



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	xv
61	List of figures	xvii
62	List of tables	xix
63	Abbreviations	xxi
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	2.2	Digital histopathological images	28
82	2.2.1	Whole Slide Imaging (WSI)	28
83	2.2.2	Regions of Interest (ROI)	29
84	2.2.3	Staining Techniques	29
85	2.3	Deep learning fundamentals	31
86	2.3.1	Learning Paradigms	32
87	2.3.2	Architecture and Training	32
88	2.3.3	Challenges and Solutions	33
89	2.3.4	Deep Learning Frameworks	33
90	2.4	Datasets and data sources	34
91	3	Chained Gaussian Processes	35
92	3.1	Gaussian processes	35
93	3.2	Chained Gaussian processes	35
94	4	Truncated Generalized Cross Entropy for segmentation	37
95	4.1	Loss functions for multiple annotators	37
96	4.1.1	Generalized Cross Entropy	38
97	4.1.2	Extension to Multiple Annotators	40
98	4.1.3	Reliability Maps and Truncated GCE	40
99	4.2	Proposed Model	42
100	4.2.1	Backbone Architecture	42
101	4.2.2	UNET Architecture	42
102	4.2.3	Reliability Map Branch	43
103	4.2.4	Integration with $TGCE_{SS}$ Loss	43
104	4.2.5	Training Process	43

105	4.3 Experiments	44
106	4.3.1 Dataset	44
107	4.3.2 Metrics	44
108	5 Chained deep learning for image segmentation	47
109	5.1 Introduction	47
110	5.2 Segmentation models	47
111	5.3 Training strategies	47
112	5.4 Evaluation metrics	47
113	5.5 Conclusion	47
114	6 Conclusions	49
115	6.1 Summary	49
116	6.2 Future work	49
117	Bibliography	50

LIST OF FIGURES

119	1-1	Estimation of the tasks and medical image types based on recent	
120		literature review (count of referenced terms).	3
121	1-2	AI and machine learning in medical imaging brief timeline.	4
122	1-3	Example of a histopathological image segmented by multiple	
123		annotators, illustrating variations in label assignment.	7
124	1-4	Summary diagram for problem Statement	10
125	1-5	Proposed framework for the approach in [López-Pérez et al., 2024]. .	17
126	2-1	Comparative Trends of the top two most popular Deep Learning	
127		Frameworks	34
128	4-1	Solution Architecture (mockup)	45

LIST OF TABLES

- 131 **CAD** Computer-Aided Diagnosis 2, 5, 6
132 **CCGP** Correlated Chained Gaussian Processes 18
133 **CCGPMA** Correlated Chained Gaussian Processes for Multiple Annotators 17, 19
134 **CE** Cross Entropy 38
135 **CGP** Chained Gaussian Processes 18
136 **CNN** Convolutional Neural Networks 3, 14, 22, 23
137 **CT** Computed Tomography 11

138 **ELBO** Evidence Lower Bound 18
139 **GCE** Generalized Cross Entropy 38
140 **GCECDL** Generalized Cross-Entropy-based Chained Deep Learning 20, 21

141 **ISS** Image Semantic segmentation 2, 3, 6, 11, 13, 21-23
142 **LF** Latent Function 18
143 **MAE** Mean Absolute Error 38, 39
144 **MITs** Medical Imaging Techniques 1
145 **ML** Machine Learning 11, 21
146 **MV** Majority Voting 11, 12

147 **OCR** Optical Character Recognition 11
148 **PET** Positron Emission Tomography 14
149 **ROI** Region of Interest 2, 6, 15, 29

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

154

155

CHAPTER

156

ONE

157

158

INTRODUCTION

159

1.1 Motivation

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

170 Along with the advances in technologies for medical images acquisition,
171 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

