



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

Brandon Lotero Londoño
2025

ABSTRACT

49 PENDING

50 **Keywords:** PENDING

53 PENDIENTE

54 **Palabras clave:** PENDIENTE

LIST OF FIGURES

LIST OF TABLES

- 60 **CAD** Computer-Aided Diagnosis 2, 5, 6
- 61 **CCGP** Correlated Chained Gaussian Processes 18
- 62 **CCGPMA** Correlated Chained Gaussian Processes for Multiple Annotators 18
- 63 **CGP** Chained Gaussian Processes 18
- 64 **CNN** Convolutional Neural Networks 3, 14, 20, 22
- 65 **CT** Computed Tomography 12
- 66 **ELBO** Evidence Lower Bound 19
- 67 **GCECDL** Generalized Cross-Entropy-based Chained Deep Learning 19, 20
- 68 **ISS** Image Semantic segmentation 2, 3, 6, 11, 13, 20–22
- 69 **LF** Latent Function 18
- 70 **MITs** Medical Imaging Techniques 1
- 71 **ML** Machine Learning 11
- 72 **MV** Majority Voting 11, 12
- 73 **OCR** Optical Character Recognition 11
- 74 **PET** Positron Emission Tomography 14
- 75 **ROI** Region of Interest 2, 6
- 76 **SLFM** Semi-Parametric Latent Factor Model 18
- 77 **SS** Semantic segmentation 3
- 78 **STAPLE** Simultaneous Truth and Performance Level Estimation 12–14
- 79 **WSI** Whole Slide Imaging 1, 5, 6, 8, 14, 16

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

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. [?]. 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 [?].

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. [?]. 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 [?].

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 [?]. 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 «CITE». 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 ??.

For solving the different requirements of tasks in medical images, a variety of computational techniques have been developed [?]. 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 application of advanced statistical tools and machine learning algorithms like Support Vector

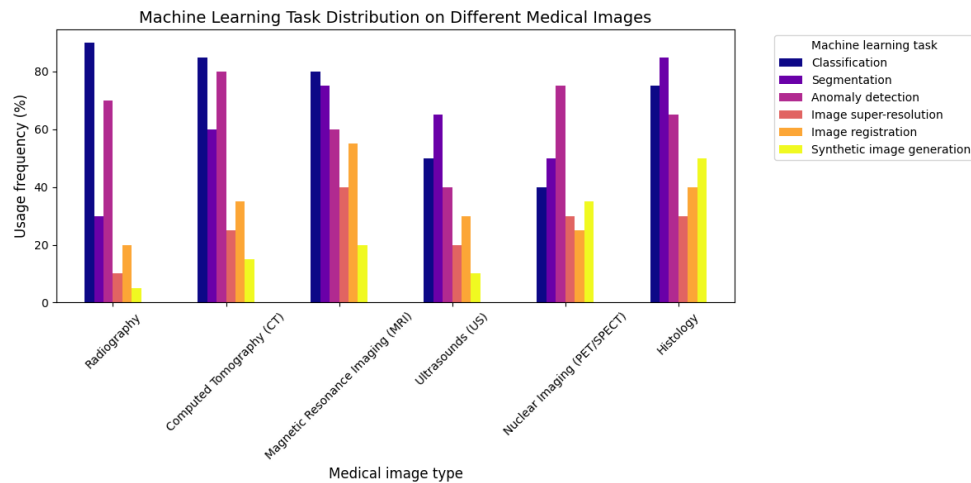


Figure 1-1 Estimation of the tasks and medical image types based on recent literature review (count of referenced terms).

Machines, Decision Trees, and SGD Neural Networks [?]. The coevolution of advances in medical image acquisition, computational power (i.e. Moore's law) and statistical/mathematical techniques have led to a convergence for merging state of the art algorithms with medical imaging [?]. Figure ?? shows a brief timeline of coevolution between some conspicuous advances in computational pattern recognition and its medical applications in different scopes (besides medical imaging) [?].

