



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 2025

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 2025

ACKNOWLEDGEMENTS

45 I would like to express my deepest gratitude to:

46 My beloved Camila, for her unwavering love, support, and patience throughout this
47 journey.

48 My family, especially my parents, for their unconditional support during challenging
49 times. I look forward to sharing many more joyful moments together.

50 My advisors, Prof. Andrés Marino Álvarez-Meza, Prof. Julián Gil, and Prof. Germán
51 Castellanos-Domínguez, for their invaluable guidance and mentorship.

52 My friends and colleagues at GCPDS, particularly Lucas, Santiago, Marcos, and
53 Rafael, for their support and camaraderie throughout this academic journey.

54 The Universidad Nacional de Colombia, for providing the resources and facilities
55 that enabled me to pursue this Master's degree. I hope this work contributes
56 meaningfully to the field of early cancer detection.

ABSTRACT

60 **Keywords:** Crowdsourcing, Multiple Annotators, Histopathology, Breast Cancer,
61 Semantic Image Segmentation, Gaussian Processes, Deep Learning, Medical Image
62 Analysis

RESUMEN

63

- 64 **Palabras clave:** Crowdsourcing, Múltiples Anotadores, Histopatología, Cáncer de
- 65 Mama, Segmentación Semántica de Imágenes, Procesos Gaussianos, Aprendizaje
- 66 Profundo, Análisis de Imágenes Médicas

67 Contents

68 Acknowledgements	vii
69 Abstract	ix
70 Resumen	xi
71 Contents	xv
72 List of figures	xviii
73 List of tables	xix
74 Abbreviations	xxi
75 1 Introduction	1
76 1.1 Motivation	1
77 1.2 Problem Statement	6
78 1.2.1 Variability in Expertise Levels	8
79 1.2.2 Technical Constraints and Image Quality	9
80 1.2.3 Research Question	9
81 1.3 Literature review	11
82 1.3.1 Facing annotation variability in medical images	12
83 1.3.2 Facing noisy annotations and low-quality data	19
84 1.4 Aims	21
85 1.4.1 General Aim	23
86 1.4.2 Specific Aims	23

87	1.5 Outline and Contributions	24
88	2 Conceptual preliminaries	25
89	2.1 Modern concept of digital image	25
90	2.1.1 Types of digital images	25
91	2.1.2 Mathematical representations	26
92	2.2 Digital histopathological images	29
93	2.2.1 Whole Slide Imaging (WSI)	29
94	2.2.2 Regions of Interest (ROI)	31
95	2.2.3 Staining Techniques	31
96	2.3 Deep learning fundamentals	32
97	2.3.1 Learning Paradigms	33
98	2.3.2 Architecture and Training	34
99	2.3.3 Challenges and Solutions	34
100	2.3.4 Deep Learning Frameworks	35
101	2.4 Datasets and data sources	35
102	2.4.1 Datasets with emulated noisy annotations	37
103	2.4.2 Real histopathology datasets	40
104	3 Chained Gaussian Processes	43
105	3.1 Background and Related Methods	43
106	3.1.1 Kernel Alignment-Based Annotator Relevance Analysis (KAAR)	43
107	3.1.2 Localized Kernel Alignment-Based Annotator Relevance Analysis (LKAAR)	44
108	3.1.3 Regularized Chained Deep Neural Network (RCDNN)	45
109	3.1.4 Chained Gaussian Processes Model	46
110	3.2 Tooling and Implementation	47
111	3.2.1 GPflow Overview	48
112	3.2.2 GPflux Framework	49
113	3.2.3 Implementation Considerations	50
114	3.3 Experimental Setup	51
115	3.3.1 Classification	51

