given the challenges posed by inconsistent labels in medical image segmentation,
301 this work aims to address the following research question:

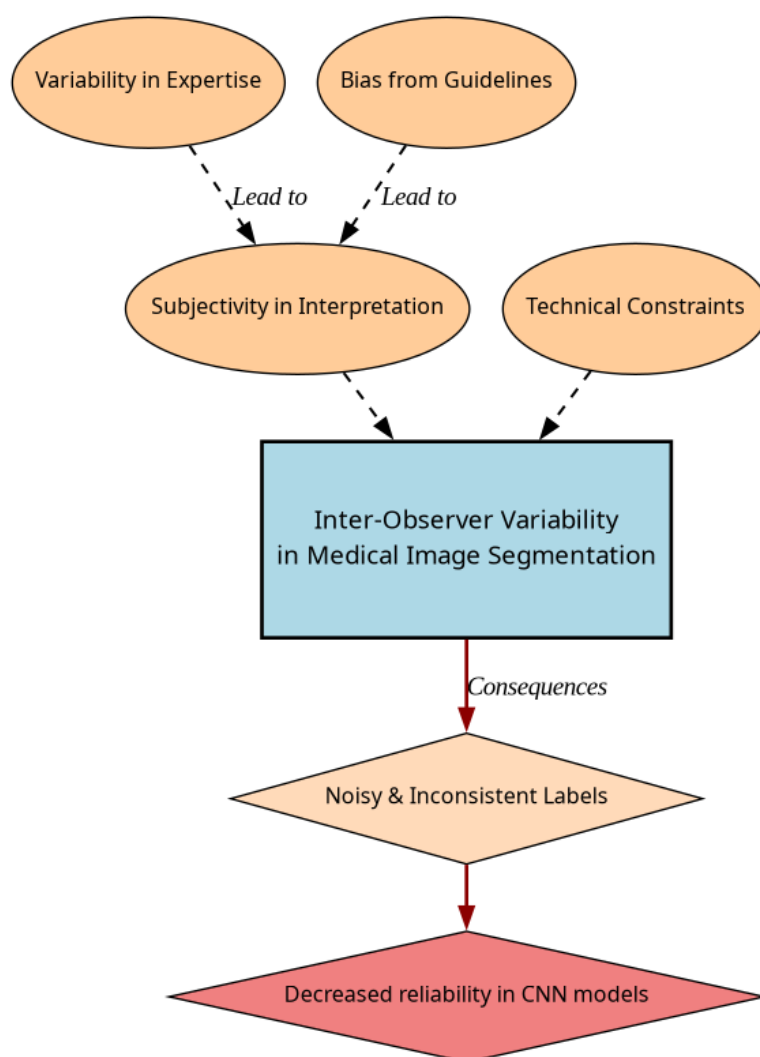


Figure 1-4 Summary diagram for problem Statement

Research Question

How can we develop a learning approach for ISS tasks in medical images that can adapt to inconsistent labels without requiring explicit supervision of labeler performance, while addressing challenges related to variability in expertise levels and technical constraints, and maintaining interpretability, generalization, and computational efficiency?

302

1.3 Literature review

303

Certainly, in general Machine Learning (ML) classification tasks³ where multiple annotators are involved, Majority Voting (MV) is by far the simplest possible approach to implement. This concept was born multiple times and divergently in multiple fields, but it was described as relevant for ML and pattern recognition labeling for classification in [Lam and Suen, 1997], in which the approach is exposed as simple, yet powerful. The authors describe the MV as a method that can be used to improve the accuracy of classification tasks by combining the labels of multiple annotators. The method is based on the assumption that the majority vote of the annotators is more likely to be correct than the vote of a single annotator. The authors also describe the method as a straightforward way to improve the accuracy of classification tasks without the need for complex algorithms or additional data. The authors also prove this method to deliver very similar results to more complicated approaches (Bayesian, logistic regression, fuzzy integral, and neural network) in the particular task of Optical Character Recognition (OCR). Despite its simplicity, modern solutions for delivering accurate medical image segmentation models still rely on Majority Voting at some stage, like [Elnakib et al., 2020], which uses a majority voting strategy for delivering a final output based on the labels of multiple models (VGG16-Segnet, Resnet-18 and Alexnet) in Computed Tomography (CT) images for Liver Tumor Segmentation, or

322

³In this work, image segmentation is considered as a particular case of classification in which target classes are assigned pixel-wise.

[López-Pérez et al., 2023], which uses MV for combining noisy annotations as an additional annotator to be included in the deep learning solution. Majority voting as a technique for setting a pseudo ground truth label is a powerful approach for its simplicity in many use cases in which the target to be labeled is not tied to an expertise related task, otherwise, the assumption of equal expertise among the labelers can be a source of bias in the final label, which is not desirable in the case of highly technical annotations like medical images. In subsection 1.3.1, we will be reviewing literature which no longer assumes the naive approach of equal expertise among labelers and face the challenge of learning from inconsistent labels.

1.3.1 Facing annotation variability in medical images

Learning from crowds approaches in general face the challenge of not having a ground truth label and hence, an intrinsic difficulty in measuring the real reliability of the labelers annotations. Some approaches assume beforehand a certain level of expertise for each labeler based on experience as an input, like in [TIAN and Zhu, 2015], which introduce the concept of max margin majority voting, using the reliability vector as weights for the weights for the binary and multiclass classifier. The crowdsourcing margin is the minimal difference between the aggregated score of the potential true label and the scores for other alternative labels. Accordingly, the annotators' reliability is estimated as generating the largest margin between the potential true labels and other alternatives. The problem introduced in this approach is assuming an stationary reliability per expert across the whole input space, which is imprecise since annotators performance may change between different tasks or even between different regions of the same image.

STAPLE Mechanism

The Simultaneous Truth and Performance Level Estimation (STAPLE) algorithm, introduced in [Warfield et al., 2004] is a probabilistic framework that estimates a

hidden true segmentation from multiple segmentations provided by different raters. It also estimates the reliability of each rater by computing their sensitivity and specificity.