Convolutional Neural Networks (CNN) have been widely used in Semantic segmentation (SS) tasks, as they have outperformed traditional machine learning algorithms in this task for both medical and non medical images [?] [?]. However, most CNN architectures are deep, which imply a necessity of a large amount of data to train them. This introduces a problem since both the acquisition and annotation of medical images are expensive and time-consuming processes. This is especially true for ISS tasks, as they require pixel-level annotations, which is taxing in terms of cost, time and logistics involved [?]. Other fashions face this problem through less expensive annotation strategies like bounding boxes or anatomical landmarks for being used in a semi-supervised strategy [?].

Many medical images datasets however, contain a high variability in class sizes

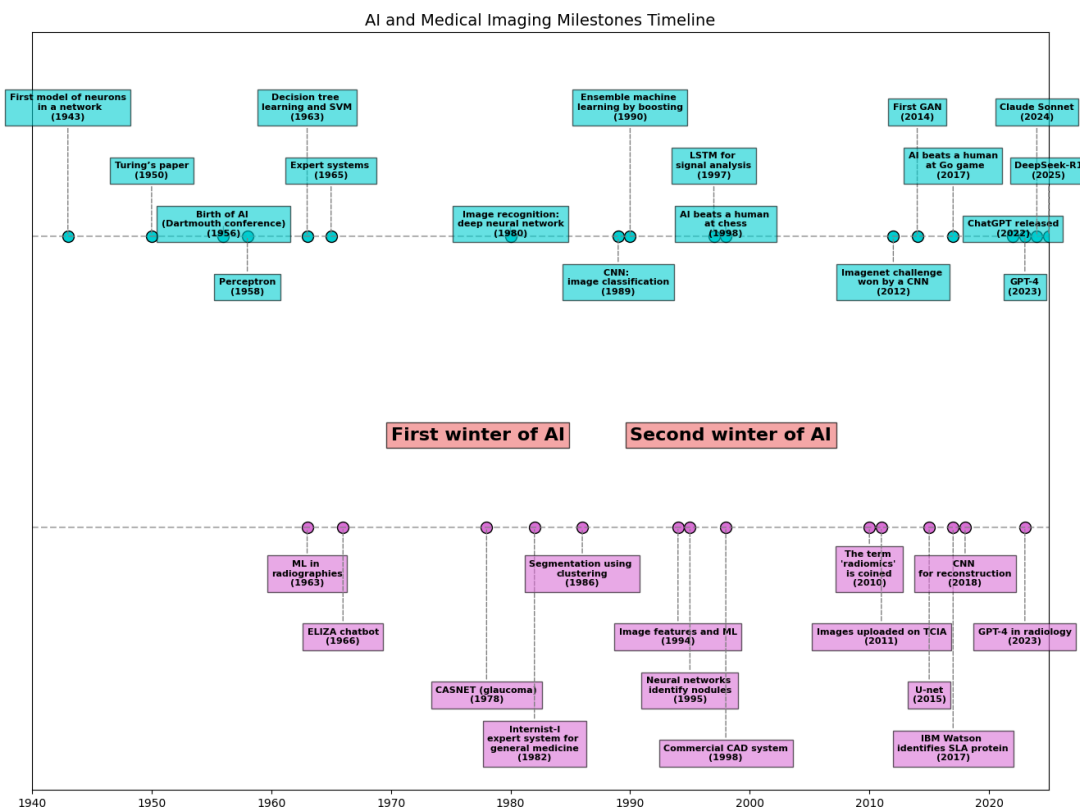


Figure 1-2 AI and machine learning in medical imaging brief timeline.

146 and variations in colors, which is specially noticeable in histopathological images
147 because of the usage of different staining and other factors which can affect the
148 color of the images. This variability can lead to a significant loss of efficiency of
149 machine learning models when using a mixed supervision strategy, as the model
150 can be biased towards the most common classes or colors in the dataset [?].