117	3.3.2 Regression	52
118	3.4 Summary	54
119	3.4.1 Comparative Analysis	54
120	3.5 Conclusion	55
121	4 Truncated Generalized Cross Entropy for segmentation	57
122	4.1 Loss functions for multiple annotators	57
123	4.1.1 Generalized Cross Entropy	58
124	4.1.2 Extension to Multiple Annotators	59
125	4.1.3 Reliability Maps and Truncated GCE	60
126	5 Chained deep learning for image segmentation	63
127	5.1 Introduction	63
128	5.2 Using U-NET as a building block	63
129	5.2.1 Backbone Architecture	64
130	5.2.2 UNET Architecture	64
131	5.2.3 Reliability Map Branch	65
132	5.2.4 Integration with TGCE_{SS} Loss	66
133	5.2.5 Training Process	66
134	6 Conclusions	69
135	6.1 Summary	69
136	6.2 Future work	70
137	Bibliography	71

LIST OF FIGURES

139	1-1	Estimation of the popularity of tasks and medical image types based on recent literature review (count of referenced terms)	3
140	1-2	AI and machine learning in medical imaging brief timeline.	4
141	1-3	Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.	7
142	1-4	Summary diagram for problem Statement	10
143	1-5	Proposed framework for the approach in [López-Pérez et al., 2024]. .	17
144	2-1	Overview of digital image concepts and their mathematical representations. The figure shows the main types of digital images (grayscale, color, multispectral, and volumetric), their mathematical representations, and the fundamental processes of sampling and quantization. Example images are included to illustrate each type. .	28
145	2-2	Histology evolution timeline. (Image from [Mazzarini et al., 2021]). .	29
146	2-3	(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]).	30
147	2-4	Common learning paradigms.	33
148	2-5	Comparative Trends of the top two most popular Deep Learning Frameworks, apparently, tendency was switched to PyTorch since 2022	36
149	2-6	Example of a noisy mask generated by naively introducing random noise into a ground truth mask. Morphological consistency is lost. .	37

162	2-7 Annotations in the Oxford-IIIT Pet data. From left to right: pet 163 image, head bounding box, and trimap segmentation (blue: 164 background region; red: ambiguous region; yellow: foreground 165 region).	39
166	2-8 Noisy mask generated by enhancing the disturbances in the 167 encoder layers weights for the Oxford-IIIT Pet Dataset. 168 Morphological consistency is preserved. From left to right, SNR 169 levels of noise in the encoder layer are 10, 5, 2, 0, -5 dB.	39
170	2-9 Different staining techniques obtained from multi-stain breast 171 cancer dataset [Weitz et al., 2023]. (a) shows H&E, (b) ER, (c) HER2, 172 (d) Ki67 and (e) PGR. (f) shows an example of a Whole Slide Imaging 173 (WSI) that was excluded since it contains multiple tissue sections.	41
174	2-10 Screenshot of the DSA and HistomicsTK web interface while creating 175 the crowdsourced annotations for the dataset presented by [Amgad 176 et al., 2019].	42
177	4-1 Working mechanism of the proposed Truncated Generalized Cross 178 Entropy for Semantic Segmentation ($TGCE_{SS}$) loss function. The loss 179 combines reliability maps Λ_r with the model's predictions $f(\mathbf{X}; \theta)$ 180 to compute a weighted loss that accounts for annotator reliability 181 across different image regions.	62
182	5-1 Original U-NET architecture.	64
183	5-2 Overview of the proposed model architecture.	65
184	5-3 Training process of the proposed model overview. Estimated 185 segmentation and reliability maps are computed for each input 186 image, and the loss function is computed for the entire batch.	66

LIST OF TABLES

188	1-1	Summary of state-of-the-art approaches for handling multiple annotators in medical image segmentation	22
189			
190	3-1	Comparison of Multiple Annotator Learning Methods	55