The **STAPLE** algorithm's goal is to maximize the log likelihood function:

$$(\mathbf{p}, \mathbf{q}) = \arg \max_{\mathbf{p}, \mathbf{q}} \ln f(\mathbf{D}, \mathbf{T} \mid \mathbf{p}, \mathbf{q}). \quad (1-1)$$

Where \mathbf{D} is the set of segmentations provided by the raters, \mathbf{T} is the hidden true segmentation, p is the sensitivity and q is the specificity of the raters.

This is achieved by using the Expectation-Maximization algorithm to maximize the log likelihood function in equation, which is done iteratively with step computations:

$$\begin{aligned} (p_j^{(k)}, q_j^{(k)}) = \arg \max_{p_j, q_j} & \sum_{i: D_{ij}=1} W_i^{(k-1)} \ln p_j \\ & + \sum_{i: D_{ij}=1} \left(1 - W_i^{(k-1)}\right) \ln(1 - q_j) \\ & + \sum_{i: D_{ij}=0} W_i^{(k-1)} \ln(1 - p_j) \\ & + \sum_{i: D_{ij}=0} \left(1 - W_i^{(k-1)}\right) \ln q_j. \end{aligned} \quad (1-2)$$

The capacity of STAPLE to accurately estimate the true segmentation, even in the presence of a majority of raters generating correlated errors, was demonstrated, which makes it theoretically a strong choice for setting a ground-truth in binary or multiclass medical **ISS** tasks.

The popularity and performance of **STAPLE** has led to its usage in modern applications medical image, 3d spatial images due to its assumption of decision

space being based on voxel-wise decisions, like the authors in [Grefve et al., 2024] which applied the algorithm on Positron Emission Tomography (PET) images. Other authors still rely heavily on STAPLE for setting a ground truth consensus for histopathological images, like [Qiu et al., 2022].

However, the STAPLE algorithm has some limitations. It assumes independent rater errors, which may not hold in practice, leading to biased estimates. STAPLE is also sensitive to low-quality annotations, potentially degrading final segmentations if the weights are not initialized correctly. The algorithm tends to over-smooth results, blurring fine details, and struggles with multi-class segmentation. Computationally, it is expensive due to its iterative EM approach. Additionally, STAPLE cannot correct systematic biases in annotations and depends on initial estimates, impacting accuracy. Lastly, the estimated performance levels lack interpretability, making it difficult to assess annotator reliability effectively.

Finally, this work contemplates STAPLE as useful for label aggregation, hence being a good support for other methods, but not that useful for providing annotations of structures on new and unlabeled images.

U-shaped CNNs

Since the introduction of U-Net [Ronneberger et al., 2015] in 2015 for biomedical image segmentation, U-shaped CNNs have become a prevalent architecture in medical image segmentation tasks. The U-Net’s success stems from its ability to capture both global and local information through its contracting and expanding paths, making it particularly effective for complex and heterogeneous structures, even with limited annotated data. This architecture has been successfully applied to various medical image segmentation tasks, including organ segmentation, tumor segmentation, and brain structure segmentation.

The U-Net architecture consists of a symmetric encoder-decoder structure with skip connections. The encoder path progressively reduces spatial dimensions

while increasing feature channels through a series of convolutional and max-pooling layers, capturing high-level semantic information. The decoder path uses transposed convolutions to gradually recover spatial resolution while reducing feature channels. Skip connections between corresponding encoder and decoder layers preserve fine-grained details by concatenating high-resolution features from the encoder with upsampled features in the decoder, enabling precise localization of structures.

U-Net based approaches

In [López-Pérez et al., 2024] two networks are trained for delivering a final segmentation. One network is trained to estimate the annotators reliability and another one is trained to segment the image. The first network is a deep neural network that takes as input features of image and the labelers id encoded as one-hot and outputs a reliability map across the image feature space. This map is then used to weight the contribution of each annotator to the final segmentation. The second network is the U-Net used for segmentation.

In this approach, it is assumed that the images are labeled for at least one labeler and not all of them, which is closer to a real world scenario, in which it is common to have images with variability in the amount of annotations, per patch. Hence, the input data can be modeled as:

$$\mathcal{D} = (\mathbf{X}, \tilde{\mathbf{Y}}) = \{(\mathbf{x}_n, \tilde{\mathbf{y}}_n^r) : n = 1, \dots, N; r \in R_n\}, \quad (1-3)$$

Where every \mathbf{x}_n is an input patch from a ROI in one WSI, $\tilde{\mathbf{y}}_n$ is the noisy annotation from the r labeler, N is the number of patches in the dataset and $R_n \subset \{1, \dots, R\}$ is the set of labelers that annotated the image \mathbf{x}_n .

The authors then assume the annotator network to deliver a reliability map $\{\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})\}_{r \in R_n}$ with different dimensions:

- 416 • CR global: a single reliability vector per labeler with dimensions C which
417 represent global reliability of the labeler across all input space.
- 418 • CR image: a single reliability vector per image per labeler with dimensions C
419 which represent local reliability of the labeler across the image.
- 420 • CR pixel: a reliability matrix per image per labeler, with dimensions C which
421 represent local reliability of the labeler across all the pixels in the image.

422 These differences in dimensions are determined by the feature extraction space
423 from segmentation network which feed the input of the annotator network, which
424 the authors vary for experimentation purposes.

425 Being $\mathbf{p}_\theta(\mathbf{x}_n)$ the estimation of the latent (ground truth) segmentation delivered by
426 the segmentation UNet network, thus, the estimated segmentation probability
427 mask for each annotator is given by the product:

$$\mathbf{p}_{\theta,\phi}^{(r)}(\mathbf{x}_n) := \mathbf{A}_\phi^{(r)}(\mathbf{x}) \odot \mathbf{p}_\theta(\mathbf{x}_n), \quad (1-4)$$

428 where \odot is the element-wise product and ϕ and θ are the parameters of the
429 annotator network and the segmentation UNet network, respectively, being the
430 latter initialized with a ResNet34 backbone pre-trained on ImageNet.