151 This is where other solutions arise to tackle the problem of the weak image
152 annotation while maintaining low costs. One of these solutions is crowdsourcing
153 strategy, which consists of having multiple annotators labeling the same image,
154 and then combining the labels to obtain a consensus label [?]. This strategy can
155 lead to a labeling cost reduction when different levels of expertise are combined,
156 since the crowd may be composed of both experts and laymen, being the latter
157 less expensive to hire [?].

158 Recently, diagnosis, prognosis and treatment of cancer have heavily relied on
159 histopathology, where tissue samples are obtained through biopsies or surgical
160 resections and critical information that helps pathologists determine the presence
161 and severity of malignancies [?]. The segmentation of histopathological images
162 enables precise identification of structures such as nuclei, glands, and tumors,
163 which are essential for assessing disease progression and treatment response [?].
164 Accurate segmentation is particularly crucial in digital pathology, where
165 whole-slide images (WSI) are analyzed using AI-powered CAD systems to support
166 clinical decision-making [?].

167 A major challenge in histopathological image segmentation arises from the
168 variability in annotations provided by different pathologists. Unlike natural
169 images, where object boundaries are often well-defined, histological structures
170 may have ambiguous borders, leading to inconsistencies among annotators [?].
171 Because of this, crowdsourcing labeling is one of the most popular approaches, as
172 illustrated in Figure ??, an example of how histopathological images are segmented
173 by multiple experts, showing some variations in label assignment ¹. These

¹obtained from a real world Triple Negative Breast Cancer (TNBC) dataset published in [?]

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 [?]. 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 [?].

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 [?]. 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.

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 [?]. Despite the help of these tools, the labeling process in WSI can have high costs, as it requires long hours of work from

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

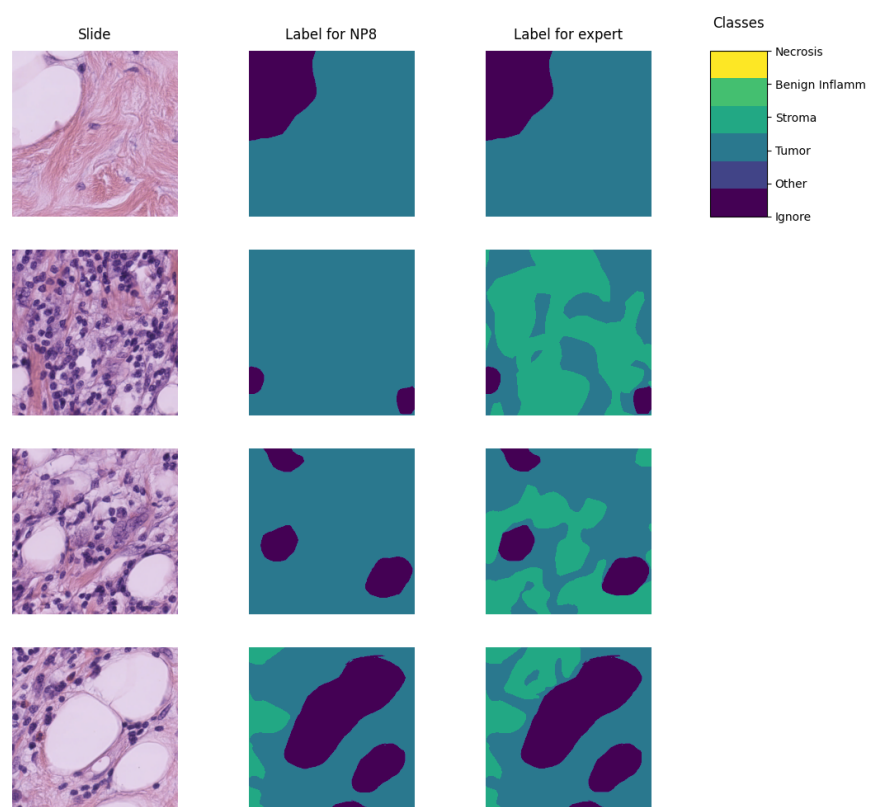


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

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 [?]. These inconsistencies are mostly represented in the subsections ?? and ??.

