



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	20
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 Chained Gaussian Processes	25
78	3 Deep Learning for Image Segmentation	27
79	4 Truncated Generalized Cross Entropy for segmentation	29
80	4.1 Loss functions for multiple annotators	29
81	4.1.1 Generalized Cross Entropy	30
82	4.1.2 Extension to Multiple Annotators	32
83	4.1.3 Reliability Maps and Truncated GCE	32
84	4.2 Proposed Model	34
85	4.2.1 Backbone Architecture	34
86	4.2.2 UNET Architecture	34
87	4.2.3 Reliability Map Branch	35
88	4.2.4 Integration with $TGCE_{SS}$ Loss	35
89	4.2.5 Training Process	36
90	4.3 Experiments	36
91	4.3.1 Dataset	36
92	4.3.2 Metrics	36
93	Bibliography	38

LIST OF FIGURES

95 **1-1** Estimation of the tasks and medical image types based on recent
96 literature review (count of referenced terms). 3
97 **1-2** AI and machine learning in medical imaging brief timeline. 4
98 **1-3** Example of a histopathological image segmented by multiple
99 annotators, illustrating variations in label assignment. 7
100 **1-4** Summary diagram for problem Statement 10
101 **1-5** Original U-Net architecture. 15
102 **1-6** Proposed framework for the approach in [López-Pérez et al., 2024]. . 18
103 **4-1** Solution Architecture (mockup) 37

LIST OF TABLES

106	CAD	Computer-Aided Diagnosis 2, 5, 6
107	CCGP	Correlated Chained Gaussian Processes 18
108	CCGPMA	Correlated Chained Gaussian Processes for Multiple Annotators 18, 19
109	CE	Cross Entropy 30
110	CGP	Chained Gaussian Processes 18
111	CNN	Convolutional Neural Networks 3, 14, 22, 23
112	CT	Computed Tomography 12
113	ELBO	Evidence Lower Bound 19
114	GCE	Generalized Cross Entropy 30
115	GCECDL	Generalized Cross-Entropy-based Chained Deep Learning 21
116	ISS	Image Semantic segmentation 2, 3, 6, 11, 13, 21-23
117	LF	Latent Function 18
118	MAE	Mean Absolute Error 30, 31
119	MITs	Medical Imaging Techniques 1
120	ML	Machine Learning 11, 21
121	MV	Majority Voting 11, 12
122	OCR	Optical Character Recognition 11
123	PET	Positron Emission Tomography 14
124	ROI	Region of Interest 2, 6
125	SLFM	Semi-Parametric Latent Factor Model 18
126	SS	Semantic segmentation 3
127	STAPLE	Simultaneous Truth and Performance Level Estimation 12-14
128	WSI	Whole Slide Imaging 1, 5, 6, 8, 16

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

CHAPTER

ONE

INTRODUCTION

1.1 Motivation

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

Along with the advances in technologies for medical images acquisition, 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

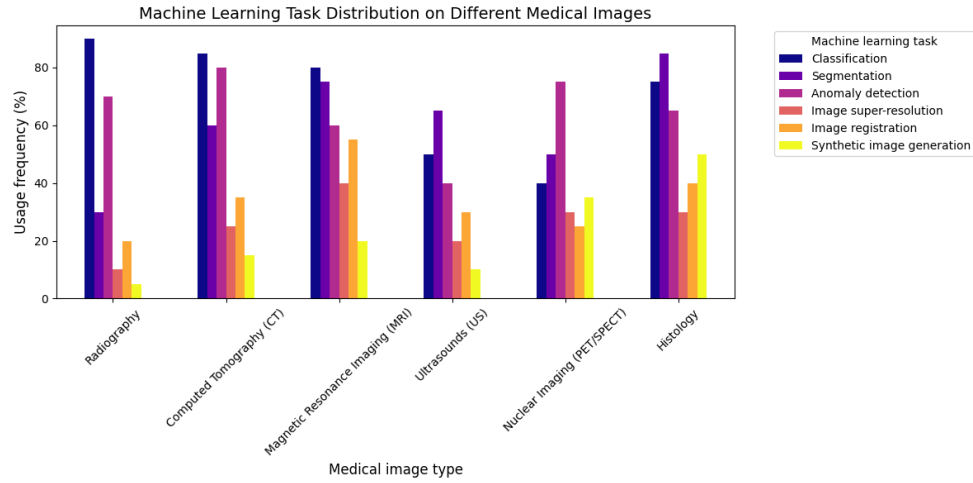


Figure 1-1 Estimation of the tasks and medical image types based on recent literature review (count of referenced terms).