ABBREVIATIONS

- 192 **CAD** Computer-Aided Diagnosis 2, 5, 6
193 **CCGP** Correlated Chained Gaussian Processes 18
194 **CCGPMA** Correlated Chained Gaussian Processes for Multiple Annotators 17, 19, 66
195 **CE** Cross Entropy 54, 57
196 **CGP** Chained Gaussian Processes 18
197 **CNN** Convolutional Neural Networks 3, 14, 21, 24, 59, 66
198 **CT** Computed Tomography 1, 2, 11
199 **ELBO** Evidence Lower Bound 18
200 **GCE** Generalized Cross Entropy 54, 57
201 **GCECDL** Generalized Cross-Entropy-based Chained Deep Learning 20
202 **ISS** Image Semantic Segmentation 2, 3, 6, 11, 13, 21, 23, 24, 35, 65, 66
203 **LF** Latent Function 18
204 **MAE** Mean Absolute Error 54, 57
205 **MITs** Medical Imaging Techniques 1
206 **ML** Machine Learning 11, 21
207 **MRI** Magnetic Resonance Imaging 1, 2
208 **MV** Majority Voting 11, 12
209 **OCR** Optical Character Recognition 11
210 **PET** Positron Emission Tomography 14
211 **ROI** Region of Interest 2, 6, 15, 31
212 **SLFM** Semi-Parametric Latent Factor Model 18
213 **SS** Semantic segmentation 3
214 **STAPLE** Simultaneous Truth and Performance Level Estimation 12–14, 36
215 **WSI** Whole Slide Imaging xviii, 1, 5, 6, 8, 15, 29, 41

216

217

CHAPTER

218

219

ONE

220

INTRODUCTION

221 1.1 Motivation

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

233 Along with the advances in technologies for medical images acquisition,
234 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

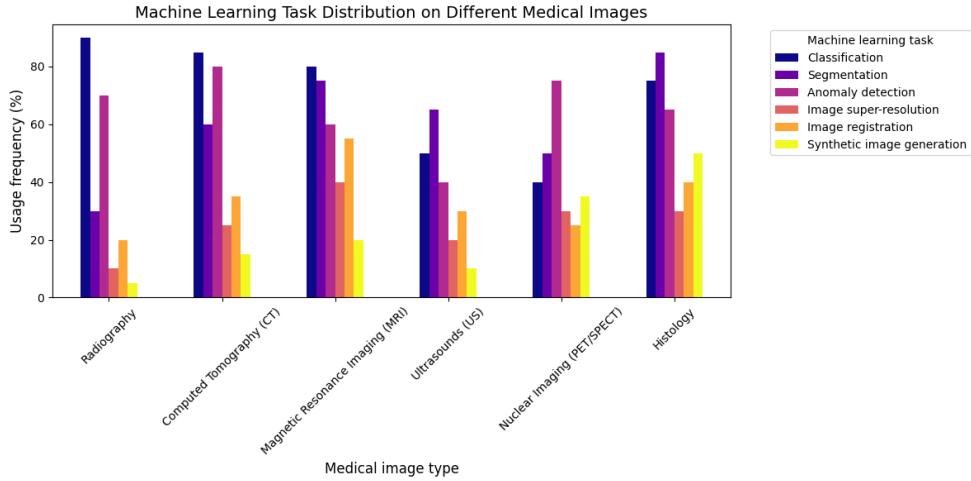


Figure 1-1 Estimation of the popularity of tasks and medical image types based on recent literature review (count of referenced terms).