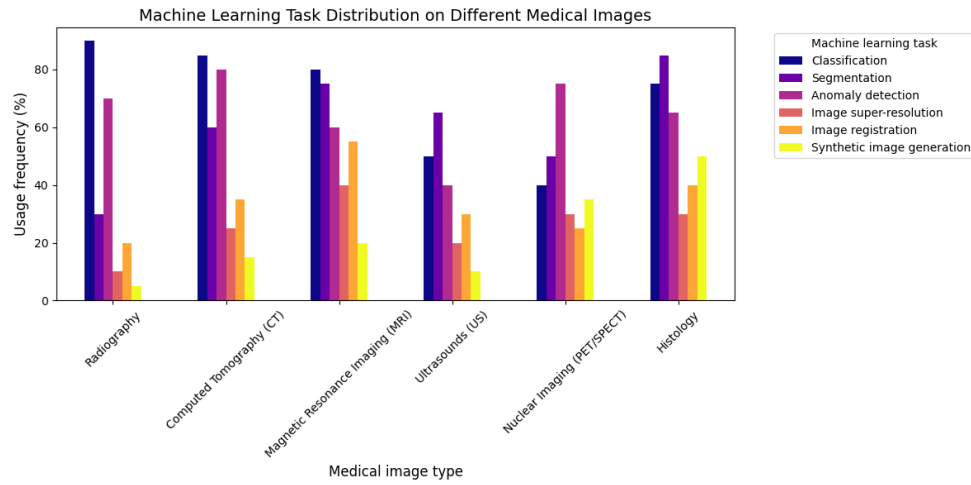


Figure 1-1 Estimation of the tasks and medical image types based on recent literature review (count of referenced terms).

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

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

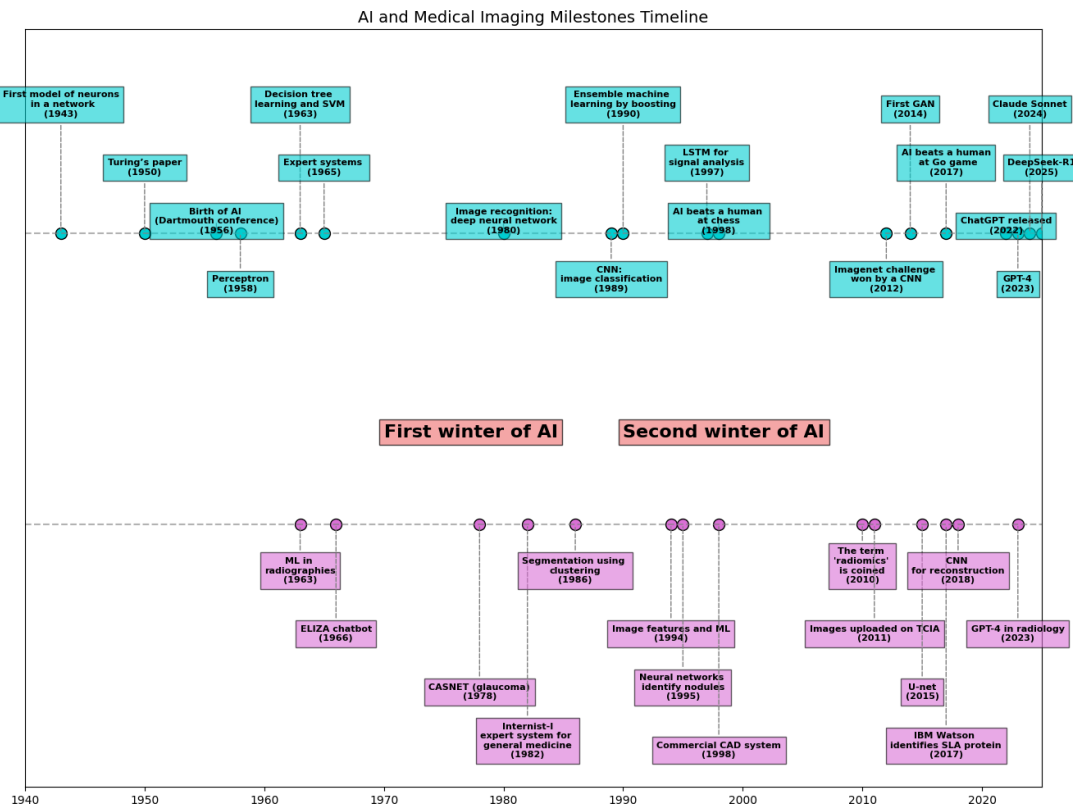


Figure 1-2 AI and machine learning in medical imaging brief timeline.

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

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

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

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

245 A major challenge in histopathological image segmentation arises from the
246 variability in annotations provided by different pathologists. Unlike natural
247 images, where object boundaries are often well-defined, histological structures
248 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