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 [?]. 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 [?].

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 [?]. 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 [?].

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

1.2.2 Technical Constraints and Image Quality

Technical constraints in medical imaging, such as resolution differences, noise levels, and contrast variations, can significantly impact segmentation accuracy. Lower-resolution images may obscure fine structures, leading to inconsistencies in boundary delineation [?].

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

1.2.3 Research Question

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

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?

1.3 Literature review

Certainly, in general Machine Learning (ML) classification tasks³ where multiple annotators are involved, Majority Voting (MV) is by far the simplest possible

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

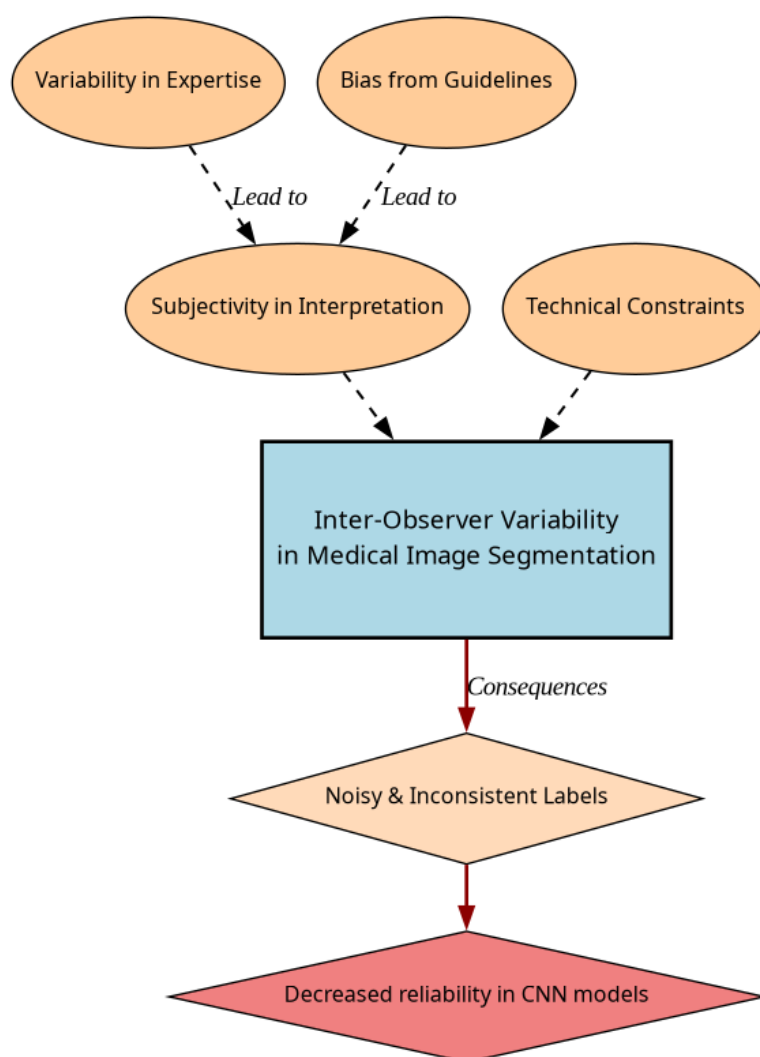


Figure 1-4 Summary diagram for problem Statement

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 [?], 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 [?], 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 [?], 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 ??, 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 [?], 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 [?] 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:

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

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

297 The popularity and performance of STAPLE has led to its usage in modern
 298 applications medical image, 3d spatial images due to its assumption of decision
 299 space being based on voxel-wise decisions, like the authors in [?] which applied
 300 the algorithm on Positron Emission Tomography (PET) images. Other authors still
 301 rely heavily on STAPLE for setting a ground truth consensus for histopathological
 302 images, like [?].