431 The authors propose a loss function involving cross-entropy and a trace based
432 regularization on the reliability map, originally proposed in [Zhang et al., 2020]
433 which combined, looks like:

$$\mathcal{L}(\theta, \phi) := \sum_{n=1}^N \sum_{r=1}^R \mathbb{I}(\tilde{\mathbf{y}}_n^{(r)} \in R_n) \cdot \left[\text{CE} \left(\mathbf{A}_\phi^{(r)}(\mathbf{x}_n) \cdot \mathbf{p}_\theta(\mathbf{x}_n), \tilde{\mathbf{y}}_n^{(r)} \right) + \lambda \cdot \text{tr} \left(\mathbf{A}_\phi^{(r)}(\mathbf{x}_n) \right) \right] \quad (1-5)$$

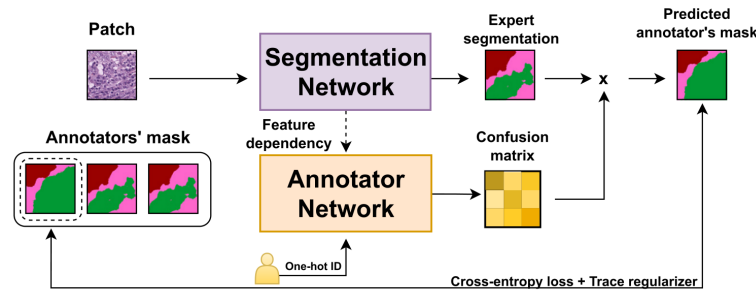


Figure 1-5 Proposed framework for the approach in [López-Pérez et al., 2024].

434 Being \mathbb{I} the indicator function, CE the cross-entropy loss, and λ the regularization
 435 parameter.

436 When evaluated on a Triple Negative Breast Cancer dataset, this approach achieves
 437 a Dice coefficient of 0.7827, outperforming STAPLE (0.7039) and matching expert-
 438 supervised performance (0.7723). The CR image reliability modeling proved most
 439 effective, as CR pixel, while potentially offering finer-grained reliability estimation,
 440 requires significantly more training data.

441 Despite the decent performance of the approach, solving the problem of multiple
 442 labelers with two networks can be overwhelming for the optimization process,
 443 requiring large amounts of annotated data to properly codify the annotators
 444 spatial reliabilities, which could be managed by a single model with an appropriate
 445 loss function.

446 Bayesian models

447 Bayesian approaches are a good choice for handling label noise and uncertainty in
 448 the labelers. In [Julián and Álvarez Meza Andrés Marino, 2023] the authors
 449 propose a novel approach from Gaussian Processes to model the relationship
 450 between the annotators' reliability and the input data, while also preserving the
 451 interdependencies among the annotators. This is achieved by introducing
 452 Correlated Chained Gaussian Processes for Multiple Annotators (CCGPMA), a

framework based on the well known **Chained Gaussian Processes (CGP)**. CGP on itself cannot consider inter-annotator dependencies, thus, the authors introduce the **Correlated Chained Gaussian Processes (CCGP)** to model correlations between the GP latent functions, which are supposed to be generated from a **Semi-Parametric Latent Factor Model (SLFM)**:

$$f_j(\mathbf{x}_n) = \sum_{q=1}^Q w_{j,q} \mu_q(\mathbf{x}_n), \quad (1-6)$$

where $f_j : \mathcal{X} \rightarrow \mathbb{R}$ is a **Latent Function (LF)**, $\mu_q(\cdot) \sim \mathcal{GP}(0, k_q(\cdot, \cdot))$ with $k_q : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ being a kernel function, and $w_{j,q} \in \mathbb{R}$ is a combination coefficient ($Q \in \mathbb{N}$). This leads to a joint distribution of the form:

$$p(\mathbf{y}, \hat{\mathbf{f}}, u | \mathbf{X}) = p(\mathbf{y} | \boldsymbol{\theta}) \prod_{j=1}^J p(\mathbf{f}_j | \mathbf{u}) p(\mathbf{u}), \quad (1-7)$$

where \mathbf{y} is the vector of noisy labels, $\hat{\mathbf{f}}$ is the vector of latent functions, u represents the inducing points, and \mathbf{X} is the input data.

Combined with inducing-variables based methods for sparse GP approximations, and maximizing an **Evidence Lower Bound (ELBO)** for the estimation of the variational parameters, the authors reach a model whose variational expectations are not analytically tractable, and hence, the authors derive a Gaussian-Hermite quadrature approach.

Finally, the authors extend this approach for being applied to classification and regression, reaching the only known approach to involve chained gaussian processes in multiple annotators classification and regression tasks while

471 preserving the interdependencies among the annotators, and also outperforming
 472 GPC-MV⁴, MA-LFC-C⁵, MA-DGRL⁶, MA-GPC⁷, MA-GPCV⁸, MA-DL⁹, KAAR¹⁰.

473 **CCGPMA** on itself proposes a good approach for handling label noise and
 474 uncertainty in the labelers for regression and classification tasks, while also
 475 preserving the interdependencies among the annotators, however, it does not face
 476 the image segmentation problem, which is the main focus of this works, however,
 477 it does not face the image segmentation problem, which is the main focus of this
 478 work. Besides, handling so many latent functions during the optimization process
 479 is computationally expensive, making it on itself infeasible for large and high
 480 resolution datasets.

481 1.3.2 Facing noisy annotations and low-quality data

482 The problem of low-quality data and noisy annotations has been tackled with
 483 various strategies. One such approach is the use of deep learning models that
 484 incorporate loss functions designed to mitigate the effects of unreliable labels.
 485 Traditional methods such as Majority Voting (MV) or Expectation-Maximization
 486 (EM) have been widely used for aggregating multiple annotators' inputs. However,
 487 they assume a homogeneous reliability of annotators, which may not hold in
 488 real-world scenarios.

⁴A GPC using the MV of the labels as the ground truth.

⁵A LRC with constant parameters across the input space.

⁶A multi-labeler approach that considers as latent variables the annotator performance.

⁷A multi-labeler GPC, which is an extension of MA-LFC.

⁸An extension of MA-GPC that includes variational inference and priors over the labelers' parameters.

⁹A Crowd Layer for DL, where the annotators' parameters are constant across the input space.

¹⁰A kernel-based approach that employs a convex combination of classifiers and codes labelers dependencies.