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

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

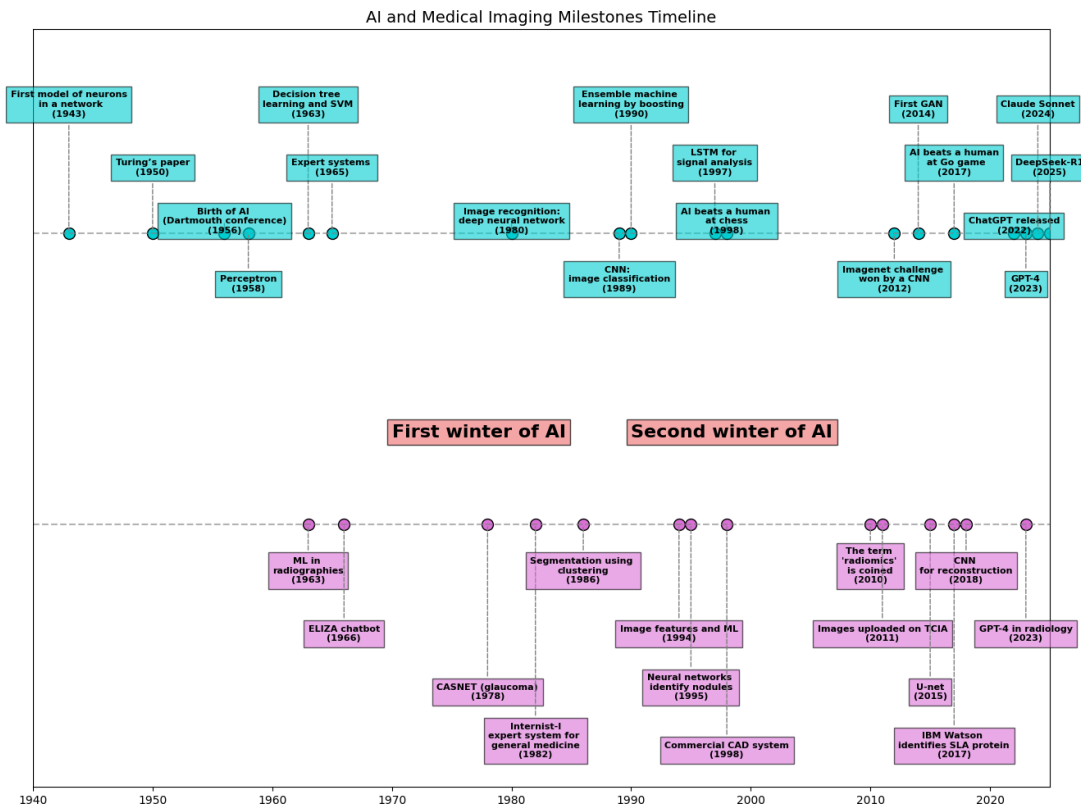


Figure 1-2 AI and machine learning in medical imaging brief timeline.

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

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

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

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

220 A major challenge in histopathological image segmentation arises from the
221 variability in annotations provided by different pathologists. Unlike natural
222 images, where object boundaries are often well-defined, histological structures
223 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 assignment¹. 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 segments².

