



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

44

ACKNOWLEDGEMENTS

45 PENDING

46

Brandon Lotero Londoño

47

2025

48

ABSTRACT

49 PENDING

50 **Keywords:** PENDING

51

RESUMEN

52

53 PENDIENTE

54 **Palabras clave:** PENDIENTE

55

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)	30
84	2.2.3 Staining Techniques	30
85	2.3 Deep learning fundamentals	32
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	35
91	2.4.1 Datasets with emulated noisy annotations	35
92	3 Chained Gaussian Processes	37
93	3.1 Gaussian processes	37
94	3.2 Chained Gaussian processes	37
95	4 Truncated Generalized Cross Entropy for segmentation	39
96	4.1 Loss functions for multiple annotators	39
97	4.1.1 Generalized Cross Entropy	40
98	4.1.2 Extension to Multiple Annotators	42
99	4.1.3 Reliability Maps and Truncated GCE	42
100	4.2 Proposed Model	44
101	4.2.1 Backbone Architecture	44
102	4.2.2 UNET Architecture	44
103	4.2.3 Reliability Map Branch	45
104	4.2.4 Integration with TGCE_{SS} Loss	45

105	4.2.5 Training Process	45
106	4.3 Experiments	46
107	4.3.1 Dataset	46
108	4.3.2 Metrics	46
109	5 Chained deep learning for image segmentation	49
110	5.1 Introduction	49
111	5.2 Segmentation models	49
112	5.3 Training strategies	49
113	5.4 Evaluation metrics	49
114	5.5 Conclusion	49
115	6 Conclusions	51
116	6.1 Summary	51
117	6.2 Future work	51
118	Bibliography	52

LIST OF FIGURES

120	1-1	Estimation of the tasks and medical image types based on recent 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 annotators, illustrating variations in label assignment.	7
123	1-4	Summary diagram for problem Statement	10
124	1-5	Proposed framework for the approach in [López-Pérez et al., 2024]. .	17
125	2-1	Histology evolution timeline. (Image from [Mazzarini et al., 2021]). .	28
126	2-2	(Above) Whole slide imaging system by Omnyx for slide digitization. (Below) Comprehensive digital pathology interface from Omnyx designed to streamline pathologists' diagnostic workflow. (From [Farahani et al., 2015]).	29
127	2-3	Different staining techniques obtained from multi-stain breast cancer dataset [Weitz et al., 2023]. (a) shows H&E, (b) ER, (c) HER2, (d) Ki67 and (e) PGR. (f) shows an example of a WSI that was excluded since it contains multiple tissue sections.	31
128	2-4	Comparative Trends of the top two most popular Deep Learning Frameworks, apparently, tendency was switched to PyTorch since 2022	34
129	4-1	Solution Architecture (mockup)	47

LIST OF TABLES

ABBREVIATIONS

- ¹⁴² **CAD** Computer-Aided Diagnosis 2, 5, 6
¹⁴³ **CCGP** Correlated Chained Gaussian Processes 18
¹⁴⁴ **CCGPMA** Correlated Chained Gaussian Processes for Multiple Annotators 17, 19
¹⁴⁵ **CE** Cross Entropy 40
¹⁴⁶ **CGP** Chained Gaussian Processes 18
¹⁴⁷ **CNN** Convolutional Neural Networks 3, 14, 22, 23
¹⁴⁸ **CT** Computed Tomography 11
¹⁴⁹ **ELBO** Evidence Lower Bound 18
¹⁵⁰ **GCE** Generalized Cross Entropy 40
¹⁵¹ **GCECDL** Generalized Cross-Entropy-based Chained Deep Learning 20, 21
¹⁵² **ISS** Image Semantic segmentation 2, 3, 6, 11, 13, 21–23, 35
¹⁵³ **LF** Latent Function 18
¹⁵⁴ **MAE** Mean Absolute Error 40, 41
¹⁵⁵ **MITs** Medical Imaging Techniques 1
¹⁵⁶ **ML** Machine Learning 11, 21
¹⁵⁷ **MV** Majority Voting 11, 12
¹⁵⁸ **OCR** Optical Character Recognition 11
¹⁵⁹ **PET** Positron Emission Tomography 14
¹⁶⁰ **ROI** Region of Interest 2, 6, 15, 30
¹⁶¹ **SLFM** Semi-Parametric Latent Factor Model 18
¹⁶² **SS** Semantic segmentation 3
¹⁶³ **STAPLE** Simultaneous Truth and Performance Level Estimation 12–14, 35
¹⁶⁴ **WSI** Whole Slide Imaging 1, 5, 6, 8, 15

165

166 CHAPTER

167

168 ONE

169

INTRODUCTION

170

1.1 Motivation

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

181 Along with the advances in technologies for medical images acquisition,
182 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

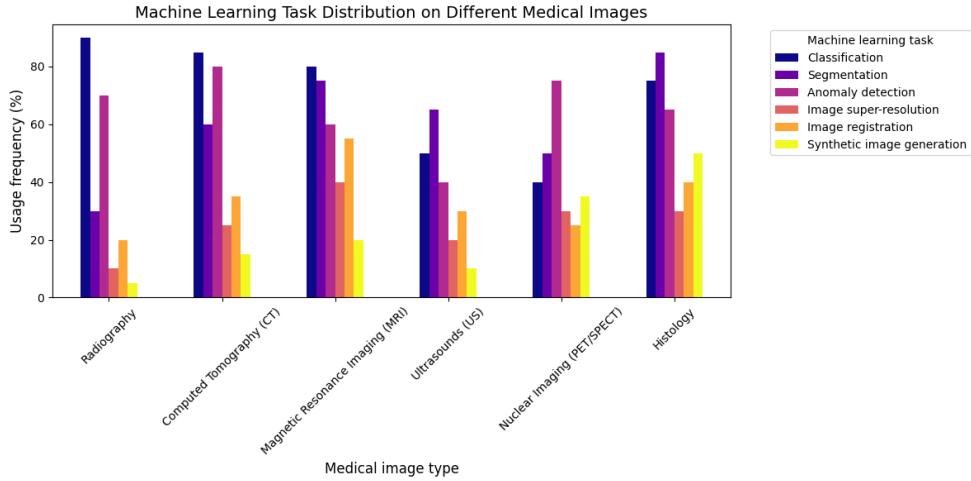


Figure 1-1 Estimation of the tasks and medical image types based on recent literature review (count of referenced terms).