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

312 Finally, this work contemplates **STAPLE** as useful for ground truth estimation given
313 the existence of multiple labelers for an input **WSI**, but not that useful for providing
314 annotations of structures on new and unlabeled images, hence being a good support
315 for other methods.

316 **U-shaped CNNs**

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

325 The U-Net architecture consists of a symmetric encoder-decoder structure with
326 skip connections. The encoder path progressively reduces spatial dimensions
327 while increasing feature channels through a series of convolutional and
328 max-pooling layers, capturing high-level semantic information. The decoder path
329 uses transposed convolutions to gradually recover spatial resolution while
330 reducing feature channels. Skip connections between corresponding encoder and
331 decoder layers preserve fine-grained details by concatenating high-resolution
332 features from the encoder with upsampled features in the decoder, enabling
333 precise localization of structures. The architecture overview can be seen in figure
334 ??.

335 **U-Net based approaches**

336 In [?] two networks are trained for delivering a final segmentation. One network is
337 trained to estimate the annotators reliability and another one is trained to segment

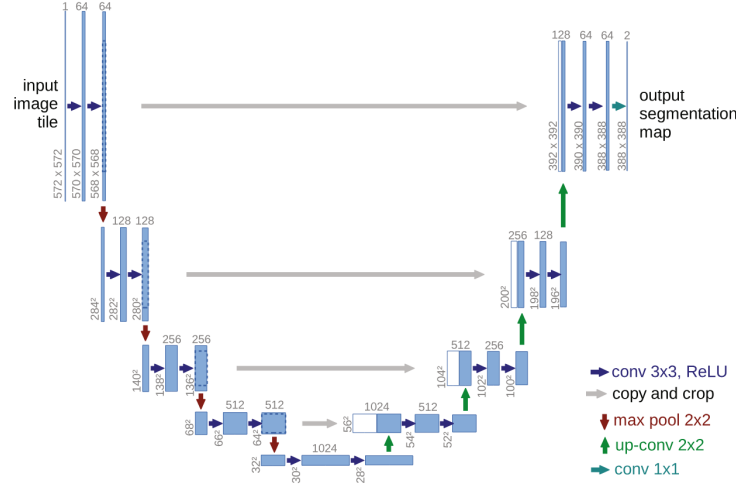


Figure 1-5 Original U-Net architecture.

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:

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

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

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

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

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

367 The authors propose a loss function involving cross-entropy and a trace based
368 regularization on the reliability map, originally proposed in [?] which combined,
369 looks like:

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

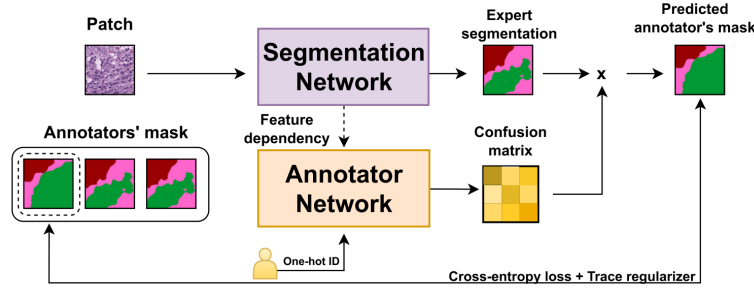


Figure 1-6 Proposed framework for the approach in [?].