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

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

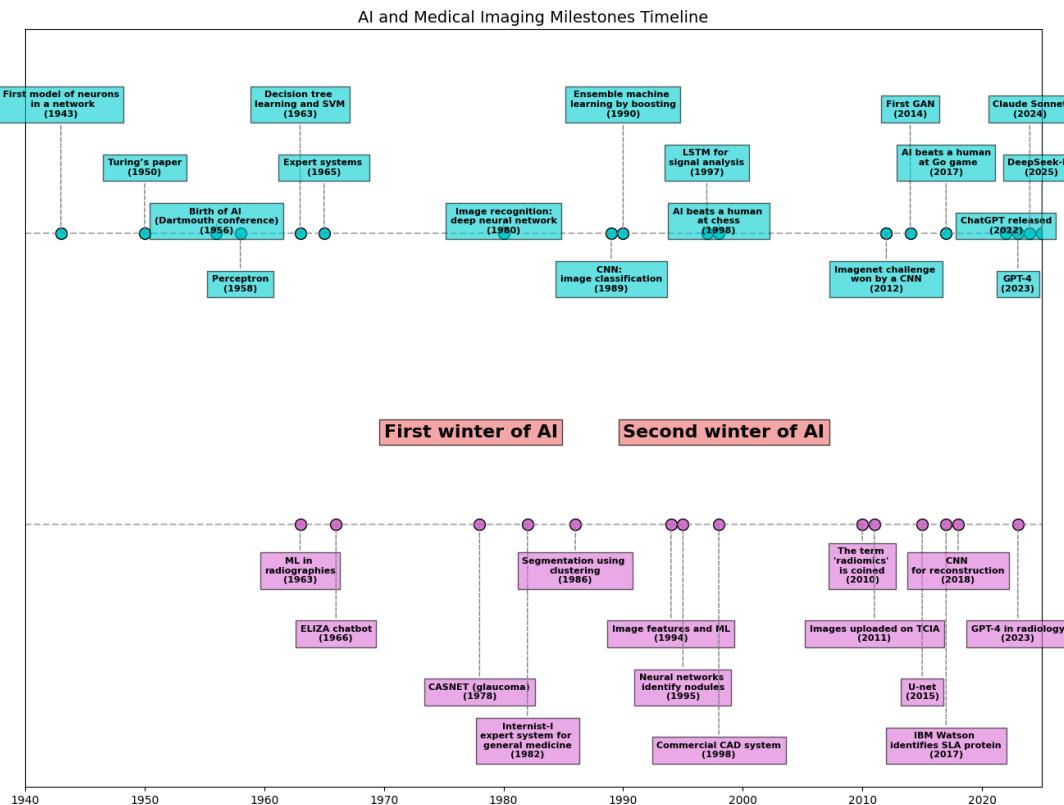


Figure 1-2 AI and machine learning in medical imaging brief timeline.

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

284 Many medical images datasets however, contain a high variability in class sizes
285 and variations in colors, which is specially noticeable in histopathological images
286 because of the usage of different staining and other factors which can affect the
287 color of the images (see section 2.2.3). This variability can lead to a significant loss
288 of efficiency of machine learning models when using a mixed supervision strategy,
289 as the model can be biased towards the most common classes or colors in the
290 dataset [Shah et al., 2018].

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

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

308 A major challenge in histopathological image segmentation arises from the
309 variability in annotations provided by different pathologists. Unlike natural
310 images, where object boundaries are often well-defined, histological structures
311 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