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

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

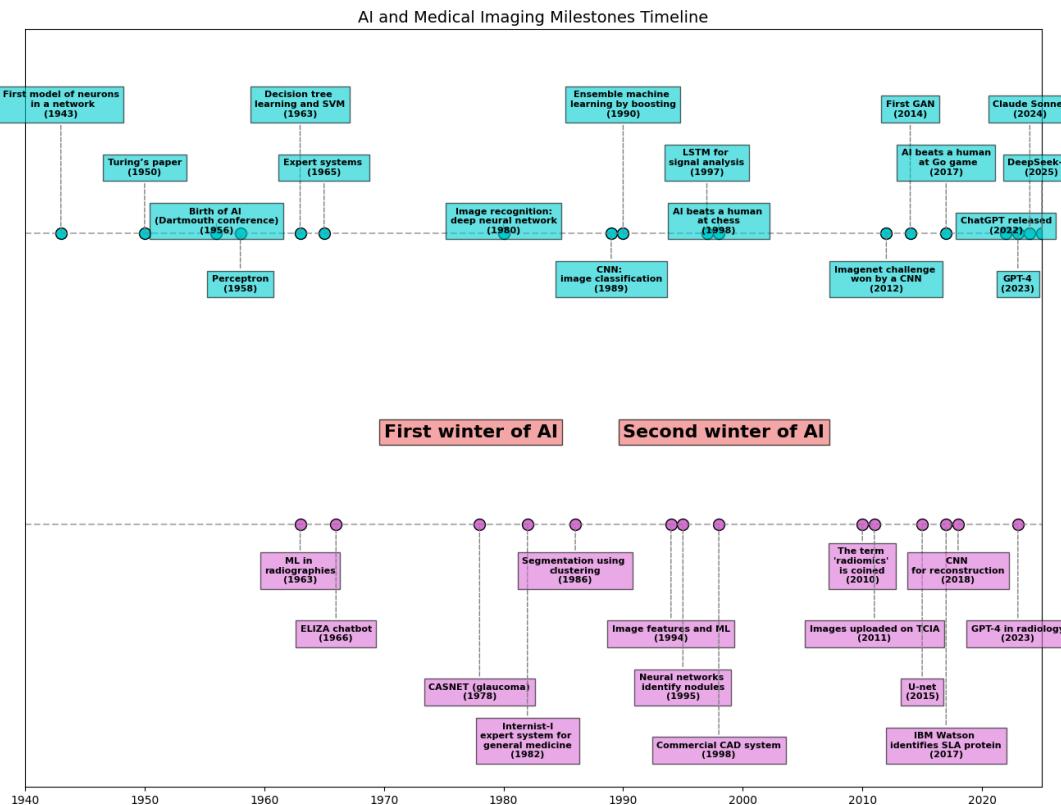


Figure 1-2 AI and machine learning in medical imaging brief timeline.

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

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

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

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

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

[López-Pérez et al., 2023]. Because of this, crowdsourcing labeling is one of the most popular approaches, as illustrated in Figure 1-3, an example of how histopathological images are segmented by multiple experts, showing some variations in label assignment¹. These discrepancies highlight the need for models that can handle annotation uncertainty effectively. Leveraging crowdsourcing strategies and machine learning techniques that infer annotator reliability can enhance segmentation performance while reducing costs.

1.2 Problem Statement

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

ISS task consists of assigning a label to each pixel in an image according to the object it belongs to. Accurate segmentation is essential for the development of CAD systems, as it allows the identification of regions of interest (ROI) in the images, which can be used to detect and classify diseases and hence, treatment planning [Sarvamangala and Kulkarni, 2022]. However, modern computational solutions for ISS tasks involve the use of deep learning, which mostly rely large amounts of labeled data to train the models on supervised learning techniques. This means that the model is trained on a dataset with ground-truth labels, which are assumed to be correct and consistent across all samples. In practice, this assumption is often violated due to the high technical complexity of labeling these segments².