489 **Loss functions in deep learning models**

490 Loss functions are fundamental components in deep learning models that quantify
491 how well a model's predictions match the ground truth. They serve as the
492 objective function that guides the learning process by measuring the discrepancy
493 between predicted and actual values. In classification tasks, the most common
494 loss functions are Cross-Entropy (CE) and Mean Absolute Error (MAE). CE is
495 particularly effective for classification as it heavily penalizes confident but wrong
496 predictions, though it can be sensitive to noisy labels. MAE, on the other hand, is
497 more robust to outliers and assigns equal weights to all mistakes, but typically
498 requires more training iterations. For image segmentation tasks, specialized loss
499 functions have been developed to handle the unique challenges of pixel-wise
500 classification. The Dice loss, which measures the overlap between predicted and
501 ground truth regions, is widely used in medical image segmentation. More
502 recently, the Generalized Cross Entropy (GCE) loss has emerged as a robust
503 alternative that combines the benefits of both CE and MAE, allowing for better
504 handling of noisy labels through a tunable parameter that controls sensitivity to
505 outliers. In multi-annotator scenarios, where multiple experts provide potentially
506 inconsistent segmentations, novel loss functions like the Truncated Generalized
507 Cross Entropy for Semantic Segmentation (TGCE_{SS}) have been developed to
508 account for varying annotator reliability across different image regions. These loss
509 functions are crucial for training accurate segmentation models, especially in
510 medical imaging where precise delineation of anatomical structures is essential for
511 diagnosis and treatment planning.

512 **Generalized Cross-Entropy for multiple annotators classification**

513 A more recent approach was proposed by [Triana-Martinez et al., 2023],
514 introducing a Generalized Cross-Entropy-based Chained Deep Learning (GCECDL)
515 framework. This method addresses the limitations of traditional label aggregation
516 techniques by modeling each annotator's reliability as a function of the input data.

The approach effectively mitigates the impact of noisy labels by using a noise-robust loss function, balancing Mean Absolute Error (MAE) and Categorical Cross-Entropy (CE). Unlike prior approaches, **GCECDL** accounts for the dependencies among annotators while encoding their non-stationary behavior across different data samples. Their experiments on multiple datasets demonstrated superior predictive performance compared to state-of-the-art methods, particularly in cases where annotations were highly inconsistent.

The strategy of the authors effectively unlocks the potential of **ML** models to handle low-quality data and noisy annotations, but it is bounded to classifications tasks only, not being by itself applicable to segmentation tasks. The TGCE equation for handling multiple annotators is defined as:

$$\text{TGCE}(\mathbf{y}, f(\mathbf{x}); \tilde{\lambda}_x, \tilde{C}) = \tilde{\lambda}_x \frac{1 - (\mathbf{1}^\top (\mathbf{y} \odot f(\mathbf{x})))^q}{q} + (1 - \tilde{\lambda}_x) \frac{1 - (\tilde{C})^q}{q}, \quad (1-8)$$

where $\tilde{\lambda}_x$ represents the annotator reliability, \tilde{C} is a constant, q is a parameter that controls the balance between MAE and CE behavior, \mathbf{y} is the annotation vector, and $f(\mathbf{x})$ is the model prediction. This approach is more deeply discussed in chapter 4.

1.4 Aims

With the mentioned considerations in section 1.3 in mind, this work proposes a novel approach for **ISS** tasks in medical images, which aims to train a model whose learning approach is adaptive to the labeler performance. This is done by introducing a loss function capable of inferring the best possible segmentation without needing separate inputs about the labeler performance. This loss function is designed to implicitly weigh the labelers based on their performance, with the presence of an intermediate reliability map allowing the model to learn from the

most reliable labelers and ignore the noisy labels. This approach differs from existing CNN-based segmentation models, as it does not require explicit supervision of the labeler performance, making it more generalizable and adaptable to different datasets and labelers.

1.4.1 General Aim

The main purpose of this work is to develop a novel approach for ISS tasks in medical images, which can adaptively infer the best possible segmentation without needing separate inputs about the labeler performance. This approach is expected to outperform the segmentation performance of other state of the art approaches, correctly facing the labeler performance inconsistency across the annotators space and the variability of images quality.

1.4.2 Specific Aims

- To develop a novel loss function for ISS tasks in medical images, capable of inferring the best possible segmentation without needing separate inputs about the labeler performance.
- Introducing a tensor map which codifies the reliability of each labeler, allowing the model to implicitly weigh the labelers based on their performance across the mask and classes space.
- To develop and test a deep learning model for ISS tasks in medical images, which can learn from inconsistent labels and improve the segmentation performance compared to other solutions in state of the art.

1.5 Outline and Contributions

As an output of this work, some contributions were made to the field of ISS in medical images. The main contributions are:

- A python package for using the proposed loss function in CNN models for ISS tasks in medical images. ¹¹
- Datasets mapping as lazy loaders for the proposed loss function. ¹²
- A public Github repository with the code used in this work. ¹³

¹¹https://pypi.org/project/seg_tgce/

¹²<https://seg-tgce.readthedocs.io/en/latest/experiments.html>

¹³https://github.com/blotero/seg_tgce

567

568

CHAPTER

569

TWO

570

571

CONCEPTUAL PRELIMINARIES

572

2.1 Modern concept of digital image

573

574

575

576

577

A digital image is a numerical representation of a visual scene, captured through various imaging devices and stored in a computer. From a mathematical perspective, a digital image can be represented as a function $f(x, y)$ that maps spatial coordinates (x, y) to intensity values. In the discrete domain, this function is sampled at regular intervals, creating a matrix of values known as pixels (picture elements).

578

2.1.1 Types of digital images

579

Grayscale images

580

581

582

583

584

Grayscale images are the simplest form of digital images, where each pixel represents a single intensity value. Mathematically, a grayscale image can be represented as a 2D matrix I of size $M \times N$, where each element $I(i, j)$ represents the intensity at position (i, j) . The intensity values typically range from 0 (black) to 255 (white) in 8-bit images, or from 0 to 65535 in 16-bit images.