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

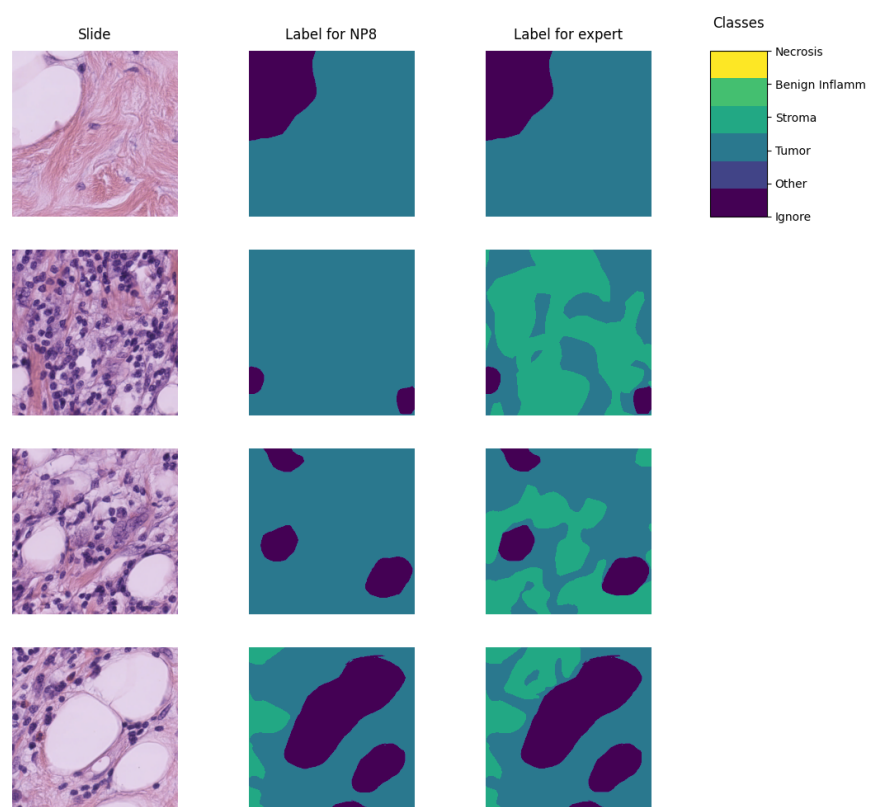


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

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].

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

281 1.2.2 Technical Constraints and Image Quality

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

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

292 1.2.3 Research Question

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

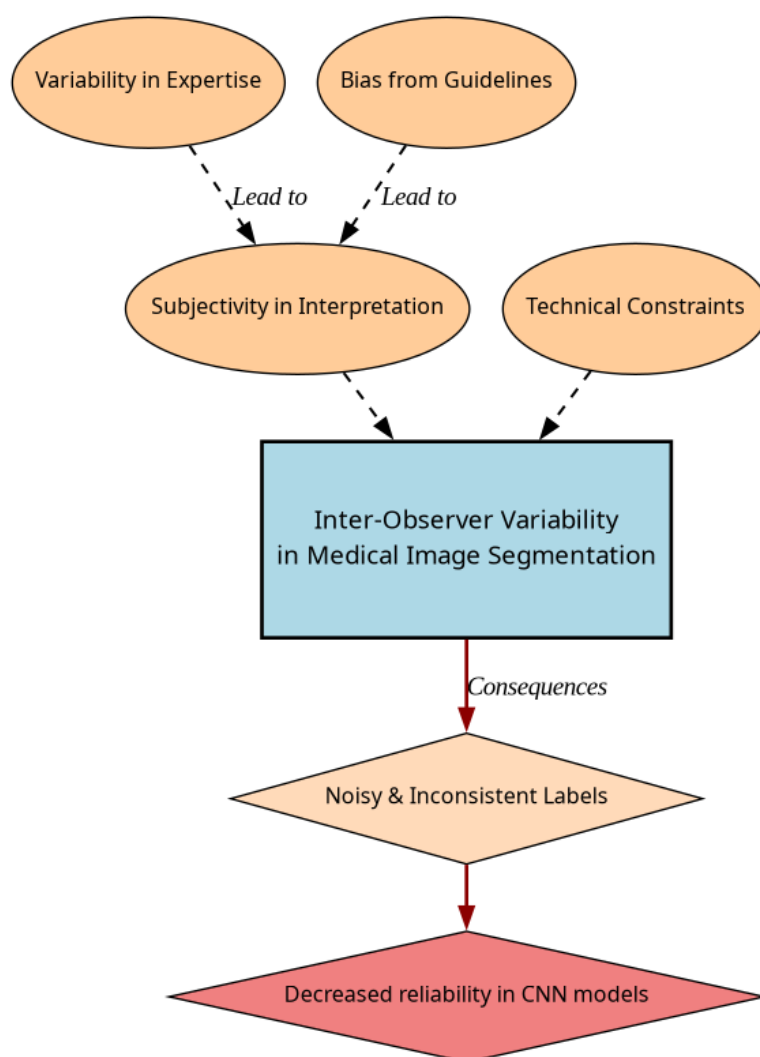


Figure 1-4 Summary diagram for problem Statement

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? Can such approach face problems related to the variability in expertise levels and technical constraints while preserving interpretability, generalization and computational efficiency?