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

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

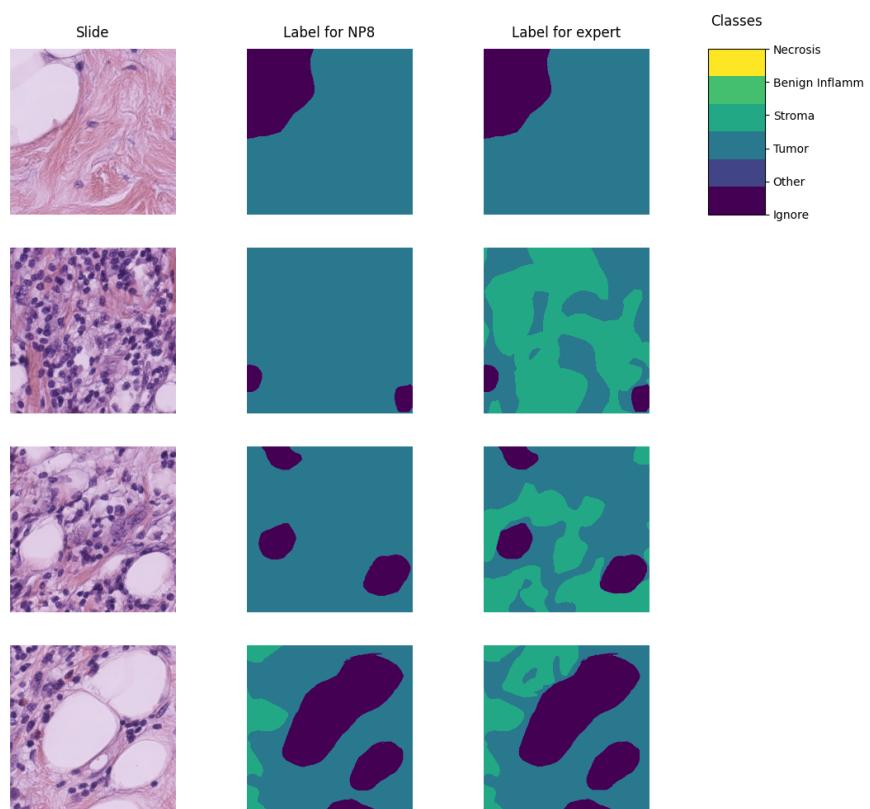


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

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

295 1.2.1 Variability in Expertise Levels

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

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

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

317 **1.2.2 Technical Constraints and Image Quality**

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

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

328 **1.2.3 Research Question**

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

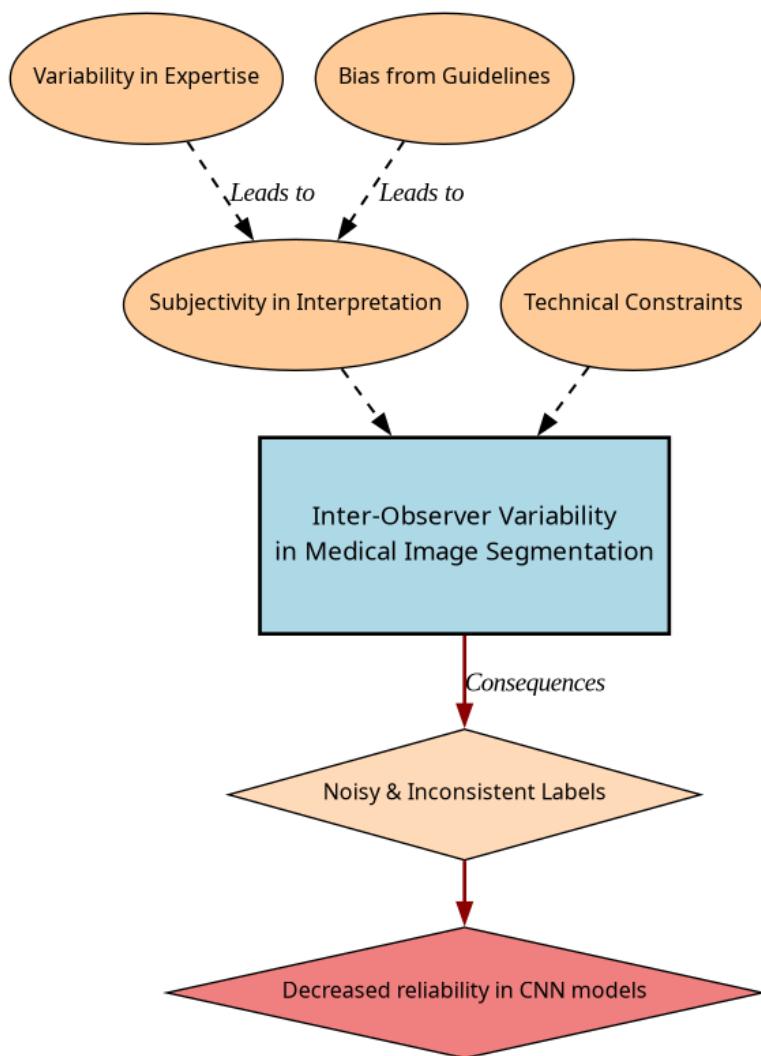


Figure 1-4 Summary diagram for problem Statement

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?

331

332 1.3 Literature review

333 Certainly, in general Machine Learning (ML) classification tasks ³ where multiple
334 annotators are involved, Majority Voting (MV) is by far the simplest possible
335 approach to implement. This concept was born multiple times and divergently in
336 multiple fields, but it was described as relevant for ML and pattern recognition
337 labeling for classification in [Lam and Suen, 1997], in which the approach is
338 exposed as simple, yet powerful. The authors describe the MV as a method that
339 can be used to improve the accuracy of classification tasks by combining the labels
340 of multiple annotators. The method is based on the assumption that the majority
341 vote of the annotators is more likely to be correct than the vote of a single
342 annotator. The authors also describe the method as a straightforward way to
343 improve the accuracy of classification tasks without the need for complex
344 algorithms or additional data. The authors also prove this method to deliver very
345 similar results to more complicated approaches (Bayesian, logistic regression,
346 fuzzy integral, and neural network) in the particular task of Optical Character
347 Recognition (OCR). Despite its simplicity, modern solutions for delivering accurate
348 medical image segmentation models still rely on Majority Voting at some stage,
349 like [Elnakib et al., 2020], which uses a majority voting strategy for delivering a
350 final output based on the labels of multiple models (VGG16-Segnet, Resnet-18 and
351 Alexnet) in Computed Tomography (CT) images for Liver Tumor Segmentation, or

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

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

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

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

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

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

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

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

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

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

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

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

410 U-shaped CNNs

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

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

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

428 **U-Net based approaches**

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

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

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

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

443 The authors then assume the annotator network to deliver a reliability map
 444 $\{\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})\}_{r \in R_n}$ with different dimensions:

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

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

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

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

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

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

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

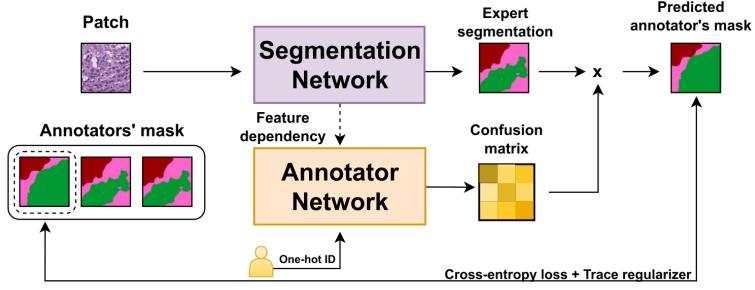


Figure 1-5 Proposed framework for the approach in [López-Pérez et al., 2024].

- 463 Being \mathbb{I} the indicator function, CE the cross-entropy loss, and λ the regularization
 464 parameter.
- 465 When evaluated on a Triple Negative Breast Cancer dataset, this approach achieves
 466 a Dice coefficient of 0.7827, outperforming STAPLE (0.7039) and matching expert-
 467 supervised performance (0.7723). The CR image reliability modeling proved most
 468 effective, as CR pixel, while potentially offering finer-grained reliability estimation,
 469 requires significantly more training data.
- 470 Despite the decent performance of the approach, solving the problem of multiple
 471 labelers with two networks can be overwhelming for the optimization process,
 472 requiring large amounts of annotated data to properly codify the annotators
 473 spatial reliabilities, which could be managed by a single model with an appropriate
 474 loss function.

475 Bayesian models

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

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

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

487 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 $k_q : \mathcal{X} \times \mathcal{X} \rightarrow$
 488 \mathbb{R} being a kernel function, and $w_{j,q} \in \mathbb{R}$ is a combination coefficient ($Q \in \mathbb{N}$). This
 489 leads to a joint distribution of the form:

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

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

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

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

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

502 CCGPMA on itself proposes a good approach for handling label noise and
503 uncertainty in the labelers for regression and classification tasks, while also
504 preserving the interdependencies among the annotators, however, it does not face
505 the image segmentation problem, which is the main focus of this works, however,
506 it does not face the image segmentation problem, which is the main focus of this
507 work. Besides, handling so many latent functions during the optimization process
508 is computationally expensive, making it on itself infeasible for large and high
509 resolution datasets.

510 1.3.2 Facing noisy annotations and low-quality data

511 The problem of low-quality data and noisy annotations has been tackled with
512 various strategies. One such approach is the use of deep learning models that
513 incorporate loss functions designed to mitigate the effects of unreliable labels.
514 Traditional methods such as Majority Voting (MV) or Expectation-Maximization
515 (EM) have been widely used for aggregating multiple annotators' inputs. However,
516 they assume a homogeneous reliability of annotators, which may not hold in
517 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.

518 **Loss functions in deep learning models**

519 Loss functions are fundamental components in deep learning models that quantify
520 how well a model's predictions match the ground truth. They serve as the
521 objective function that guides the learning process by measuring the discrepancy
522 between predicted and actual values. In classification tasks, the most common
523 loss functions are Cross-Entropy (CE) and Mean Absolute Error (MAE). CE is
524 particularly effective for classification as it heavily penalizes confident but wrong
525 predictions, though it can be sensitive to noisy labels. MAE, on the other hand, is
526 more robust to outliers and assigns equal weights to all mistakes, but typically
527 requires more training iterations. For image segmentation tasks, specialized loss
528 functions have been developed to handle the unique challenges of pixel-wise
529 classification. The Dice loss, which measures the overlap between predicted and
530 ground truth regions, is widely used in medical image segmentation. More
531 recently, the Generalized Cross Entropy (GCE) loss has emerged as a robust
532 alternative that combines the benefits of both CE and MAE, allowing for better
533 handling of noisy labels through a tunable parameter that controls sensitivity to
534 outliers. In multi-annotator scenarios, where multiple experts provide potentially
535 inconsistent segmentations, novel loss functions like the Truncated Generalized
536 Cross Entropy for Semantic Segmentation (TGCE_{SS}) have been developed to
537 account for varying annotator reliability across different image regions. These loss
538 functions are crucial for training accurate segmentation models, especially in
539 medical imaging where precise delineation of anatomical structures is essential for
540 diagnosis and treatment planning.