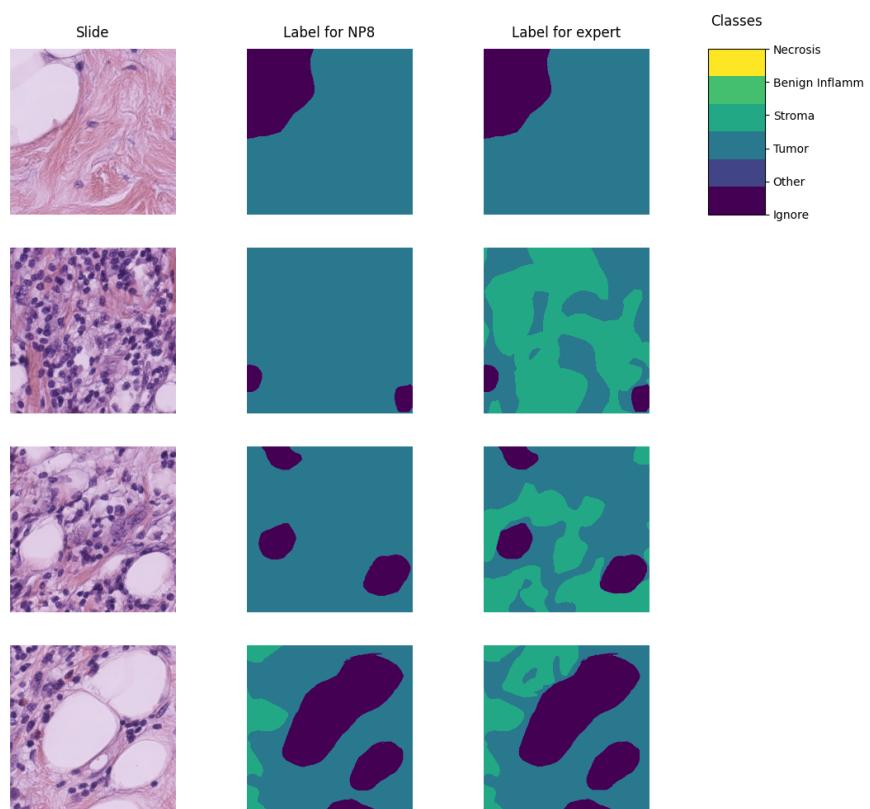


Figure 1-3 Example of a histopathological image segmented by multiple annotators, illustrating variations in label assignment.

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

347 1.2.1 Variability in Expertise Levels

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

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

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

369 **1.2.2 Technical Constraints and Image Quality**

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

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

380 **1.2.3 Research Question**

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

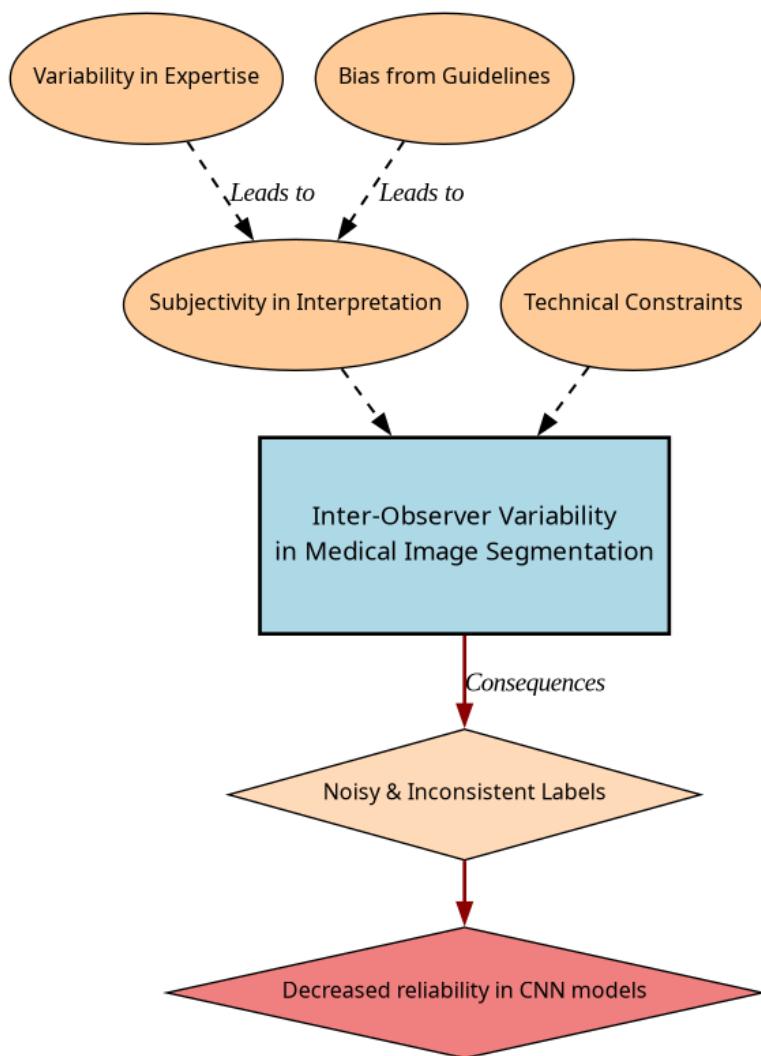


Figure 1-4 Summary diagram for problem Statement

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?