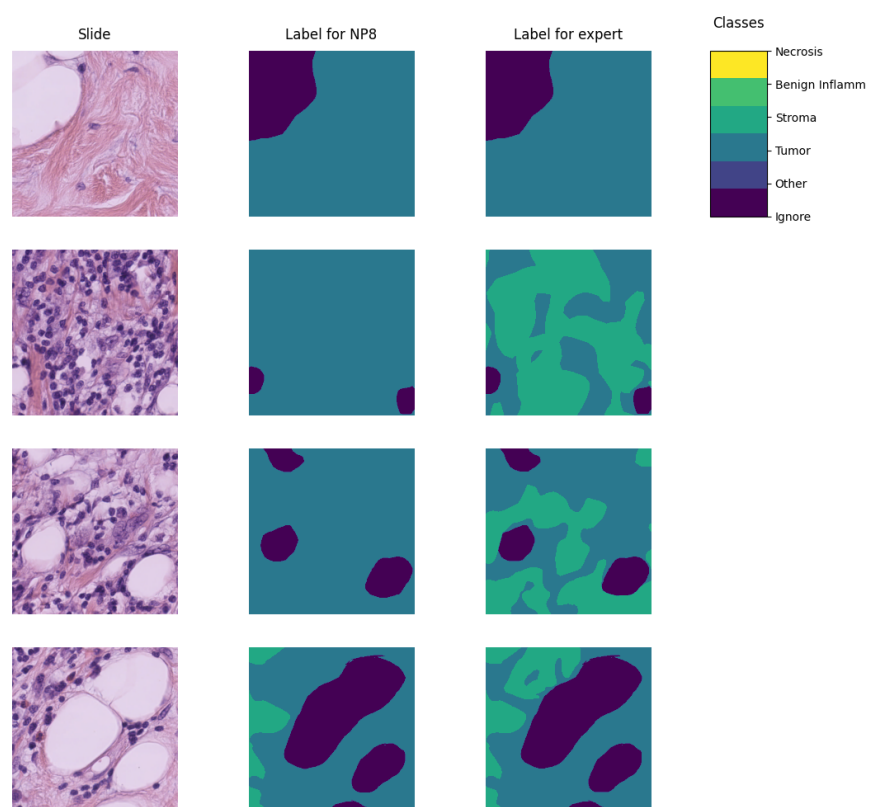


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

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].

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

306 1.2.2 Technical Constraints and Image Quality

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

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

317 1.2.3 Research Question

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

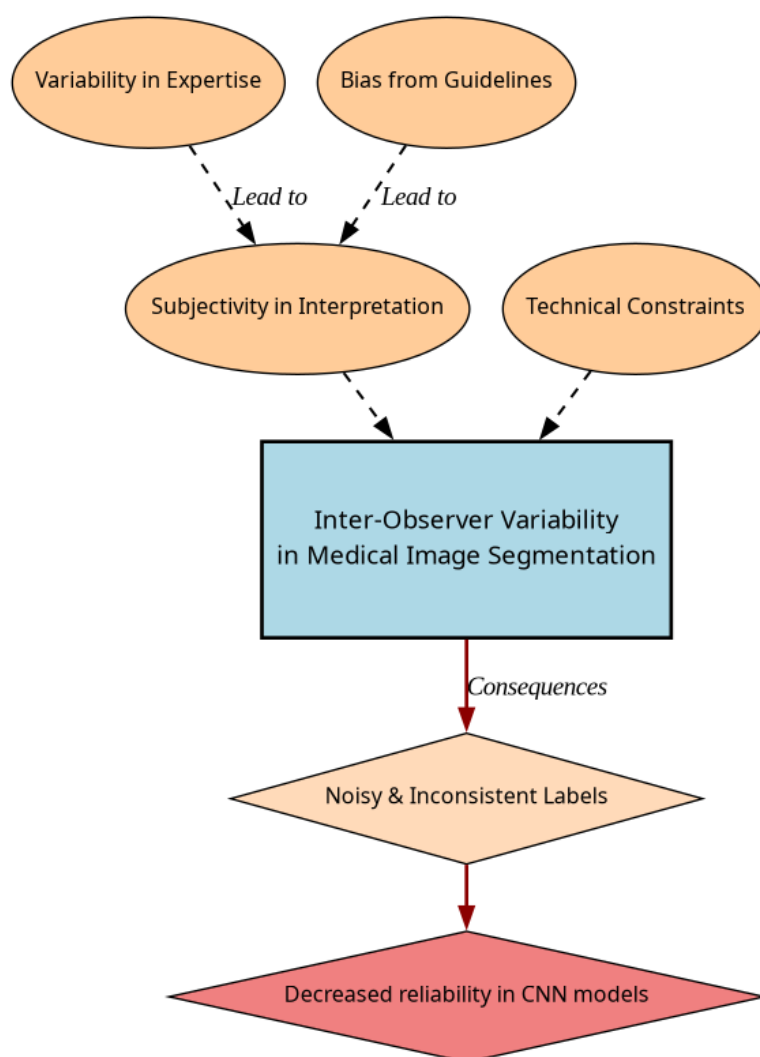


Figure 1-4 Summary diagram for problem Statement

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?