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

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

- 546 The approach effectively mitigates the impact of noisy labels by using a
 547 noise-robust loss function, balancing Mean Absolute Error (MAE) and Categorical
 548 Cross-Entropy (CE). Unlike prior approaches, GCECDL accounts for the
 549 dependencies among annotators while encoding their non-stationary behavior
 550 across different data samples. Their experiments on multiple datasets
 551 demonstrated superior predictive performance compared to state-of-the-art
 552 methods, particularly in cases where annotations were highly inconsistent.
- 553 The strategy of the authors effectively unlocks the potential of ML models to handle
 554 low-quality data and noisy annotations, but it is bounded to classifications tasks
 555 only, not being by itself applicable to segmentation tasks. The TGCE equation for
 556 handling multiple annotators is defined as:

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

- 557 where $\tilde{\lambda}_x$ represents the annotator reliability, \tilde{C} is a constant, q is a parameter that
 558 controls the balance between MAE and CE behavior, \mathbf{y} is the annotation vector, and
 559 $f(\mathbf{x})$ is the model prediction. This approach is more deeply discussed in chapter 4.

560 1.4 Aims

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

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

572 **1.4.1 General Aim**

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

579 **1.4.2 Specific Aims**

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

589 1.5 Outline and Contributions

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

- 592 • A python package for using the proposed loss function in CNN models for ISS
593 tasks in medical images. ¹¹
- 594 • Datasets mapping as lazy loaders for the proposed loss function. ¹²
- 595 • 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

596

CHAPTER

597

598

TWO

599

600

CONCEPTUAL PRELIMINARIES

601

2.1 Modern concept of digital image

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

607

2.1.1 Types of digital images

608

Grayscale images

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

614 **Color images**

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

621 **Multispectral images**

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

628 **3D images and volumetric data**

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

634 **2.1.2 Mathematical representations**

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

- 636 • **Sampling:** The process of converting a continuous image into a discrete
637 representation. According to the Nyquist-Shannon sampling theorem, the
638 sampling frequency must be at least twice the highest frequency present in
639 the image to avoid aliasing.
- 640 • **Quantization:** The process of converting continuous intensity values into
641 discrete levels. The number of quantization levels determines the image's
642 bit depth and affects its quality and storage requirements.
- 643 • **Resolution:** The number of pixels per unit length in an image, typically
644 measured in pixels per inch (PPI) or dots per inch (DPI).
- 645 • **Dynamic range:** The ratio between the maximum and minimum measurable
646 light intensities in an image, often expressed in decibels (dB).

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

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

648 where $I(x, y)$ is the digital image, $f(i, j)$ represents the intensity values, and $\delta(x -$
649 $i, y - j)$ is the Kronecker delta function.

650 For color images, the representation extends to:

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

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

652 2.2 Digital histopathological images

653 Digital histopathology represents a significant advancement in medical imaging,
 654 where traditional glass slides containing tissue samples are digitized using
 655 specialized scanning devices. This transformation has revolutionized the field of
 656 pathology by enabling remote diagnosis, computer-aided analysis, and digital
 657 archiving of tissue samples [Amgad et al., 2019]. This process has evolved
 658 significantly over the past few decades, as shown in Figure 2-1.

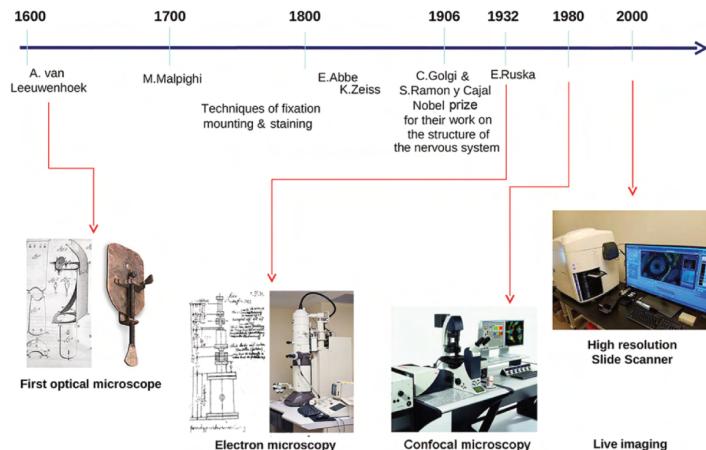


Figure 2-1 Histology evolution timeline. (Image from [Mazzarini et al., 2021]).

659 2.2.1 Whole Slide Imaging (WSI)

660 Whole Slide Imaging (WSI) is the process of digitizing entire glass slides at high
 661 resolution, creating a digital representation that can be viewed, analyzed, and
 662 shared electronically. Modern WSI scanners use sophisticated optical systems that
 663 capture multiple fields of view at high magnification, which are then stitched
 664 together to create a seamless digital image [Hu et al., 2025]. These systems
 665 incorporate high-resolution objectives with magnifications ranging from 20x to
 666 40x, precise motorized stages for accurate slide positioning, automated focus

systems to maintain image quality, and high-quality cameras equipped with large sensor arrays. The resulting digital slides can reach sizes of several gigabytes, containing billions of pixels that capture the microscopic details of tissue samples [Hu et al., 2025]. Figure 2-2 shows a whole slide imaging system by Omnyx for slide digitization and a comprehensive digital pathology interface from Omnyx designed to streamline pathologists' diagnostic workflow [Farahani et al., 2015].

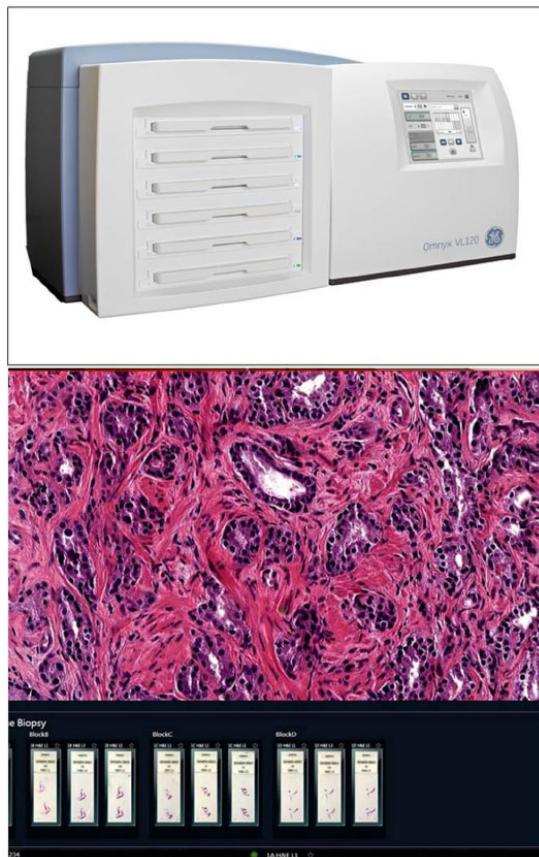


Figure 2-2 (Above) Whole slide imaging system by Omnyx for slide digitization. (Below) Comprehensive digital pathology interface from Omnyx designed to streamline pathologists' diagnostic workflow. (From [Farahani et al., 2015]).