295

1.3 Literature review

296

297 Certainly, in general Machine Learning (ML) classification tasks ³ where multiple
298 annotators are involved, Majority Voting (MV) is by far the simplest possible
299 approach to implement. This concept was born multiple times and divergently in
300 multiple fields, but it was described as relevant for ML and pattern recognition
301 labeling for classification in [Lam and Suen, 1997], in which the approach is
302 exposed as simple, yet powerful. The authors describe the MV as a method that
303 can be used to improve the accuracy of classification tasks by combining the labels
304 of multiple annotators. The method is based on the assumption that the majority
305 vote of the annotators is more likely to be correct than the vote of a single
306 annotator. The authors also describe the method as a straightforward way to
307 improve the accuracy of classification tasks without the need for complex
308 algorithms or additional data. The authors also prove this method to deliver very
309 similar results to more complicated approaches (Bayesian, logistic regression,
310 fuzzy integral, and neural network) in the particular task of Optical Character
311 Recognition (OCR). Despite its simplicity, modern solutions for delivering accurate
312 medical image segmentation models still rely on Majority Voting at some stage,
313 like [Elnakib et al., 2020], which uses a majority voting strategy for delivering a
314 final output based on the labels of multiple models (VGG16-Segnet, Resnet-18 and

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

Alexnet) in **Computed Tomography (CT)** images for Liver Tumor Segmentation, or [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. The architecture overview can be seen in figure 1-5.

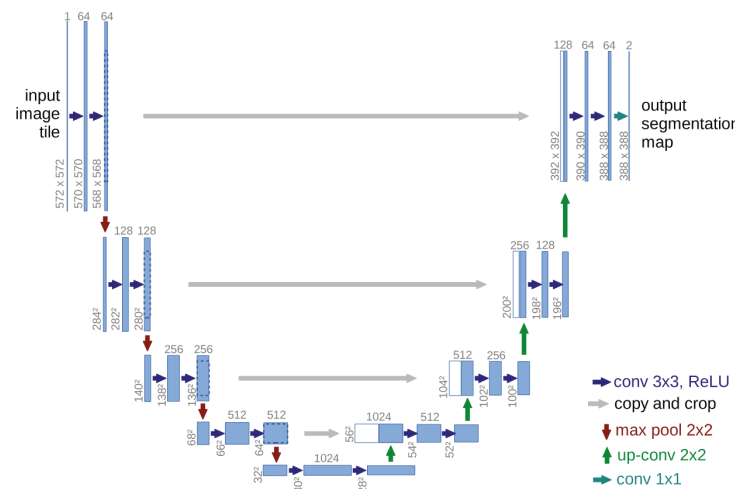


Figure 1-5 Original U-Net architecture.

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:

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

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

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

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

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

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

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

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

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

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

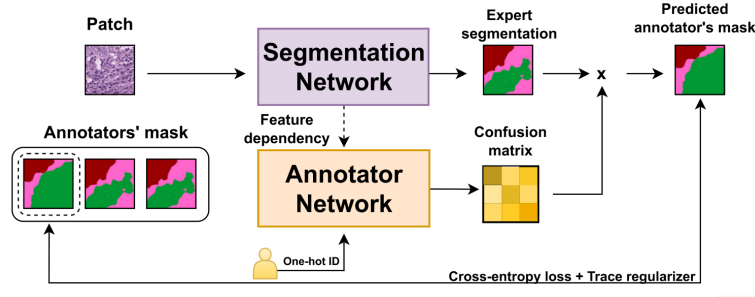


Figure 1-6 Proposed framework for the approach in [López-Pérez et al., 2024].

Bayesian models

Bayesian approaches are a good choice for handling label noise and uncertainty in the labelers. In [Julián and Álvarez Meza Andrés Marino, 2023] the authors propose a novel approach from Gaussian Processes to model the relationship between the annotators' reliability and the input data, while also preserving the interdependencies among the annotators. This is achieved by introducing **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 preserving the interdependencies among the annotators, and also outperforming GPC-MV⁴, MA-LFC-C⁵, MA-DGRL⁶, MA-GPC⁷, MA-GPCV⁸, MA-DL⁹, KAAR¹⁰.