585 **Color images**

586 Color images extend the grayscale concept by representing each pixel with multiple
587 channels, typically Red, Green, and Blue (RGB). A color image can be represented
588 as a 3D matrix I of size $M \times N \times 3$, where $I(i, j, k)$ represents the intensity of the
589 k -th color channel at position (i, j) . Other color spaces like HSV (Hue, Saturation,
590 Value) or CMYK (Cyan, Magenta, Yellow, Key) are also commonly used in different
591 applications.

592 **Multispectral images**

593 Multispectral images capture information across multiple wavelength bands
594 beyond the visible spectrum. These images can be represented as a 3D matrix I of
595 size $M \times N \times B$, where B is the number of spectral bands. Each band $I(i, j, b)$
596 represents the intensity at position (i, j) for the b -th spectral band. This
597 representation is particularly useful in medical imaging, remote sensing, and
598 scientific applications.

599 **3D images and volumetric data**

600 Three-dimensional images extend the concept of pixels to voxels (volume elements).
601 A 3D image can be represented as a 3D matrix V of size $M \times N \times D$, where D
602 represents the depth dimension. Each voxel $V(i, j, k)$ represents the intensity at
603 position (i, j, k) in the 3D space. This representation is fundamental in medical
604 imaging (CT, MRI), scientific visualization, and computer graphics.

605 **2.1.2 Mathematical representations**

606 The mathematical foundation of digital images relies on several key concepts:

- 607 • **Sampling:** The process of converting a continuous image into a discrete
608 representation. According to the Nyquist-Shannon sampling theorem, the
609 sampling frequency must be at least twice the highest frequency present in
610 the image to avoid aliasing.
- 611 • **Quantization:** The process of converting continuous intensity values into
612 discrete levels. The number of quantization levels determines the image's
613 bit depth and affects its quality and storage requirements.
- 614 • **Resolution:** The number of pixels per unit length in an image, typically
615 measured in pixels per inch (PPI) or dots per inch (DPI).
- 616 • **Dynamic range:** The ratio between the maximum and minimum measurable
617 light intensities in an image, often expressed in decibels (dB).

618 The mathematical representation of a digital image can be expressed as:

$$I(x, y) = \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} f(i, j) \cdot \delta(x - i, y - j) \quad (2-1)$$

619 where $I(x, y)$ is the digital image, $f(i, j)$ represents the intensity values, and $\delta(x -$
620 $i, y - j)$ is the Kronecker delta function.

621 For color images, the representation extends to:

$$I(x, y) = \begin{bmatrix} I_R(x, y) \\ I_G(x, y) \\ I_B(x, y) \end{bmatrix} \quad (2-2)$$

622 where I_R , I_G , and I_B represent the red, green, and blue channels respectively.

623 **2.1.3 Digital histopathological images**

624 **2.2 Deep learning fundamentals**

625 **2.3 Datasets and data sources**

626

627

CHAPTER

628

THREE

629

630

CHAINED GAUSSIAN PROCESSES

631

3.1 Gaussian processes

632

3.2 Chained Gaussian processes

633

634

CHAPTER

635

FOUR

636

637

TRUNCATED GENERALIZED CROSS ENTROPY FOR

638

SEGMENTATION

639

4.1 Loss functions for multiple annotators

640

As mentioned in Section ??, a loss function is a key element for defining the

641

objective function of a deep learning model. The categorical cross-entropy loss is a

642

common loss function for classification tasks. However, in the case of multiple

643

annotators, the categorical cross-entropy loss is not able to handle the varying

644

reliability of the annotators. In this section, we will propose a loss function that is

645

able to handle multiple annotators' segmentation masks while accounting for their

646

varying reliability across different regions of the image.

647 4.1.1 Generalized Cross Entropy

648 The Generalized Cross Entropy (GCE) loss function was first introduced by [Zhang
649 and Sabuncu, 2018] as a robust alternative to the standard cross-entropy loss,
650 particularly effective in handling noisy labels. Let us first consider the Cross
651 Entropy (CE) and Mean Absolute Error (MAE) loss functions:

$$MAE(\mathbf{y}, f(\mathbf{x})) = \|\mathbf{y} - f(\mathbf{x})\|_1 \quad (4-1)$$

$$CE(\mathbf{y}, f(\mathbf{x})) = \sum_{k=1}^K y_k \log(f_k(\mathbf{x})) \quad (4-2)$$

652 where $y_k \in \mathbf{y}$, $f_k(\mathbf{x}) \in f(\mathbf{x})$, and $\|\cdot\|_1$ stands for the l_1 -norm. Of note, $\mathbf{1}^\top \mathbf{y} =$
653 $\mathbf{1}^\top f(\mathbf{x}) = 1$, $\mathbf{1} \in \{1\}^K$ being an all-ones vector. In addition, the MAE loss can be
654 rewritten for softmax outputs, yielding:

$$MAE(\mathbf{y}, f(\mathbf{x})) = 2(1 - \mathbf{1}^\top (\mathbf{y} \odot f(\mathbf{x}))) \quad (4-3)$$

655 where \odot stands for the element-wise product.

656 The CE is characterized by the following properties:

- 657 • It is unbounded from above.
- 658 • It heavily penalizes confident but wrong predictions.
- 659 • It is more sensitive to noisy labels.

660 On the other hand, the MAE is characterized by the following properties:

- 661 • It is bounded and more robust to outliers.
- 662 • It assigns equal weights to all mistakes regardless of confidence.
- 663 • It is symmetric in softmax based representations.
- 664 • It is more robust to noisy labels but slower to train.