673 2.2.2 Regions of Interest (ROI)

674 In digital histopathology, ROIs are specific areas within a whole slide image that
675 contain diagnostically relevant information. These regions can be manually
676 annotated by pathologists, automatically detected using computer vision
677 algorithms, or defined based on specific tissue characteristics or abnormalities.
678 The importance of ROIs lies in their ability to focus computational analysis on
679 relevant areas, reduce computational complexity in automated systems, facilitate
680 targeted diagnosis and research, and enable efficient storage and transmission of
681 critical information.

682 2.2.3 Staining Techniques

683 Histopathological analysis relies heavily on various staining techniques to enhance
684 the visibility of different tissue components and cellular structures. The choice of
685 staining method depends on the specific diagnostic requirements and the type of
686 tissue being examined.

687 Hematoxylin and Eosin (H&E)

688 Hematoxylin and Eosin (H&E) staining is the most widely used technique in
689 histopathology, particularly in breast cancer diagnosis [Pan et al., 2021]. This
690 staining method provides essential visualization through two components:
691 hematoxylin, which stains cell nuclei blue/purple to highlight nuclear morphology,
692 and eosin, which stains cytoplasm and extracellular matrix pink/red to reveal
693 tissue architecture.

694 The popularity of H&E staining in breast cancer histopathology stems from its
695 ability to clearly visualize tumor architecture and growth patterns, distinguish
696 between different types of breast cancer, identify important diagnostic features

like nuclear pleomorphism, and assess tumor grade and stage. Beyond breast cancer, H&E staining finds extensive application across various medical specialties including general pathology, dermatology, gastroenterology, neurology, and oncology.

Special Stains

In addition to H&E, various special stains are used for specific diagnostic purposes. Immunohistochemistry (IHC) uses antibodies to detect specific proteins, playing a crucial role in subtyping breast cancers. Key IHC stains include Estrogen Receptor (ER) staining for detecting estrogen receptors, Progesterone Receptor (PGR) staining for assessing progesterone receptor status, Human Epidermal Growth Factor Receptor 2 (HER2) staining for evaluating HER2 protein expression, and Ki67 staining for measuring cellular proliferation rates. These markers are particularly crucial in breast cancer diagnosis and treatment planning, as they help determine the molecular subtype of the cancer and guide personalized therapeutic approaches. Other specialized stains include Periodic Acid-Schiff (PAS) for highlighting carbohydrates and basement membranes, Masson's Trichrome for distinguishing between collagen and muscle fibers, and silver stains for detecting microorganisms and nerve fibers. These specialized staining techniques complement H&E by providing additional diagnostic information that is crucial for accurate diagnosis and treatment planning [Weitz et al., 2023].

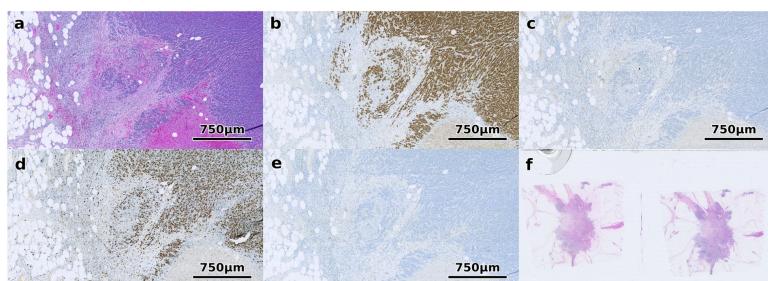


Figure 2-3 Different staining techniques obtained from multi-stain breast cancer dataset [Weitz et al., 2023]. (a) shows H&E, (b) ER, (c) HER2, (d) Ki67 and (e) PGR. (f) shows an example of a WSI that was excluded since it contains multiple tissue sections.

⁷¹⁷ 2.3 Deep learning fundamentals

⁷¹⁸ Deep learning has emerged as a powerful subset of machine learning,
⁷¹⁹ revolutionizing the field of artificial intelligence. Its roots can be traced back to the
⁷²⁰ early development of artificial neural networks in the 1940s and 1950s, with
⁷²¹ significant milestones including the perceptron in 1958 and the backpropagation
⁷²² algorithm in the 1980s. However, it wasn't until the early 21st century, with the
⁷²³ advent of more powerful computational resources and the availability of large
⁷²⁴ datasets, that deep learning truly began to flourish.

⁷²⁵ 2.3.1 Learning Paradigms

⁷²⁶ Deep learning systems can be categorized into three main learning paradigms. The
⁷²⁷ most common approach is supervised learning, where models learn from labeled
⁷²⁸ data by mapping inputs to known outputs. This paradigm requires a large amount
⁷²⁹ of labeled training data, which can be expensive and time-consuming to acquire.
⁷³⁰ Semi-supervised learning offers a hybrid approach that leverages both labeled and
⁷³¹ unlabeled data, proving particularly useful when labeled data is scarce but
⁷³² unlabeled data is abundant. Finally, unsupervised learning enables models to
⁷³³ discover patterns and structures from unlabeled data without explicit guidance,
⁷³⁴ making it valuable for tasks like clustering and dimensionality reduction.

⁷³⁵ 2.3.2 Architecture and Training

⁷³⁶ Deep learning architectures are characterized by their layered structure, where
⁷³⁷ each layer progressively extracts and transforms features from the input data. The
⁷³⁸ early layers typically focus on low-level feature extraction, such as edges, textures,
⁷³⁹ and basic patterns in the case of image processing. As information flows through

740 the network, middle layers combine these basic features into more complex
741 representations. The final layers perform high-level reasoning and make the
742 ultimate predictions or classifications.

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

749 **2.3.3 Challenges and Solutions**

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

761 **2.3.4 Deep Learning Frameworks**

762 The development of powerful open-source frameworks has significantly
763 accelerated deep learning research and applications. TensorFlow, developed by

764 Google, provides a comprehensive ecosystem for building and deploying machine
 765 learning models. PyTorch, created by Facebook's AI Research lab, offers dynamic
 766 computation graphs and has become particularly popular in research settings.
 767 Caffe, known for its speed and modularity, is widely used in computer vision
 768 applications.

769 These frameworks have democratized deep learning by providing efficient
 770 implementations of common operations, automatic differentiation for gradient
 771 computation, and GPU acceleration for faster training. They also offer pre-trained
 772 models and transfer learning capabilities, along with active communities for
 773 support and knowledge sharing. The combination of these frameworks with
 774 modern hardware has enabled researchers and practitioners to develop
 775 increasingly sophisticated models, pushing the boundaries of what's possible in
 776 artificial intelligence. As shown in Figure 2-4, which presents data from Google
 777 Trends over the last five years (as of April 2025), TensorFlow and PyTorch have
 778 emerged as the two most prominent frameworks in the deep learning landscape.