CCGPMA on itself proposes a good approach for handling label noise and uncertainty in the labelers for regression and classification tasks, while also preserving the interdependencies among the annotators, however, it does not face the image segmentation problem, which is the main focus of this work, however, it does not face the image segmentation problem, which is the main focus of this

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.

work. Besides, handling so many latent functions during the optimization process is computationally expensive, making it on itself infeasible for large and high resolution datasets.

1.3.2 Facing noisy annotations and low-quality data

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

Loss functions in deep learning models

Loss functions are fundamental components in deep learning models that quantify how well a model's predictions match the ground truth. They serve as the objective function that guides the learning process by measuring the discrepancy between predicted and actual values. In classification tasks, the most common loss functions are Cross-Entropy (CE) and Mean Absolute Error (MAE). CE is particularly effective for classification as it heavily penalizes confident but wrong predictions, though it can be sensitive to noisy labels. MAE, on the other hand, is more robust to outliers and assigns equal weights to all mistakes, but typically requires more training iterations. For image segmentation tasks, specialized loss functions have been developed to handle the unique challenges of pixel-wise classification. The Dice loss, which measures the overlap between predicted and ground truth regions, is widely used in medical image segmentation. More recently, the Generalized Cross Entropy (GCE) loss has emerged as a robust alternative that combines the benefits of both CE and MAE, allowing for better

handling of noisy labels through a tunable parameter that controls sensitivity to outliers. In multi-annotator scenarios, where multiple experts provide potentially inconsistent segmentations, novel loss functions like the Truncated Generalized Cross Entropy for Semantic Segmentation ($TGCE_{SS}$) have been developed to account for varying annotator reliability across different image regions. These loss functions are crucial for training accurate segmentation models, especially in medical imaging where precise delineation of anatomical structures is essential for diagnosis and treatment planning.

Generalized Cross-Entropy for multiple annotators classification

A more recent approach was proposed by [Triana-Martinez et al., 2023], introducing a Generalized Cross-Entropy-based Chained Deep Learning (GCECDL) framework. This method addresses the limitations of traditional label aggregation 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.

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, eliminate the need for explicit labeler supervision, and enhance automation in medical image analysis.

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

557

558

559

560

CHAPTER

TWO

561

CHAINED GAUSSIAN PROCESSES

562

563

564

565

CHAPTER

THREE

566

DEEP LEARNING FOR IMAGE SEGMENTATION

567

568

CHAPTER

569

FOUR

570

571

TRUNCATED GENERALIZED CROSS ENTROPY FOR

572

SEGMENTATION

573

4.1 Loss functions for multiple annotators

574

As mentioned in Section 3, a loss function is a key element for defining the objective function of a deep learning model. The categorical cross-entropy loss is a common loss function for classification tasks. However, in the case of multiple annotators, the categorical cross-entropy loss is not able to handle the varying reliability of the annotators. In this section, we will propose a loss function that is able to handle multiple annotators' segmentation masks while accounting for their varying reliability across different regions of the image.

575

576

577

578

579

580

4.1.1 Generalized Cross Entropy

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

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

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

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, $\mathbf{1}^\top \mathbf{y} = \mathbf{1}^\top f(\mathbf{x}) = 1$, $\mathbf{1} \in \{1\}^K$ being an all-ones vector. In addition, the MAE loss can be rewritten for softmax outputs, yielding:

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

where \odot stands for the element-wise product.

The CE is characterized by the following properties:

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

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

- 595 • It is bounded and more robust to outliers.
- 596 • It assigns equal weights to all mistakes regardless of confidence.
- 597 • It is symmetric in softmax based representations.
- 598 • It is more robust to noisy labels but slower to train.