320

1.3 Literature review

321

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

340

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:

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

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

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

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

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

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

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

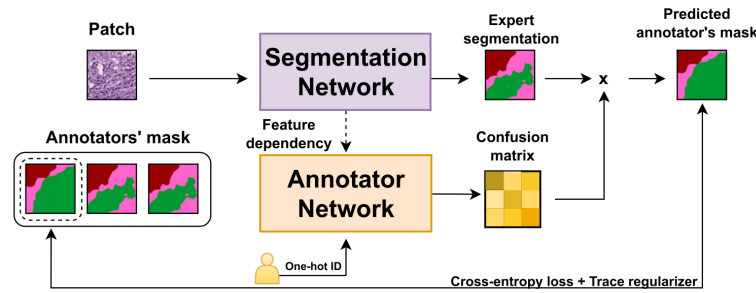


Figure 1-5 Proposed framework for the approach in [López-Pérez et al., 2024].

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.

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 works, however, it does not face the image segmentation problem, which is the main focus of this 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.

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.

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. 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

585

586

CHAPTER

587

TWO

588

589

CONCEPTUAL PRELIMINARIES

590

2.1 Modern concept of digital image

591

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).

592

593

594

595

596

2.1.1 Types of digital images

597

Grayscale images

598

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.

599

600

601

602

603 **Color images**

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

610 **Multispectral images**

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

617 **3D images and volumetric data**

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

623 **2.1.2 Mathematical representations**

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

- 625 • **Sampling:** The process of converting a continuous image into a discrete
626 representation. According to the Nyquist-Shannon sampling theorem, the
627 sampling frequency must be at least twice the highest frequency present in
628 the image to avoid aliasing.
- 629 • **Quantization:** The process of converting continuous intensity values into
630 discrete levels. The number of quantization levels determines the image's
631 bit depth and affects its quality and storage requirements.
- 632 • **Resolution:** The number of pixels per unit length in an image, typically
633 measured in pixels per inch (PPI) or dots per inch (DPI).
- 634 • **Dynamic range:** The ratio between the maximum and minimum measurable
635 light intensities in an image, often expressed in decibels (dB).

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

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

637 where $I(x, y)$ is the digital image, $f(i, j)$ represents the intensity values, and $\delta(x -$
638 $i, y - j)$ is the Kronecker delta function.

639 For color images, the representation extends to:

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

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

641 **2.2 Digital histopathological images**

642 Digital histopathology represents a significant advancement in medical imaging,
643 where traditional glass slides containing tissue samples are digitized using
644 specialized scanning devices. This transformation has revolutionized the field of
645 pathology by enabling remote diagnosis, computer-aided analysis, and digital
646 archiving of tissue samples «CITE».

647 **2.2.1 Whole Slide Imaging (WSI)**

648 Whole Slide Imaging (WSI) is the process of digitizing entire glass slides at high
649 resolution, creating a digital representation that can be viewed, analyzed, and
650 shared electronically. Modern WSI scanners use sophisticated optical systems that
651 capture multiple fields of view at high magnification, which are then stitched
652 together to create a seamless digital image «CITE». These systems typically
653 employ:

- 654 • High-resolution objectives (20x to 40x magnification)
- 655 • Precise motorized stages for accurate slide positioning
- 656 • Automated focus systems to maintain image quality
- 657 • High-quality cameras with large sensor arrays

658 The resulting digital slides can reach sizes of several gigabytes, containing billions
659 of pixels that capture the microscopic details of tissue samples «CITE».

660 2.2.2 Regions of Interest (ROI)

661 In digital histopathology, ROIs are specific areas within a whole slide image that
662 contain diagnostically relevant information. These regions can be:

- 663 • Manually annotated by pathologists
- 664 • Automatically detected using computer vision algorithms
- 665 • Defined based on specific tissue characteristics or abnormalities

666 ROIs are particularly important for:

- 667 • Focusing computational analysis on relevant areas
- 668 • Reducing computational complexity in automated systems
- 669 • Facilitating targeted diagnosis and research
- 670 • Enabling efficient storage and transmission of critical information

671 2.2.3 Staining Techniques

672 Histopathological analysis relies heavily on various staining techniques to enhance
673 the visibility of different tissue components and cellular structures. The choice of
674 staining method depends on the specific diagnostic requirements and the type of
675 tissue being examined.

676 Hematoxylin and Eosin (H&E)

677 Hematoxylin and Eosin (H&E) staining is the most widely used technique in
678 histopathology, particularly in breast cancer diagnosis «CITE». This staining
679 method provides:

- 680 • Hematoxylin: Stains cell nuclei blue/purple, highlighting nuclear morphology
- 681 • Eosin: Stains cytoplasm and extracellular matrix pink/red, revealing tissue
682 architecture

683 The popularity of H&E staining in breast cancer histopathology stems from its ability
684 to:

- 685 • Clearly visualize tumor architecture and growth patterns
- 686 • Distinguish between different types of breast cancer
- 687 • Identify important diagnostic features like nuclear pleomorphism
- 688 • Assess tumor grade and stage