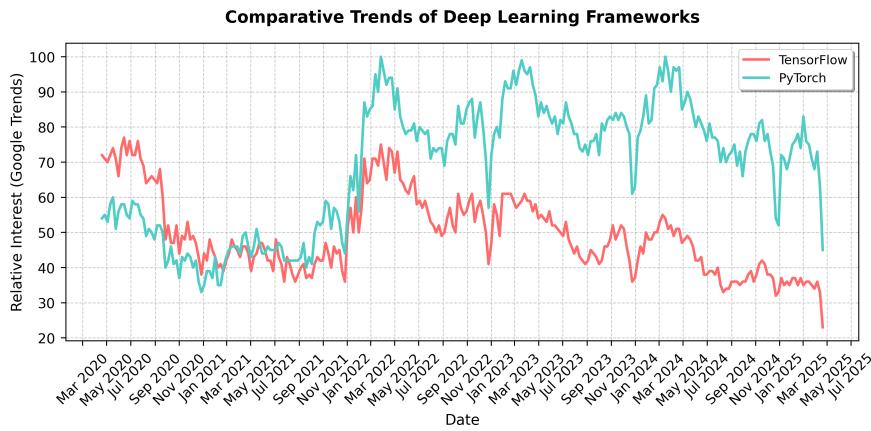


Figure 2-4 Comparative Trends of the top two most popular Deep Learning Frameworks, apparently, tendency was switched to PyTorch since 2022

779 2.4 Datasets and data sources

780 Throughout the development of this work, multiple datasets were used for
781 evaluation of ISS models. The common elements of all these datasets are that they
782 contain RGB images and are crowdsourced with multiple labelers, where not
783 necessarily all labeler label all images.

784 As it has been mentioned in Chapter 1, the main goal of this work is mainly
785 focused on crowdsourced histopathology images semantic segmentation, however,
786 these datasets present the following challenges:

- 787 • Distribution of segmentation labels is not uniform across the image, since
788 some tissues and structures are more common than others.
- 789 • Visualization of performance of the models in debug time (like per epoch
790 analysis) is not simple for non experts in the subject, which makes it hard to
791 evaluate whether the model is overfitting or not at a glance.

792 For these reasons, multiple datasets were created in the pursuit of an initial
793 evaluation of performance of the models against more traditional and familiar
794 images before the focus on histopathology images. Once a decent performance in
795 metrics like Dice coefficient was achieved, the focus was shifted to histopathology
796 images and further tunings on the models were performed if needed.

797 In any case, both the emulated noisy annotations datasets and the histopathology
798 datasets somehow contained ground truth aggregation, either from the original
799 source (in the case of emulated noisy annotations), the expert annotation (if
800 available) or from the aggregation of multiple labelers ¹.

801 2.4.1 Datasets with emulated noisy annotations

802 A challenge arises for the creation of emulated noisy annotations datasets,

¹STAPLE in the case of histopathology datasets with no expert annotations available

803

804

CHAPTER

805

THREE

806

807

CHAINED GAUSSIAN PROCESSES

808

3.1 Gaussian processes

809

3.2 Chained Gaussian processes

810

811 CHAPTER

812

FOUR

813

814

TRUNCATED GENERALIZED CROSS ENTROPY FOR SEGMENTATION

815

816

4.1 Loss functions for multiple annotators

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

824 4.1.1 Generalized Cross Entropy

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

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

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

829 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, $\mathbf{1}^\top \mathbf{y} =$
 830 $\mathbf{1}^\top f(\mathbf{x}) = 1$, $\mathbf{1} \in \{1\}^K$ being an all-ones vector. In addition, the MAE loss can be
 831 rewritten for softmax outputs, yielding:

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

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

833 The CE is characterized by the following properties:

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

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

- 838 • It is bounded and more robust to outliers.
- 839 • It assigns equal weights to all mistakes regardless of confidence.
- 840 • It is symmetric in softmax based representations.
- 841 • It is more robust to noisy labels but slower to train.

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

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

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

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

846 The GCE loss exhibits several desirable properties:

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

850 4.1.2 Extension to Multiple Annotators

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

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

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

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

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

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

857 4.1.3 Reliability Maps and Truncated GCE

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

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

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

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

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

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

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

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

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

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

876 4.2 Proposed Model

877 Our proposed model architecture combines the strengths of UNET with a ResNet-
878 34 backbone, specifically designed to work with the TGCE_{SS} loss function. The
879 architecture is illustrated in Figure ??.

880 4.2.1 Backbone Architecture

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

889 4.2.2 UNET Architecture

890 The UNET architecture follows a traditional encoder-decoder structure with skip
891 connections, where the encoder path implements the ResNet-34 structure. The
892 decoder path employs transposed convolutions for upsampling, creating a
893 symmetrical architecture that effectively captures both high-level and low-level
894 features. The architecture incorporates four downsampling stages in the encoder,
895 corresponding to the ResNet-34 blocks, and four upsampling stages in the decoder.
896 These stages are connected through skip connections that bridge corresponding
897 encoder and decoder stages, allowing the network to preserve fine-grained details.
898 Each convolution operation is followed by batch normalization and ReLU
899 activation to ensure stable training and effective feature learning.

900 4.2.3 Reliability Map Branch

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

910 4.2.4 Integration with TGCE_{SS} Loss

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

917 4.2.5 Training Process

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

923 effectively handle multiple annotators' inputs while accounting for their varying
924 reliability.

925 The model's architecture is specifically designed to address the challenges of
926 multi-annotator segmentation. Through the ResNet-34 backbone, it learns robust
927 segmentation features that capture high-level patterns in the data. The UNET's
928 skip connections enable the preservation of fine-grained details, while the parallel
929 reliability branch allows the model to adapt to annotator-specific characteristics.
930 This comprehensive design enables the model to effectively handle multiple
931 annotators' inputs while maintaining high segmentation accuracy and reliability
932 estimation.

933 **4.3 Experiments**

934 **4.3.1 Dataset**

935 **4.3.2 Metrics**

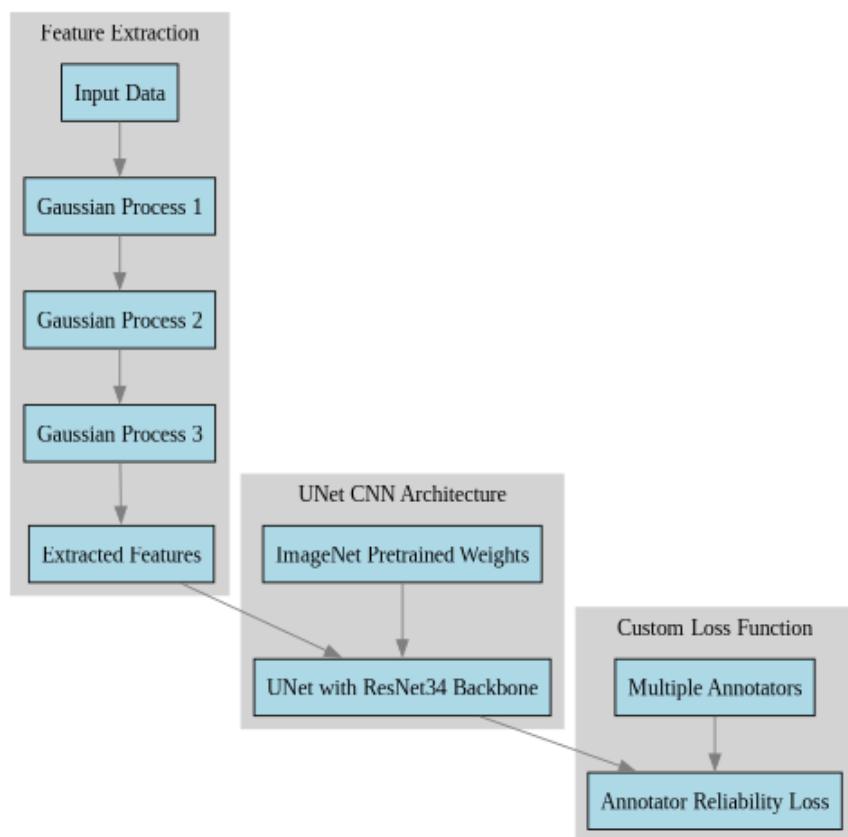


Figure 4-1 Solution Architecture (mockup)