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

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

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

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

603 The GCE loss exhibits several desirable properties:

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

607 4.1.2 Extension to Multiple Annotators

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

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

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

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

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

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

614 4.1.3 Reliability Maps and Truncated GCE

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

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

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

- 620 • Weight the contribution of each annotator differently across the image
- 621 • Adapt to varying levels of expertise in different regions
- 622 • Handle cases where annotators might be more reliable in certain areas than
- 623 others

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

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

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

- 628 • The first term weighted by Λ_r represents the GCE loss for regions where the
- 629 annotator is considered reliable
- 630 • The second term weighted by $(1 - \Lambda_r)$ provides a uniform prior for regions
- 631 where the annotator is considered unreliable

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

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

633 4.2 Proposed Model

634 Our proposed model architecture combines the strengths of UNET with a ResNet-
635 34 backbone, specifically designed to work with the $TGCE_{SS}$ loss function. The
636 architecture is illustrated in Figure ??.

637 4.2.1 Backbone Architecture

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

- 640 • Efficient feature extraction through residual connections
- 641 • Pre-trained weights that capture rich visual representations
- 642 • Stable gradient flow during training

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

646 4.2.2 UNET Architecture

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

- 650 • Four downsampling stages in the encoder (ResNet-34 blocks)
- 651 • Four upsampling stages in the decoder
- 652 • Skip connections between corresponding encoder and decoder stages
- 653 • Batch normalization and ReLU activation after each convolution

654 4.2.3 Reliability Map Branch

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

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

661 4.2.4 Integration with TGCE_{SS} Loss

662 The model outputs two components:

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

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

668 4.2.5 Training Process

669 The training process involves:

- 670 • Initializing the ResNet-34 backbone with pre-trained weights
- 671 • Training the entire network end-to-end
- 672 • Using the Adam optimizer with a learning rate of 10^{-4}
- 673 • Applying the $TGCE_{SS}$ loss to update both the segmentation and reliability
- 674 branches

675 The model's architecture allows it to:

- 676 • Learn robust segmentation features through the ResNet-34 backbone
- 677 • Capture fine-grained details through UNET's skip connections
- 678 • Adapt to annotator reliability through the parallel reliability branch
- 679 • Handle multiple annotators' inputs effectively

680 4.3 Experiments

681 4.3.1 Dataset

682 4.3.2 Metrics

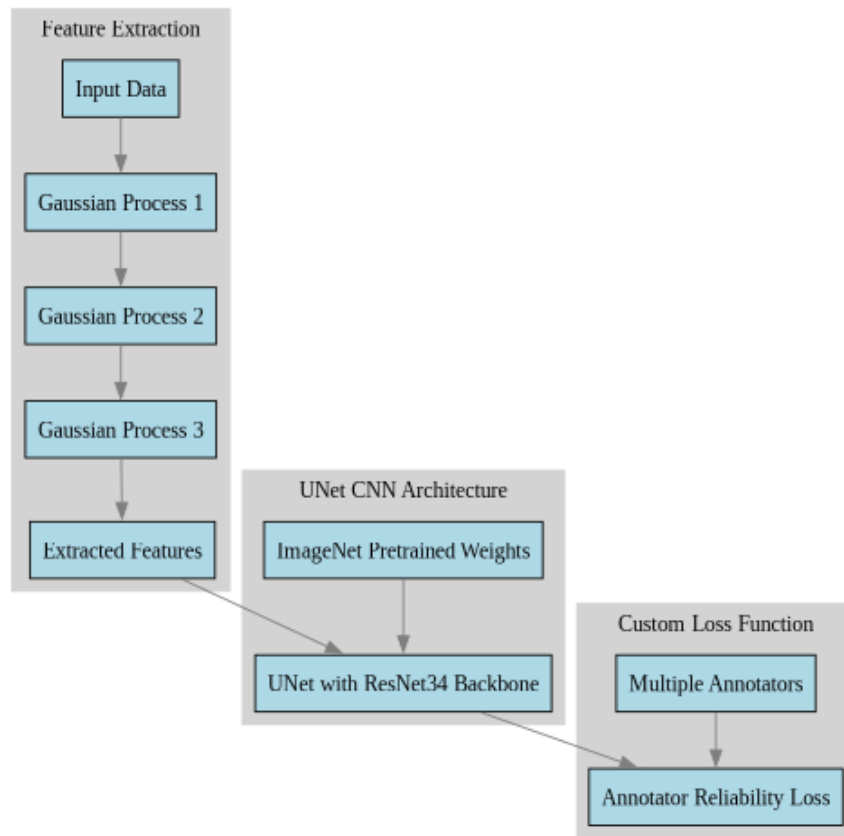


Figure 4-1 Solution Architecture (mockup)