383

384 1.3 Literature review

385 Certainly, in general Machine Learning (ML) classification tasks ³ where multiple
386 annotators are involved, Majority Voting (MV) is by far the simplest possible
387 approach to implement. This concept was born multiple times and divergently in
388 multiple fields, but it was described as relevant for ML and pattern recognition
389 labeling for classification in [Lam and Suen, 1997], in which the approach is
390 exposed as simple, yet powerful. The authors describe the MV as a method that
391 can be used to improve the accuracy of classification tasks by combining the labels
392 of multiple annotators. The method is based on the assumption that the majority
393 vote of the annotators is more likely to be correct than the vote of a single
394 annotator. The authors also describe the method as a straightforward way to
395 improve the accuracy of classification tasks without the need for complex
396 algorithms or additional data. The authors also prove this method to deliver very
397 similar results to more complicated approaches (Bayesian, logistic regression,
398 fuzzy integral, and neural network) in the particular task of Optical Character
399 Recognition (OCR). Despite its simplicity, modern solutions for delivering accurate
400 medical image segmentation models still rely on Majority Voting at some stage,
401 like [Elnakib et al., 2020], which uses a majority voting strategy for delivering a
402 final output based on the labels of multiple models (VGG16-Segnet, Resnet-18 and
403 Alexnet) in CT images for Liver Tumor Segmentation, or [López-Pérez et al., 2023],

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

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

413 1.3.1 Facing annotation variability in medical images

414 Learning from crowds approaches in general face the challenge of not having a
415 ground truth label and hence, an intrinsic difficulty in measuring the real reliability
416 of the labelers annotations. Some approaches assume beforehand a certain level of
417 expertise for each labeler based on experience as an input, like in [TIAN and Zhu,
418 2015], which introduce the concept of max margin majority voting, using the
419 reliability vector as weights for the weights for the binary and multiclass classifier.
420 The crowdsourcing margin is the minimal difference between the aggregated score
421 of the potential true label and the scores for other alternative labels. Accordingly,
422 the annotators' reliability is estimated as generating the largest margin between
423 the potential true labels and other alternatives. The problem introduced in this
424 approach is assuming an stationary reliability per expert across the whole input
425 space, which is imprecise since annotators performance may change between
426 different tasks or even between different regions of the same image.

427 STAPLE Mechanism

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

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

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

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

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

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

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

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

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

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

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

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

U-shaped CNNs

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

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

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

479 **U-Net based approaches**

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

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

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

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

494 The authors then assume the annotator network to deliver a reliability map
 495 $\{\hat{\mathbf{A}}_\phi^{(r)}(\mathbf{x})\}_{r \in R_n}$ with different dimensions:

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

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

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

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

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

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

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

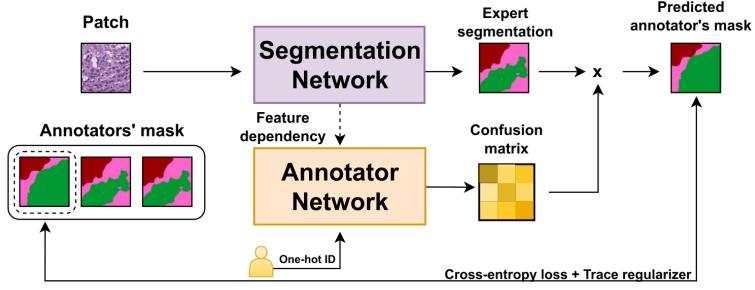


Figure 1-5 Proposed framework for the approach in [López-Pérez et al., 2024].