936

937

CHAPTER

938

939

FIVE

940

CHAINED DEEP LEARNING FOR IMAGE SEGMENTATION

941

942

5.1 Introduction

943

5.2 Segmentation models

944

5.3 Training strategies

945

5.4 Evaluation metrics

946

5.5 Conclusion

947

948

CHAPTER

949

950

SIX

951

CONCLUSIONS

952

6.1 Summary

953

6.2 Future work

BIBLIOGRAPHY

954

- 955 [Amgad et al., 2019] Amgad, M., Elfandy, H., Hussein, H., Atteya, L. A., Elsebaie,
956 M. A. T., Abo Elnasr, L. S., Sakr, R. A., Salem, H. S. E., Ismail, A. F., Saad,
957 A. M., Ahmed, J., Elsebaie, M. A. T., Rahman, M., Ruhban, I. A., Elgazar, N. M.,
958 Alagha, Y., Osman, M. H., Alhusseiny, A. M., Khalaf, M. M., Younes, A.-A. F.,
959 Abdulkarim, A., Younes, D. M., Gadallah, A. M., Elkashash, A. M., Fala, S. Y., Zaki,
960 B. M., Beezley, J., Chittajallu, D. R., Manthey, D., Gutman, D. A., and Cooper, L.
961 A. D. (2019). Structured crowdsourcing enables convolutional segmentation of
962 histology images. *Bioinformatics*, 35(18):3461–3467. (page 28)
- 963 [Avanzo et al., 2024] Avanzo, M., Stancanello, J., Pirrone, G., Drigo, A., and Retico,
964 A. (2024). The evolution of artificial intelligence in medical imaging: From
965 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
966 Author Joseph Stancanello is employed by Elekta SA. The remaining authors
967 declare no commercial or financial conflicts of interest. (page 3)
- 968 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
969 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
970 (2024). Medical image segmentation review: The success of u-net. *IEEE
971 Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
972 (page 2)

- 973 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
974 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in*
975 *Radiology*, 5. (page 8)
- 976 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
977 cost minimization for medical image segmentation using suggestive mixed
978 supervision fully convolutional networks. *CoRR*, abs/1812.11302. (page 3)
- 979 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
980 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
981 literature review. (page 2)
- 982 [Carmo et al., 2025] Carmo, D. S., Pezzulo, A. A., Villacreses, R. A., Eisenbeisz,
983 M. L., Anderson, R. L., Van Dorin, S. E., Rittner, L., Lotufo, R. A., Gerard, S. E.,
984 Reinhardt, J. M., and Comellas, A. P. (2025). Manual segmentation of opacities
985 and consolidations on ct of long covid patients from multiple annotators. *Scientific*
986 *Data*, 12(1):402. (page 9)
- 987 [Elhaminia et al., 2025] Elhaminia, B., Alsalemi, A., Nasir, E., Jahanifar, M., Awan,
988 R., Young, L. S., Rajpoot, N. M., Minhas, F., and Raza, S. E. A. (2025). From
989 traditional to deep learning approaches in whole slide image registration: A
990 methodological review. (page 2)
- 991 [Elnakib et al., 2020] Elnakib, A., Elmenabawy, N., and S Moustafa, H. (2020).
992 Automated deep system for joint liver and tumor segmentation using majority
993 voting. *MEJ-Mansoura Engineering Journal*, 45(4):30–36. (page 11)
- 994 [Farahani et al., 2015] Farahani, N., Parwani, A. V., and Pantanowitz, L.
995 (2015). Whole slide imaging in pathology: advantages, limitations, and
996 emerging perspectives. *Pathology and Laboratory Medicine International*, 7:23–33.
997 (pages xvii and 29)
- 998 [Giri and Bhatia, 2024] Giri, K. and Bhatia, S. (2024). Artificial intelligence in
999 nephrology- its applications from bench to bedside. *International Journal of*
1000 *Advances in Nephrology Research*, 7(1):90–97. (page 6)

- 1001 [Grefve et al., 2024] Grefve, J., Söderkvist, K., Gunnlaugsson, A., Sandgren, K.,
1002 Jonsson, J., Keeratijarut Lindberg, A., Nilsson, E., Axelsson, J., Bergh,
1003 A., Zackrisson, B., Moreau, M., Thellenberg Karlsson, C., Olsson, L.,
1004 Widmark, A., Riklund, K., Blomqvist, L., Berg Loegager, V., Strandberg,
1005 S. N., and Nyholm, T. (2024). Histopathology-validated gross tumor
1006 volume delineations of intraprostatic lesions using psma-positron emission
1007 tomography/multiparametric magnetic resonance imaging. *Physics and Imaging in
1008 Radiation Oncology*, 31:100633. (page 14)
- 1009 [Habis, 2024] Habis, A. A. (2024). *Developing interactive artificial intelligence tools to
1010 assist pathologists with histology annotation*. Theses, Institut Polytechnique de Paris.
1011 (page 8)
- 1012 [Hu et al., 2025] Hu, D., Jiang, Z., Shi, J., Xie, F., Wu, K., Tang, K., Cao, M., Huai, J.,
1013 and Zheng, Y. (2025). Pathology report generation from whole slide images with
1014 knowledge retrieval and multi-level regional feature selection. *Computer Methods
1015 and Programs in Biomedicine*, 263:108677. (pages 2, 28, and 29)
- 1016 [Julián and Álvarez Meza Andrés Marino, 2023] Julián, G. G. and Álvarez Meza
1017 Andrés Marino (2023). A supervised learning framework in the context of
1018 multiple annotators. (page 17)
- 1019 [Karthikeyan et al., 2023] Karthikeyan, R., McDonald, A., and Mehta, R. (2023).
1020 What's in a label? annotation differences in forecasting mental fatigue using ecg
1021 data and seq2seq architectures. (page 9)
- 1022 [Kim et al., 2024] Kim, Y., Lee, E., Lee, Y., and Oh, U. (2024). Understanding
1023 novice's annotation process for 3d semantic segmentation task with human-
1024 in-the-loop. In *Proceedings of the 29th International Conference on Intelligent User
1025 Interfaces*, IUI '24, page 444–454, New York, NY, USA. Association for Computing
1026 Machinery. (page 9)
- 1027 [Lam and Suen, 1997] Lam, L. and Suen, S. (1997). Application of majority
1028 voting to pattern recognition: an analysis of its behavior and performance.
1029 *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*,
1030 27(5):553–568. (page 11)