BIBLIOGRAPHY

- 684 [Avanzo et al., 2024] Avanzo, M., Stancanella, J., Pirrone, G., Drigo, A., and Retico,
685 A. (2024). The evolution of artificial intelligence in medical imaging: From
686 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
687 Author Joseph Stancanella is employed by Elekta SA. The remaining authors
688 declare no commercial or financial conflicts of interest. (page 3)
- 689 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
690 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
691 (2024). Medical image segmentation review: The success of u-net. *IEEE*
692 *Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
693 (page 2)
- 694 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
695 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in*
696 *Radiology*, 5. (page 8)
- 697 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
698 cost minimization for medical image segmentation using suggestive mixed
699 supervision fully convolutional networks. CoRR, abs/1812.11302. (page 3)
- 700 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
701 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
702 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)

- [Julián and Álvarez Meza Andrés Marino, 2023] Julián, G. G. and Álvarez Meza Andrés Marino (2023). A supervised learning framework in the context of multiple annotators. (page 18)
- [Karthikeyan et al., 2023] Karthikeyan, R., McDonald, A., and Mehta, R. (2023). What's in a label? annotation differences in forecasting mental fatigue using ecg data and seq2seq architectures. (page 9)
- [Kim et al., 2024] Kim, Y., Lee, E., Lee, Y., and Oh, U. (2024). Understanding novice's annotation process for 3d semantic segmentation task with human-in-the-loop. In *Proceedings of the 29th International Conference on Intelligent User Interfaces, IUI '24*, page 444–454, New York, NY, USA. Association for Computing Machinery. (page 9)
- [Lam and Suen, 1997] Lam, L. and Suen, S. (1997). Application of majority voting to pattern recognition: an analysis of its behavior and performance. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 27(5):553–568. (page 11)
- [Lin et al., 2024] Lin, Y., Lian, A., Liao, M., and Yuan, S. (2024). Bcdnet: A fast residual neural network for invasive ductal carcinoma detection. (page 6)
- [López-Pérez et al., 2023] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A. D., Molina, R., and Katsaggelos, A. K. (2023). Crowdsourcing segmentation of histopathological images using annotations provided by medical students. In Juarez, J. M., Marcos, M., Stiglic, G., and Tucker, A., editors, *Artificial Intelligence in Medicine*, pages 245–249, Cham. Springer Nature Switzerland. (pages 5, 6, 8, and 12)
- [Lu et al., 2023] Lu, X., Ratcliffe, D., Kao, T.-T., Tikhonov, A., Litchfield, L., Rodger, C., and Wang, K. (2023). Rethinking quality assurance for crowdsourced multi-roi image segmentation. *Proceedings of the AAAI Conference on Human Computation 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 18)
- [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)

- 788 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018).
789 Ms-net: Mixed-supervision fully-convolutional networks for full-resolution
790 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
791 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer*
792 *Assisted Intervention – MICCAI 2018*, pages 379–387, Cham. Springer International
793 Publishing. (page 5)
- 794 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
795 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
796 *Sciences*, 378(2166):20190061. (page 3)
- 797 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
798 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
799 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
800 Curran Associates, Inc. (page 12)
- 801 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
802 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
803 Chained deep learning using generalized cross-entropy for multiple annotators
804 classification. *Sensors*, 23(7). (page 21)
- 805 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
806 truth and performance level estimation (staple): an algorithm for the validation
807 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
808 (page 12)
- 809 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
810 Advances in medical image segmentation: A comprehensive review of traditional,
811 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 812 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
813 Semi-supervised gland segmentation via feature-enhanced contrastive learning
814 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
815 pages 1–11. (page 2)

- [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Cicarrelli, O., Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 15750–15762. Curran Associates, Inc. (page 17)
- [Zhang and Sabuncu, 2018] Zhang, Z. and Sabuncu, M. R. (2018). Generalized cross entropy loss for training deep neural networks with noisy labels. (pages 30 and 31)
- [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S., Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M. (2021). A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. *Proceedings of the IEEE*, 109(5):820–838. (pages 1 and 2)
- [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L. (2024). Image quality evaluation in deep-learning-based ct noise reduction using virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical Physics*, 51(8):5399–5413. (page 9)