665 The GCE loss function is defined by the authors in [Zhang and Sabuncu, 2018] as:

$$GCE(\mathbf{y}, f(\mathbf{x})) = 2 \frac{1 - (\mathbf{1}^\top (\mathbf{y} \odot f(\mathbf{x})))^q}{q}, \quad (4-4)$$

666 with $q \in (0, 1]$. Remarkably, the limiting case for $q \rightarrow 0$ in GCE is equivalent to the
 667 CE expression, and when $q = 1$, it equals the MAE loss. In addition, the GCE holds
 668 the following gradient with regard to θ :

$$\frac{\partial GCE(\mathbf{y}, f(\mathbf{x}; \theta)|_k)}{\partial \theta} = -f_k(\mathbf{x}; \theta)^{q-1} \nabla_\theta f_k(\mathbf{x}; \theta). \quad (4-5)$$

669 The GCE loss exhibits several desirable properties:

- 670 • It is more robust to label noise compared to standard cross-entropy
- 671 • The truncation parameter q allows for controlling the sensitivity to outliers
- 672 • It preserves the convexity property for optimization

673 4.1.2 Extension to Multiple Annotators

674 In the context of multiple annotators, we need to consider the varying reliability
 675 of each annotator across different regions of the image. Let's consider a k -class
 676 multiple annotators segmentation problem with the following data representation:

$$\mathbf{X} \in \mathbb{R}^{W \times H}, \{\mathbf{Y}_r \in \{0, 1\}^{W \times H \times K}\}_{r=1}^R; \quad \mathbf{\Psi} \in [0, 1]^{W \times H \times K} = f(\mathbf{X}) \quad (4-6)$$

677 where the segmentation mask function maps the input to output as:

$$f : \mathbb{R}^{W \times H} \rightarrow [0, 1]^{W \times H \times K} \quad (4-7)$$

678 The segmentation masks \mathbf{Y}_r satisfy the following condition for being a softmax-like
 679 representation:

$$\mathbf{Y}_r[w, h, :] \mathbf{1}_k^\top = 1; \quad w \in W, h \in H \quad (4-8)$$

680 4.1.3 Reliability Maps and Truncated GCE

681 The key innovation in our approach is the introduction of reliability maps Λ_r for
 682 each annotator:

$$\left\{ \Lambda_r(\mathbf{X}; \theta) \in [0, 1]^{W \times H} \right\}_{r=1}^R \quad (4-9)$$

683 These reliability maps estimate the confidence of each annotator at every spatial
 684 location (w, h) in the image. The maps are learned jointly with the segmentation
 685 model, allowing the network to:

- 686 • Weight the contribution of each annotator differently across the image
- 687 • Adapt to varying levels of expertise in different regions
- 688 • Handle cases where annotators might be more reliable in certain areas than
- 689 others

690 The proposed Truncated Generalized Cross Entropy for Semantic Segmentation
 691 (TGCE_{SS}) combines the robustness of GCE with the flexibility of reliability maps:

$$\begin{aligned}
 TGCE_{SS}(\mathbf{Y}_r, f(\mathbf{X}; \theta)|_r(\mathbf{X}; \theta)) = \mathbb{E}_r \left\{ \mathbb{E}_{w,h} \left\{ \Lambda_r(\mathbf{X}; \theta) \circ \mathbb{E}_k \left\{ \mathbf{Y}_r \circ \left(\frac{\mathbf{1}_{W \times H \times K} - f(\mathbf{X}; \theta)^{\circ q}}{q} \right); k \in K \right\} + \right. \right. \\
 \left. \left. (\mathbf{1}_{W \times H} - \Lambda_r(\mathbf{X}; \theta)) \circ \left(\frac{\mathbf{1}_{W \times H} - (\frac{1}{k} \mathbf{1}_{W \times H})^{\circ q}}{q} \right); w \in W, h \in H \right\}; r \in R \right\}
 \end{aligned}
 \tag{4-10}$$

692 where $q \in (0, 1)$ controls the truncation level. The loss function consists of two
 693 main components:

- 694 • The first term weighted by Λ_r represents the GCE loss for regions where the
- 695 annotator is considered reliable
- 696 • The second term weighted by $(1 - \Lambda_r)$ provides a uniform prior for regions
- 697 where the annotator is considered unreliable

698 For a batch containing N samples, the total loss is computed as:

$$\mathcal{L}(\mathbf{Y}_r[n], f(\mathbf{X}[n]; \theta)|_r(\mathbf{X}[n]; \theta)) = \frac{1}{N} \sum_n TGCE_{SS}(\mathbf{Y}_r[n], f(\mathbf{X}[n]; \theta)|_r(\mathbf{X}[n]; \theta))
 \tag{4-11}$$

699 4.2 Proposed Model

700 Our proposed model architecture combines the strengths of UNET with a ResNet-
701 34 backbone, specifically designed to work with the $TGCE_{SS}$ loss function. The
702 architecture is illustrated in Figure ??.

703 4.2.1 Backbone Architecture

704 The model uses a pre-trained ResNet-34 as its encoder backbone. ResNet-34's deep
705 residual learning framework provides several advantages:

- 706 • Efficient feature extraction through residual connections
- 707 • Pre-trained weights that capture rich visual representations
- 708 • Stable gradient flow during training

709 The ResNet-34 backbone is modified to serve as the encoder in our UNET
710 architecture. We remove the final fully connected layer and use the feature maps
711 from different stages of the network for skip connections.

712 4.2.2 UNET Architecture

713 The UNET architecture consists of an encoder-decoder structure with skip
714 connections. The encoder path follows the ResNet-34 structure, while the decoder
715 path uses transposed convolutions for upsampling. The architecture includes:

- 716 • Four downsampling stages in the encoder (ResNet-34 blocks)
- 717 • Four upsampling stages in the decoder
- 718 • Skip connections between corresponding encoder and decoder stages
- 719 • Batch normalization and ReLU activation after each convolution

4.2.3 Reliability Map Branch

A key innovation in our architecture is the addition of a parallel branch for estimating reliability maps. This branch:

- Takes the same encoder features as input
- Uses a series of 1×1 convolutions to reduce channel dimensions
- Produces R reliability maps Λ_r for each annotator
- Applies a sigmoid activation to ensure values in $[0, 1]$

4.2.4 Integration with TGCE_{SS} Loss

The model outputs two components:

- Segmentation masks $\mathbf{\hat{Y}} = f(\mathbf{X}; \theta)$
- Reliability maps $\{\Lambda_r(\mathbf{X}; \theta)\}_{r=1}^R$

These outputs are used to compute the TGCE_{SS} loss as described in Section ?? . The loss function guides the learning of both the segmentation masks and reliability maps simultaneously.