689 Beyond breast cancer, H&E staining is extensively used in various medical
690 specialties including:

- 691 • General pathology
- 692 • Dermatology
- 693 • Gastroenterology
- 694 • Neurology
- 695 • Oncology

696 Special Stains

697 In addition to H&E, various special stains are used for specific diagnostic purposes:

- 698 • **Immunohistochemistry (IHC):** Uses antibodies to detect specific proteins,
699 crucial for:
 - 700 – Subtyping breast cancers (ER, PR, HER2)
 - 701 – Identifying tumor markers
 - 702 – Determining treatment options
- 703 • **Periodic Acid-Schiff (PAS):** Highlights carbohydrates and basement
704 membranes
- 705 • **Masson's Trichrome:** Distinguishes between collagen and muscle fibers
- 706 • **Silver stains:** Used for detecting microorganisms and nerve fibers

707 These specialized staining techniques complement H&E by providing additional
708 diagnostic information that is crucial for accurate diagnosis and treatment
709 planning «CITE».

710 2.3 Deep learning fundamentals

711 Deep learning has emerged as a powerful subset of machine learning,
712 revolutionizing the field of artificial intelligence. Its roots can be traced back to the
713 early development of artificial neural networks in the 1940s and 1950s, with
714 significant milestones including the perceptron in 1958 and the backpropagation
715 algorithm in the 1980s. However, it wasn't until the early 21st century, with the
716 advent of more powerful computational resources and the availability of large
717 datasets, that deep learning truly began to flourish.

718 2.3.1 Learning Paradigms

719 Deep learning systems can be categorized into three main learning paradigms. The
720 most common approach is supervised learning, where models learn from labeled
721 data by mapping inputs to known outputs. This paradigm requires a large amount
722 of labeled training data, which can be expensive and time-consuming to acquire.
723 Semi-supervised learning offers a hybrid approach that leverages both labeled and
724 unlabeled data, proving particularly useful when labeled data is scarce but
725 unlabeled data is abundant. Finally, unsupervised learning enables models to
726 discover patterns and structures from unlabeled data without explicit guidance,
727 making it valuable for tasks like clustering and dimensionality reduction.

728 2.3.2 Architecture and Training

729 Deep learning architectures are characterized by their layered structure, where
730 each layer progressively extracts and transforms features from the input data. The
731 early layers typically focus on low-level feature extraction, such as edges, textures,
732 and basic patterns in the case of image processing. As information flows through
733 the network, middle layers combine these basic features into more complex
734 representations. The final layers perform high-level reasoning and make the
735 ultimate predictions or classifications.

736 The training process relies heavily on the gradient descent algorithm, which
737 iteratively adjusts the model's parameters to minimize a loss function. This loss
738 function serves as a crucial component of the learning process, quantifying how
739 well the model's predictions match the actual targets. By providing a measure of
740 the model's performance, the loss function guides the optimization process,
741 enabling the network to learn meaningful patterns from the training data.

742 2.3.3 Challenges and Solutions

743 Despite their power, deep learning systems face several significant challenges. One
744 of the most prominent issues is overfitting, where models may memorize training
745 data instead of learning generalizable patterns. This challenge is typically
746 addressed through various regularization techniques such as dropout, L1/L2
747 regularization, and early stopping. Another critical challenge is the substantial
748 data requirements; deep learning models often need massive amounts of training
749 data to achieve good performance, which can be a limiting factor in many
750 applications. Additionally, the complex, layered nature of deep learning models
751 makes them difficult to interpret, often referred to as "black boxes." This lack of
752 transparency can be particularly problematic in critical applications where
753 understanding the decision-making process is essential.

754 2.3.4 Deep Learning Frameworks

755 The development of powerful open-source frameworks has significantly
756 accelerated deep learning research and applications. TensorFlow, developed by
757 Google, provides a comprehensive ecosystem for building and deploying machine
758 learning models. PyTorch, created by Facebook's AI Research lab, offers dynamic
759 computation graphs and has become particularly popular in research settings.
760 Caffe, known for its speed and modularity, is widely used in computer vision
761 applications.

762 These frameworks have democratized deep learning by providing efficient
763 implementations of common operations, automatic differentiation for gradient
764 computation, and GPU acceleration for faster training. They also offer pre-trained
765 models and transfer learning capabilities, along with active communities for
766 support and knowledge sharing. The combination of these frameworks with
767 modern hardware has enabled researchers and practitioners to develop

increasingly sophisticated models, pushing the boundaries of what’s possible in artificial intelligence. As shown in Figure 2-1, which presents data from Google Trends over the last five years (as of April 2025), TensorFlow and PyTorch have emerged as the two most prominent frameworks in the deep learning landscape.