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

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

377 Bayesian models

378 Bayesian approaches are a good choice for handling label noise and uncertainty in
 379 the labelers. In [?] the authors propose a novel approach from Gaussian Processes
 380 to model the relationship between the annotators' reliability and the input data,
 381 while also preserving the interdependencies among the annotators. This is
 382 achieved by introducing **Correlated Chained Gaussian Processes for Multiple**
 383 **Annotators (CCGPMA)**, a framework based on the well known **Chained Gaussian**
 384 **Processes (CGP)**. CGP on itself cannot consider inter-annotator dependencies, thus,
 385 the authors introduce the **Correlated Chained Gaussian Processes (CCGP)** to model
 386 correlations between the GP latent functions, which are supposed to be generated
 387 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,

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.

1.3.2 Strategies for handling low-quality images

The problem of low-quality images 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.

A more recent approach was proposed by [?], 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 image regions. Their experiments on multiple datasets demonstrated superior predictive performance compared to state-of-the-art methods, particularly in cases where annotations were highly inconsistent.

This strategy is especially relevant for handling low-quality medical images, where expert annotations may be inconsistent, and traditional consensus-based approaches fail to account for varying expertise levels. By leveraging deep learning with robust noise-handling loss functions, the reliability of segmentation models can be significantly improved.

1.4 Aims

With the mentioned considerations in section ?? 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

463

464

465

466

CHAPTER

TWO

467

468

TRUNCATED GENERALIZED CROSS ENTROPY FOR SEGMENTATION

469

2.1 Proposed Loss Function

470

471

472

473

The development of our proposed loss function stems from the need to handle multiple annotators' segmentation masks while accounting for their varying reliability across different regions of the image. We begin by examining the foundation of our approach: the Generalized Cross Entropy (GCE).

474

2.1.1 Generalized Cross Entropy

475

476

477

The Generalized Cross Entropy (GCE) loss function was introduced as a robust alternative to the standard cross-entropy loss, particularly effective in handling noisy labels. The GCE loss for a single annotator can be expressed as:

$$GCE(\mathbf{Y}, f(\mathbf{X}; \theta)) = \frac{1}{q} \left(1 - \sum_{k=1}^K \mathbf{Y}_k f(\mathbf{X}; \theta)_k^q \right) \quad (2-1)$$

where $q \in (0, 1)$ is a hyperparameter that controls the truncation level, \mathbf{Y} is the ground truth label, and $f(\mathbf{X}; \theta)$ is the model's prediction. The GCE loss exhibits several desirable properties:

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

2.1.2 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:

$$\mathbf{X} \in \mathbb{R}^{W \times H}, \{\mathbf{Y}_r \in \{0, 1\}^{W \times H \times K}\}_{r=1}^R; \quad \mathbf{\Psi} \in [0, 1]^{W \times H \times K} = f(\mathbf{X}) \quad (2-2)$$

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

$$f : \mathbb{R}^{W \times H} \rightarrow [0, 1]^{W \times H \times K} \quad (2-3)$$

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

$$\mathbf{Y}_r[w, h, :] \mathbf{1}_k^\top = 1; \quad w \in W, h \in H \quad (2-4)$$

2.1.3 Reliability Maps and Truncated GCE

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

$$\left\{ \Lambda_r(\mathbf{X}; \theta) \in [0, 1]^{W \times H} \right\}_{r=1}^R \quad (2-5)$$

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:

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

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

$$\begin{aligned} TGCE_{SS}(\mathbf{Y}_r, f(\mathbf{X}; \theta) | \Lambda_r(\mathbf{X}; \theta)) = & \mathbb{E}_r \left\{ \mathbb{E}_{w,h} \left\{ \Lambda_r(\mathbf{X}; \theta) \circ \mathbb{E}_k \left\{ \mathbf{Y}_r \circ \left(\frac{\mathbf{1}_{W \times H \times K} - f(\mathbf{X}; \theta)^{\circ q}}{q} \right); k \in K \right\} + \right. \right. \\ & \left. \left. (\mathbf{1}_{W \times H} - \Lambda_r(\mathbf{X}; \theta)) \circ \left(\frac{\mathbf{1}_{W \times H} - (\frac{1}{k} \mathbf{1}_{W \times H})^{\circ q}}{q} \right); w \in W, h \in H \right\}; r \in R \right\} \end{aligned} \quad (2-6)$$