734 4.2.5 Training Process

735 The training process involves:

- 736 • Initializing the ResNet-34 backbone with pre-trained weights
- 737 • Training the entire network end-to-end
- 738 • Using the Adam optimizer with a learning rate of 10^{-4}
- 739 • Applying the $TGCE_{SS}$ loss to update both the segmentation and reliability
- 740 branches

741 The model's architecture allows it to:

- 742 • Learn robust segmentation features through the ResNet-34 backbone
- 743 • Capture fine-grained details through UNET's skip connections
- 744 • Adapt to annotator reliability through the parallel reliability branch
- 745 • Handle multiple annotators' inputs effectively

746 4.3 Experiments

747 4.3.1 Dataset

748 4.3.2 Metrics

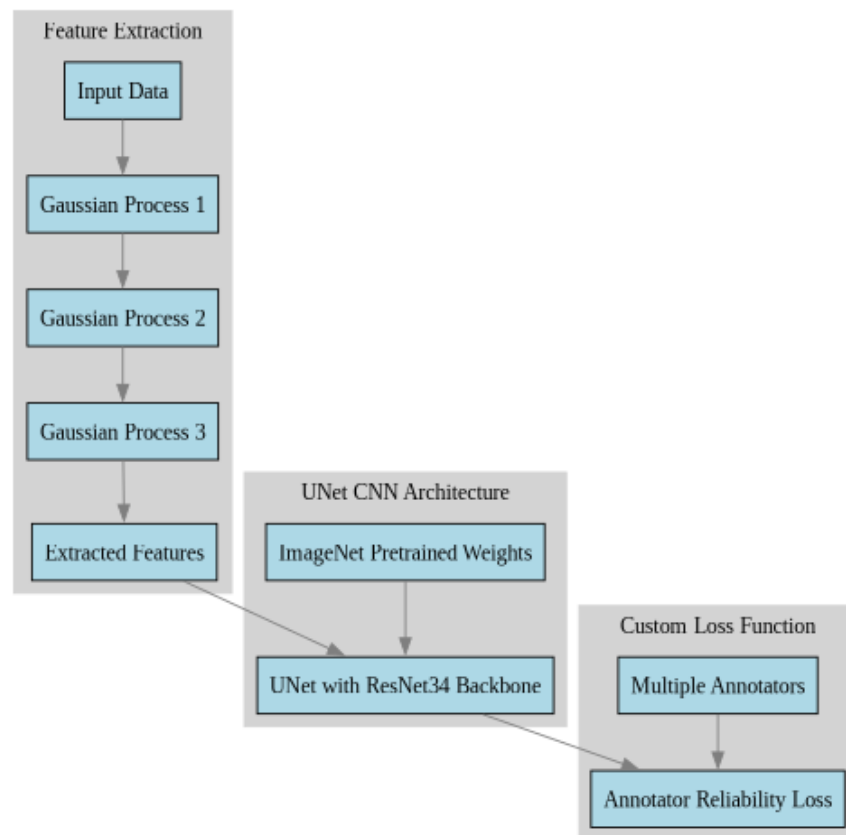


Figure 4-1 Solution Architecture (mockup)

BIBLIOGRAPHY

- 750 [Avanzo et al., 2024] Avanzo, M., Stancanella, J., Pirrone, G., Drigo, A., and Retico,
751 A. (2024). The evolution of artificial intelligence in medical imaging: From
752 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
753 Author Joseph Stancanella is employed by Elekta SA. The remaining authors
754 declare no commercial or financial conflicts of interest. (page 3)
- 755 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
756 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
757 (2024). Medical image segmentation review: The success of u-net. *IEEE*
758 *Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
759 (page 2)
- 760 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
761 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in*
762 *Radiology*, 5. (page 8)
- 763 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
764 cost minimization for medical image segmentation using suggestive mixed
765 supervision fully convolutional networks. *CoRR*, abs/1812.11302. (page 3)
- 766 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
767 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
768 literature review. (page 2)

- [Carmo et al., 2025] Carmo, D. S., Pezzulo, A. A., Villacreses, R. A., Eisenbeisz, M. L., Anderson, R. L., Van Dorin, S. E., Rittner, L., Lotufo, R. A., Gerard, S. E., Reinhardt, J. M., and Comellas, A. P. (2025). Manual segmentation of opacities and consolidations on ct of long covid patients from multiple annotators. *Scientific Data*, 12(1):402. (page 9)
- [Elhaminia et al., 2025] Elhaminia, B., Alsalemi, A., Nasir, E., Jahanifar, M., Awan, R., Young, L. S., Rajpoot, N. M., Minhas, F., and Raza, S. E. A. (2025). From traditional to deep learning approaches in whole slide image registration: A methodological review. (page 2)
- [Elnakib et al., 2020] Elnakib, A., Elmenabawy, N., and S Moustafa, H. (2020). Automated deep system for joint liver and tumor segmentation using majority voting. *MEJ-Mansoura Engineering Journal*, 45(4):30–36. (page 11)
- [Giri and Bhatia, 2024] Giri, K. and Bhatia, S. (2024). Artificial intelligence in nephrology- its applications from bench to bedside. *International Journal of Advances in Nephrology Research*, 7(1):90–97. (page 6)
- [Grefve et al., 2024] Grefve, J., Söderkvist, K., Gunnlaugsson, A., Sandgren, K., Jonsson, J., Keeratijarut Lindberg, A., Nilsson, E., Axelsson, J., Bergh, A., Zackrisson, B., Moreau, M., Thellenberg Karlsson, C., Olsson, L., Widmark, A., Riklund, K., Blomqvist, L., Berg Loegager, V., Strandberg, S. N., and Nyholm, T. (2024). Histopathology-validated gross tumor volume delineations of intraprostatic lesions using psma-positron emission tomography/multiparametric magnetic resonance imaging. *Physics and Imaging in Radiation Oncology*, 31:100633. (page 14)
- [Habis, 2024] Habis, A. A. (2024). *Developing interactive artificial intelligence tools to assist pathologists with histology annotation*. Theses, Institut Polytechnique de Paris. (page 8)
- [Hu et al., 2025] Hu, D., Jiang, Z., Shi, J., Xie, F., Wu, K., Tang, K., Cao, M., Huai, J., and Zheng, Y. (2025). Pathology report generation from whole slide images with knowledge retrieval and multi-level regional feature selection. *Computer Methods and Programs in Biomedicine*, 263:108677. (page 2)