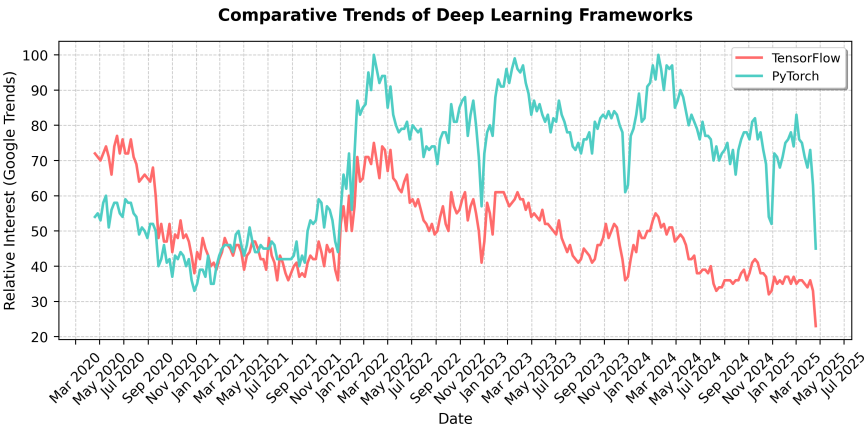


Figure 2-1 Comparative Trends of the top two most popular Deep Learning Frameworks

2.4 Datasets and data sources

773

774

775

776

CHAPTER

THREE

777

CHAINED GAUSSIAN PROCESSES

778

3.1 Gaussian processes

779

3.2 Chained Gaussian processes

780

781

CHAPTER

782

FOUR

783

784

TRUNCATED GENERALIZED CROSS ENTROPY FOR

785

SEGMENTATION

786

4.1 Loss functions for multiple annotators

787

As mentioned in Section ??, 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.

788

789

790

791

792

793

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.

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

- 808 • It is bounded and more robust to outliers.
- 809 • It assigns equal weights to all mistakes regardless of confidence.
- 810 • It is symmetric in softmax based representations.
- 811 • It is more robust to noisy labels but slower to train.

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

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

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

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

816 The GCE loss exhibits several desirable properties:

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

4.1.2 Extension to Multiple Annotators

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

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

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

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

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

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

4.1.3 Reliability Maps and Truncated GCE

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

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

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

- 833 • Weight the contribution of each annotator differently across the image
- 834 • Adapt to varying levels of expertise in different regions
- 835 • Handle cases where annotators might be more reliable in certain areas than
- 836 others

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

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

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

- 841 • The first term weighted by Λ_r represents the GCE loss for regions where the
- 842 annotator is considered reliable
- 843 • The second term weighted by $(1 - \Lambda_r)$ provides a uniform prior for regions
- 844 where the annotator is considered unreliable

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

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

846 4.2 Proposed Model

847 Our proposed model architecture combines the strengths of UNET with a ResNet-
848 34 backbone, specifically designed to work with the $TGCE_{SS}$ loss function. The
849 architecture is illustrated in Figure ??.

850 4.2.1 Backbone Architecture

851 The model employs a pre-trained ResNet-34 as its encoder backbone, leveraging
852 its deep residual learning framework for efficient feature extraction. The choice
853 of ResNet-34 provides several key advantages: efficient feature extraction through
854 residual connections, pre-trained weights that capture rich visual representations,
855 and stable gradient flow during training. We modify the ResNet-34 backbone to
856 serve as the encoder in our UNET architecture by removing the final fully connected
857 layer and utilizing the feature maps from different stages of the network for skip
858 connections.

859 4.2.2 UNET Architecture

860 The UNET architecture follows a traditional encoder-decoder structure with skip
861 connections, where the encoder path implements the ResNet-34 structure. The
862 decoder path employs transposed convolutions for upsampling, creating a
863 symmetrical architecture that effectively captures both high-level and low-level
864 features. The architecture incorporates four downsampling stages in the encoder,
865 corresponding to the ResNet-34 blocks, and four upsampling stages in the decoder.
866 These stages are connected through skip connections that bridge corresponding
867 encoder and decoder stages, allowing the network to preserve fine-grained details.
868 Each convolution operation is followed by batch normalization and ReLU
869 activation to ensure stable training and effective feature learning.

870 4.2.3 Reliability Map Branch

871 A key innovation in our architecture is the parallel branch dedicated to estimating
 872 reliability maps. This branch processes the same encoder features as the main
 873 segmentation path but focuses on learning the confidence of each annotator.
 874 Through a series of 1×1 convolutions, the branch reduces channel dimensions
 875 while maintaining spatial information. The final output consists of R reliability
 876 maps Λ_r , one for each annotator, with values constrained to the $[0, 1]$ range
 877 through a sigmoid activation function. This design allows the network to learn and
 878 adapt to the varying reliability of different annotators across different regions of
 879 the image.