- 1031 [Lin et al., 2024] Lin, Y., Lian, A., Liao, M., and Yuan, S. (2024). Bcdnet: A fast
1032 residual neural network for invasive ductal carcinoma detection. (page 6)
- 1033 [López-Pérez et al., 2023] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A. D.,
1034 Molina, R., and Katsaggelos, A. K. (2023). Crowdsourcing segmentation
1035 of histopathological images using annotations provided by medical students.
1036 In Juarez, J. M., Marcos, M., Stiglic, G., and Tucker, A., editors, *Artificial
1037 Intelligence in Medicine*, pages 245–249, Cham. Springer Nature Switzerland.
1038 (pages 5, 6, 8, and 12)
- 1039 [Lu et al., 2023] Lu, X., Ratcliffe, D., Kao, T.-T., Tikhonov, A., Litchfield, L., Rodger,
1040 C., and Wang, K. (2023). Rethinking quality assurance for crowdsourced multi-
1041 roi image segmentation. *Proceedings of the AAAI Conference on Human Computation
and Crowdsourcing*, 11(1):103–114. (pages 5 and 8)
- 1043 [López-Pérez et al., 2024] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A.,
1044 Felicelli, C., Goldstein, J., Vadasz, B., Molina, R., and Katsaggelos, A. K. (2024).
1045 Learning from crowds for automated histopathological image segmentation.
1046 *Computerized Medical Imaging and Graphics*, 112:102327. (pages xvii, 5, 15, and 17)
- 1047 [Mazzarini et al., 2021] Mazzarini, M., Falchi, M., Bani, D., and Migliaccio, A. R.
1048 (2021). Evolution and new frontiers of histology in bio-medical research.
1049 *Microscopy Research and Technique*, 84(2):217–237. (pages xvii and 28)
- 1050 [Pan et al., 2021] Pan, X., Lu, Y., Lan, R., Liu, Z., Qin, Z., Wang, H., and Liu, Z. (2021).
1051 Mitosis detection techniques in h&e stained breast cancer pathological images:
1052 A comprehensive review. *Computers & Electrical Engineering*, 91:107038. (page 30)
- 1053 [Panayides et al., 2020] Panayides, A. S., Amini, A., Filipovic, N. D., Sharma, A.,
1054 Tsaftaris, S. A., Young, A., Foran, D., Do, N., Golemati, S., Kurc, T., Huang, K.,
1055 Nikita, K. S., Veasey, B. P., Zervakis, M., Saltz, J. H., and Pattichis, C. S. (2020). Ai
1056 in medical imaging informatics: Current challenges and future directions. *IEEE
1057 Journal of Biomedical and Health Informatics*, 24(7):1837–1857. (page 2)

- 1058 [Qiu et al., 2022] Qiu, Y., Hu, Y., Kong, P., Xie, H., Zhang, X., Cao, J., Wang, T.,
1059 and Lei, B. (2022). Automatic prostate gleason grading using pyramid semantic
1060 parsing network in digital histopathology. *Frontiers in Oncology*, 12. (page 14)
- 1061 [Rashmi et al., 2021] Rashmi, R., Prasad, K., and Udupa, C. B. K. (2021).
1062 Breast histopathological image analysis using image processing techniques for
1063 diagnostic purposes: A methodological review. *Journal of Medical Systems*, 46(1):7.
1064 (pages 1 and 5)
- 1065 [Ronneberger et al., 2015] Ronneberger, O., Fischer, P., and Brox, T. (2015). U-
1066 net: Convolutional networks for biomedical image segmentation. In Navab, N.,
1067 Hornegger, J., Wells, W. M., and Frangi, A. F., editors, *Medical Image Computing*
1068 and *Computer-Assisted Intervention - MICCAI 2015*, pages 234–241, Cham. Springer
1069 International Publishing. (page 14)
- 1070 [Ryou et al., 2025] Ryou, H., Thomas, E., Wojciechowska, M., Harding, L., Tam,
1071 K. H., Wang, R., Hu, X., Rittscher, J., Cooper, R., and Royston, D. (2025). Reticulin-
1072 free quantitation of bone marrow fibrosis in mpns: Utility and applications.
1073 *eJHaem*, 6(2):e70005. (page 2)
- 1074 [Sarvamangala and Kulkarni, 2022] Sarvamangala, D. R. and Kulkarni, R. V. (2022).
1075 Convolutional neural networks in medical image understanding: a survey.
1076 *Evolutionary Intelligence*, 15(1):1–22. (pages 3 and 6)
- 1077 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018). Ms-net:
1078 Mixed-supervision fully-convolutional networks for full-resolution
1079 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
1080 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer
1081 Assisted Intervention - MICCAI 2018*, pages 379–387, Cham. Springer International
1082 Publishing. (page 5)
- 1083 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
1084 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering
1085 Sciences*, 378(2166):20190061. (page 3)

- 1086 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
1087 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
1088 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
1089 Curran Associates, Inc. (page 12)
- 1090 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
1091 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
1092 Chained deep learning using generalized cross-entropy for multiple annotators
1093 classification. *Sensors*, 23(7). (page 20)
- 1094 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
1095 truth and performance level estimation (staple): an algorithm for the validation
1096 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
1097 (page 12)
- 1098 [Weitz et al., 2023] Weitz, P., Valkonen, M., Solorzano, L., Carr, C., Kartasalo, K.,
1099 Boissin, C., Koivukoski, S., Kuusela, A., Rasic, D., Feng, Y., Sinius Pouplier,
1100 S., Sharma, A., Ledesma Eriksson, K., Latonen, L., Laenholm, A.-V., Hartman,
1101 J., Ruusuvuori, P., and Rantalainen, M. (2023). A multi-stain breast cancer
1102 histological whole-slide-image data set from routine diagnostics. *Scientific Data*,
1103 10(1):562. (pages xvii and 31)
- 1104 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
1105 Advances in medical image segmentation: A comprehensive review of traditional,
1106 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 1107 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
1108 Semi-supervised gland segmentation via feature-enhanced contrastive learning
1109 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
1110 pages 1–11. (page 2)
- 1111 [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Cicarrelli, O.,
1112 Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground
1113 truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell,
1114 R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*,
1115 volume 33, pages 15750–15762. Curran Associates, Inc. (page 16)

- 1116 [Zhang and Sabuncu, 2018] Zhang, Z. and Sabuncu, M. R. (2018). Generalized
1117 cross entropy loss for training deep neural networks with noisy labels.
1118 (pages 40 and 41)
- 1119 [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S.,
1120 Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M.
1121 (2021). A review of deep learning in medical imaging: Imaging traits, technology
1122 trends, case studies with progress highlights, and future promises. *Proceedings of
the IEEE*, 109(5):820–838.
1123 (pages 1 and 2)
- 1124 [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L.
1125 (2024). Image quality evaluation in deep-learning-based ct noise reduction using
1126 virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical
Physics*, 51(8):5399–5413.
1127 (page 9)