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

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

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

$$L(\mathbf{Y}_r[n], f(\mathbf{X}[n]; \theta)|_r(\mathbf{X}[n]; \theta)) = \frac{1}{N} \sum_n^N TGCE_{SS}(\mathbf{Y}_r[n], f(\mathbf{X}[n]; \theta)|_r(\mathbf{X}[n]; \theta)) \quad (2-7)$$

2.2 Proposed Model

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

2.2.1 Backbone Architecture

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

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

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

2.2.2 UNET Architecture

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

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

2.2.3 Reliability Map Branch

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

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

2.2.4 Integration with TGCE_{SS} Loss

The model outputs two components:

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

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

545 **2.2.5 Training Process**

546 The training process involves:

- 547 • Initializing the ResNet-34 backbone with pre-trained weights
- 548 • Training the entire network end-to-end
- 549 • Using the Adam optimizer with a learning rate of 10^{-4}
- 550 • Applying the $TGCE_{SS}$ loss to update both the segmentation and reliability
- 551 branches

552 The model's architecture allows it to:

- 553 • Learn robust segmentation features through the ResNet-34 backbone
- 554 • Capture fine-grained details through UNET's skip connections
- 555 • Adapt to annotator reliability through the parallel reliability branch
- 556 • Handle multiple annotators' inputs effectively

BIBLIOGRAPHY

- 558 [Avanzo et al., 2024] Avanzo, M., Stancanella, J., Pirrone, G., Drigo, A., and Retico,
559 A. (2024). The evolution of artificial intelligence in medical imaging: From
560 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
561 Author Joseph Stancanella is employed by Elekta SA. The remaining authors
562 declare no commercial or financial conflicts of interest. (page 3)
- 563 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
564 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
565 (2024). Medical image segmentation review: The success of u-net. *IEEE*
566 *Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
567 (page 2)
- 568 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
569 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in*
570 *Radiology*, 5. (page 8)
- 571 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
572 cost minimization for medical image segmentation using suggestive mixed
573 supervision fully convolutional networks. CoRR, abs/1812.11302. (page 3)
- 574 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
575 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
576 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 17)
- [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)

- 662 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018).
663 Ms-net: Mixed-supervision fully-convolutional networks for full-resolution
664 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
665 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer*
666 *Assisted Intervention – MICCAI 2018*, pages 379–387, Cham. Springer International
667 Publishing. (page 5)
- 668 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
669 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
670 *Sciences*, 378(2166):20190061. (page 3)
- 671 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
672 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
673 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
674 Curran Associates, Inc. (page 12)
- 675 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
676 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
677 Chained deep learning using generalized cross-entropy for multiple annotators
678 classification. *Sensors*, 23(7). (page 19)
- 679 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
680 truth and performance level estimation (staple): an algorithm for the validation
681 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
682 (page 12)
- 683 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
684 Advances in medical image segmentation: A comprehensive review of traditional,
685 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 686 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
687 Semi-supervised gland segmentation via feature-enhanced contrastive learning
688 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
689 pages 1–11. (page 2)

- 690 [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Ciccarrelli, O.,
691 Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground
692 truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell,
693 R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*,
694 volume 33, pages 15750–15762. Curran Associates, Inc. (page 17)
- 695 [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S.,
696 Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M.
697 (2021). A review of deep learning in medical imaging: Imaging traits, technology
698 trends, case studies with progress highlights, and future promises. *Proceedings of*
699 *the IEEE*, 109(5):820–838. (pages 1 and 2)
- 700 [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L.
701 (2024). Image quality evaluation in deep-learning-based ct noise reduction using
702 virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical*
703 *Physics*, 51(8):5399–5413. (page 9)