880 4.2.4 Integration with TGCE_{SS} Loss

881 The model produces two distinct outputs: segmentation masks $\mathbf{Y} = f(\mathbf{X}; \theta)$ and
 882 reliability maps $\{\Lambda_r(\mathbf{X}; \theta)\}_{r=1}^R$. These outputs work in tandem with the TGCE_{SS}
 883 loss function described in Section ???. The loss function simultaneously guides the
 884 learning of both the segmentation masks and reliability maps, ensuring that the
 885 model learns to balance the contributions of different annotators based on their
 886 estimated reliability.

887 4.2.5 Training Process

888 The training process begins with the initialization of the ResNet-34 backbone
 889 using pre-trained weights, providing a strong foundation for feature extraction.
 890 The entire network is then trained end-to-end using the Adam optimizer with a
 891 learning rate of 10^{-4} . The TGCE_{SS} loss function plays a crucial role in updating
 892 both the segmentation and reliability branches, ensuring that the model learns to

893 effectively handle multiple annotators' inputs while accounting for their varying
894 reliability.

895 The model's architecture is specifically designed to address the challenges of
896 multi-annotator segmentation. Through the ResNet-34 backbone, it learns robust
897 segmentation features that capture high-level patterns in the data. The UNET's
898 skip connections enable the preservation of fine-grained details, while the parallel
899 reliability branch allows the model to adapt to annotator-specific characteristics.
900 This comprehensive design enables the model to effectively handle multiple
901 annotators' inputs while maintaining high segmentation accuracy and reliability
902 estimation.

903 **4.3 Experiments**

904 **4.3.1 Dataset**

905 **4.3.2 Metrics**

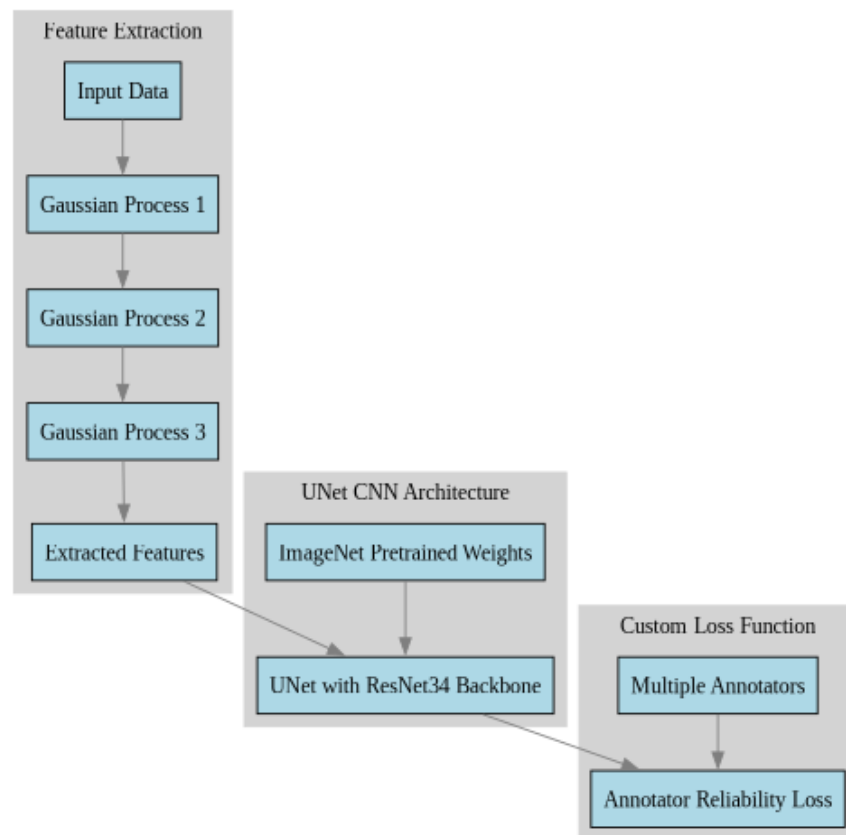


Figure 4-1 Solution Architecture (mockup)