- 799 [Julián and Álvarez Meza Andrés Marino, 2023] Julián, G. G. and Álvarez Meza
800 Andrés Marino (2023). A supervised learning framework in the context of
801 multiple annotators. (page 17)
- 802 [Karthikeyan et al., 2023] Karthikeyan, R., McDonald, A., and Mehta, R. (2023).
803 What’s in a label? annotation differences in forecasting mental fatigue using ecg
804 data and seq2seq architectures. (page 9)
- 805 [Kim et al., 2024] Kim, Y., Lee, E., Lee, Y., and Oh, U. (2024). Understanding
806 novice’s annotation process for 3d semantic segmentation task with human-
807 in-the-loop. In *Proceedings of the 29th International Conference on Intelligent User*
808 *Interfaces, IUI ’24*, page 444–454, New York, NY, USA. Association for Computing
809 Machinery. (page 9)
- 810 [Lam and Suen, 1997] Lam, L. and Suen, S. (1997). Application of majority
811 voting to pattern recognition: an analysis of its behavior and performance.
812 *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*,
813 27(5):553–568. (page 11)
- 814 [Lin et al., 2024] Lin, Y., Lian, A., Liao, M., and Yuan, S. (2024). Bcdnet: A fast
815 residual neural network for invasive ductal carcinoma detection. (page 6)
- 816 [López-Pérez et al., 2023] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A. D.,
817 Molina, R., and Katsaggelos, A. K. (2023). Crowdsourcing segmentation
818 of histopathological images using annotations provided by medical students.
819 In Juarez, J. M., Marcos, M., Stiglic, G., and Tucker, A., editors, *Artificial*
820 *Intelligence in Medicine*, pages 245–249, Cham. Springer Nature Switzerland.
821 (pages 5, 6, 8, and 12)
- 822 [Lu et al., 2023] Lu, X., Ratcliffe, D., Kao, T.-T., Tikhonov, A., Litchfield, L., Rodger,
823 C., and Wang, K. (2023). Rethinking quality assurance for crowdsourced multi-
824 roi image segmentation. *Proceedings of the AAAI Conference on Human Computation*
825 *and Crowdsourcing*, 11(1):103–114. (pages 5 and 8)

- [López-Pérez et al., 2024] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A., Felicelli, C., Goldstein, J., Vadasz, B., Molina, R., and Katsaggelos, A. K. (2024). Learning from crowds for automated histopathological image segmentation. *Computerized Medical Imaging and Graphics*, 112:102327. (pages xv, 5, 15, and 17)
- [Panayides et al., 2020] Panayides, A. S., Amini, A., Filipovic, N. D., Sharma, A., Tsiftaris, S. A., Young, A., Foran, D., Do, N., Golemati, S., Kurc, T., Huang, K., Nikita, K. S., Veasey, B. P., Zervakis, M., Saltz, J. H., and Pattichis, C. S. (2020). Ai in medical imaging informatics: Current challenges and future directions. *IEEE Journal of Biomedical and Health Informatics*, 24(7):1837–1857. (page 2)
- [Qiu et al., 2022] Qiu, Y., Hu, Y., Kong, P., Xie, H., Zhang, X., Cao, J., Wang, T., and Lei, B. (2022). Automatic prostate gleason grading using pyramid semantic parsing network in digital histopathology. *Frontiers in Oncology*, 12. (page 14)
- [Rashmi et al., 2021] Rashmi, R., Prasad, K., and Udupa, C. B. K. (2021). Breast histopathological image analysis using image processing techniques for diagnostic purposes: A methodological review. *Journal of Medical Systems*, 46(1):7. (pages 1 and 5)
- [Ronneberger et al., 2015] Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical image segmentation. In Navab, N., Hornegger, J., Wells, W. M., and Frangi, A. F., editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham. Springer International Publishing. (page 14)
- [Ryou et al., 2025] Ryou, H., Thomas, E., Wojciechowska, M., Harding, L., Tam, K. H., Wang, R., Hu, X., Rittscher, J., Cooper, R., and Royston, D. (2025). Reticulin-free quantitation of bone marrow fibrosis in mpns: Utility and applications. *eJHaem*, 6(2):e70005. (page 2)
- [Sarvamangala and Kulkarni, 2022] Sarvamangala, D. R. and Kulkarni, R. V. (2022). Convolutional neural networks in medical image understanding: a survey. *Evolutionary Intelligence*, 15(1):1–22. (pages 3 and 6)

- 854 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018).
855 Ms-net: Mixed-supervision fully-convolutional networks for full-resolution
856 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
857 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer*
858 *Assisted Intervention – MICCAI 2018*, pages 379–387, Cham. Springer International
859 Publishing. (page 5)
- 860 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
861 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
862 *Sciences*, 378(2166):20190061. (page 3)
- 863 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
864 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
865 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
866 Curran Associates, Inc. (page 12)
- 867 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
868 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
869 Chained deep learning using generalized cross-entropy for multiple annotators
870 classification. *Sensors*, 23(7). (page 20)
- 871 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
872 truth and performance level estimation (staple): an algorithm for the validation
873 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
874 (page 12)
- 875 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
876 Advances in medical image segmentation: A comprehensive review of traditional,
877 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 878 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
879 Semi-supervised gland segmentation via feature-enhanced contrastive learning
880 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
881 pages 1–11. (page 2)

- 882 [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Cicarrelli, O.,
883 Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground
884 truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell,
885 R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*,
886 volume 33, pages 15750–15762. Curran Associates, Inc. (page 16)
- 887 [Zhang and Sabuncu, 2018] Zhang, Z. and Sabuncu, M. R. (2018). Generalized
888 cross entropy loss for training deep neural networks with noisy labels.
889 (pages 32 and 33)
- 890 [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S.,
891 Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M.
892 (2021). A review of deep learning in medical imaging: Imaging traits, technology
893 trends, case studies with progress highlights, and future promises. *Proceedings of*
894 *the IEEE*, 109(5):820–838. (pages 1 and 2)
- 895 [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L.
896 (2024). Image quality evaluation in deep-learning-based ct noise reduction using
897 virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical*
898 *Physics*, 51(8):5399–5413. (page 9)