906

907

CHAPTER

908

FIVE

909

910

CHAINED DEEP LEARNING FOR IMAGE

911

SEGMENTATION

912

5.1 Introduction

913

5.2 Segmentation models

914

5.3 Training strategies

915

5.4 Evaluation metrics

916

5.5 Conclusion

917

918

919

920

CHAPTER

SIX

921

CONCLUSIONS

922

6.1 Summary

923

6.2 Future work

BIBLIOGRAPHY

- 925 [Avanzo et al., 2024] Avanzo, M., Stancanella, J., Pirrone, G., Drigo, A., and Retico,
926 A. (2024). The evolution of artificial intelligence in medical imaging: From
927 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
928 Author Joseph Stancanella is employed by Elekta SA. The remaining authors
929 declare no commercial or financial conflicts of interest. (page 3)
- 930 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
931 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
932 (2024). Medical image segmentation review: The success of u-net. *IEEE*
933 *Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
934 (page 2)
- 935 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
936 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in*
937 *Radiology*, 5. (page 8)
- 938 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
939 cost minimization for medical image segmentation using suggestive mixed
940 supervision fully convolutional networks. *CoRR*, abs/1812.11302. (page 3)
- 941 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
942 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
943 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)

- 974 [Julián and Álvarez Meza Andrés Marino, 2023] Julián, G. G. and Álvarez Meza
975 Andrés Marino (2023). A supervised learning framework in the context of
976 multiple annotators. (page 17)
- 977 [Karthikeyan et al., 2023] Karthikeyan, R., McDonald, A., and Mehta, R. (2023).
978 What’s in a label? annotation differences in forecasting mental fatigue using ecg
979 data and seq2seq architectures. (page 9)
- 980 [Kim et al., 2024] Kim, Y., Lee, E., Lee, Y., and Oh, U. (2024). Understanding
981 novice’s annotation process for 3d semantic segmentation task with human-
982 in-the-loop. In *Proceedings of the 29th International Conference on Intelligent User*
983 *Interfaces, IUI ’24*, page 444–454, New York, NY, USA. Association for Computing
984 Machinery. (page 9)
- 985 [Lam and Suen, 1997] Lam, L. and Suen, S. (1997). Application of majority
986 voting to pattern recognition: an analysis of its behavior and performance.
987 *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*,
988 27(5):553–568. (page 11)
- 989 [Lin et al., 2024] Lin, Y., Lian, A., Liao, M., and Yuan, S. (2024). Bcdnet: A fast
990 residual neural network for invasive ductal carcinoma detection. (page 6)
- 991 [López-Pérez et al., 2023] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A. D.,
992 Molina, R., and Katsaggelos, A. K. (2023). Crowdsourcing segmentation
993 of histopathological images using annotations provided by medical students.
994 In Juarez, J. M., Marcos, M., Stiglic, G., and Tucker, A., editors, *Artificial*
995 *Intelligence in Medicine*, pages 245–249, Cham. Springer Nature Switzerland.
996 (pages 5, 6, 8, and 12)
- 997 [Lu et al., 2023] Lu, X., Ratcliffe, D., Kao, T.-T., Tikhonov, A., Litchfield, L., Rodger,
998 C., and Wang, K. (2023). Rethinking quality assurance for crowdsourced multi-
999 roi image segmentation. *Proceedings of the AAAI Conference on Human Computation*
1000 *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 xvii, 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)

- 1029 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018).
1030 Ms-net: Mixed-supervision fully-convolutional networks for full-resolution
1031 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
1032 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer*
1033 *Assisted Intervention – MICCAI 2018*, pages 379–387, Cham. Springer International
1034 Publishing. (page 5)
- 1035 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
1036 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
1037 *Sciences*, 378(2166):20190061. (page 3)
- 1038 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
1039 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
1040 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
1041 Curran Associates, Inc. (page 12)
- 1042 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
1043 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
1044 Chained deep learning using generalized cross-entropy for multiple annotators
1045 classification. *Sensors*, 23(7). (page 20)
- 1046 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
1047 truth and performance level estimation (staple): an algorithm for the validation
1048 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
1049 (page 12)
- 1050 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
1051 Advances in medical image segmentation: A comprehensive review of traditional,
1052 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 1053 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
1054 Semi-supervised gland segmentation via feature-enhanced contrastive learning
1055 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
1056 pages 1–11. (page 2)

- 1057 [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Cicarrelli, O.,
1058 Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground
1059 truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell,
1060 R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*,
1061 volume 33, pages 15750–15762. Curran Associates, Inc. (page 16)
- 1062 [Zhang and Sabuncu, 2018] Zhang, Z. and Sabuncu, M. R. (2018). Generalized
1063 cross entropy loss for training deep neural networks with noisy labels.
1064 (pages 38 and 39)
- 1065 [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S.,
1066 Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M.
1067 (2021). A review of deep learning in medical imaging: Imaging traits, technology
1068 trends, case studies with progress highlights, and future promises. *Proceedings of*
1069 *the IEEE*, 109(5):820–838. (pages 1 and 2)
- 1070 [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L.
1071 (2024). Image quality evaluation in deep-learning-based ct noise reduction using
1072 virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical*
1073 *Physics*, 51(8):5399–5413. (page 9)