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

526 Bayesian models

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

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

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

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

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

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

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

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

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

553 CCGPMA on itself proposes a good approach for handling label noise and
554 uncertainty in the labelers for regression and classification tasks, while also
555 preserving the interdependencies among the annotators, however, it does not face
556 the image segmentation problem, which is the main focus of this works, however,
557 it does not face the image segmentation problem, which is the main focus of this
558 work. Besides, handling so many latent functions during the optimization process
559 is computationally expensive, making it on itself infeasible for large and high
560 resolution datasets.

561 1.3.2 Facing noisy annotations and low-quality data

562 The problem of low-quality data and noisy annotations has been tackled with
563 various strategies. One such approach is the use of deep learning models that
564 incorporate loss functions designed to mitigate the effects of unreliable labels.
565 Traditional methods such as Majority Voting (MV) or Expectation-Maximization
566 (EM) have been widely used for aggregating multiple annotators' inputs. However,
567 they assume a homogeneous reliability of annotators, which may not hold in
568 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.

569 **Loss functions in deep learning models**

570 Loss functions are fundamental components in deep learning models that quantify
571 how well a model’s predictions match the ground truth. They serve as the
572 objective function that guides the learning process by measuring the discrepancy
573 between predicted and actual values. In classification tasks, the most common
574 loss functions are Cross-Entropy (CE) and Mean Absolute Error (MAE) [Zhang and
575 Sabuncu, 2018]. CE is particularly effective for classification as it heavily penalizes
576 confident but wrong predictions, though it can be sensitive to noisy labels. MAE,
577 on the other hand, is more robust to outliers and assigns equal weights to all
578 mistakes, but typically requires more training iterations. For image segmentation
579 tasks, specialized loss functions have been developed to handle the unique
580 challenges of pixel-wise classification. The Dice loss, which measures the overlap
581 between predicted and ground truth regions, is widely used in medical image
582 segmentation [Zhao et al., 2020]. More recently, the Generalized Cross Entropy
583 (GCE) loss has emerged as a robust alternative that combines the benefits of both
584 CE and MAE, allowing for better handling of noisy labels through a tunable
585 parameter that controls sensitivity to outliers. However, to the best of our
586 knowledge, none of the current loss functions approaches have focused on the
587 multiple labelers scenario, neither on the per labeler spatial reliability estimation.

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

589 A more recent approach was proposed by [Triana-Martinez et al., 2023],
590 introducing a Generalized Cross-Entropy-based Chained Deep Learning (GCECDL)
591 framework. This method addresses the limitations of traditional label aggregation
592 techniques by modeling each annotator’s reliability as a function of the input data.
593 The approach effectively mitigates the impact of noisy labels by using a
594 noise-robust loss function, balancing Mean Absolute Error (MAE) and Categorical
595 Cross-Entropy (CE). Unlike prior approaches, GCECDL accounts for the
596 dependencies among annotators while encoding their non-stationary behavior

597 across different data samples. Their experiments on multiple datasets
 598 demonstrated superior predictive performance compared to state-of-the-art
 599 methods, particularly in cases where annotations were highly inconsistent.

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

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

604 where $\tilde{\lambda}_x$ represents the annotator reliability, \tilde{C} is a constant, q is a parameter that
 605 controls the balance between MAE and CE behavior, \mathbf{y} is the annotation vector, and
 606 $f(\mathbf{x})$ is the model prediction. This approach is more deeply discussed in chapter 4.

607 1.4 Aims

608 With the mentioned considerations in section 1.3 in mind, this work proposes a
 609 novel approach for ISS tasks in medical images, which aims to train a model whose
 610 learning approach is adaptive to the labeler performance. This is done by
 611 introducing a loss function capable of inferring the best possible segmentation
 612 without needing separate inputs about the labeler performance. This loss function
 613 is designed to implicitly weigh the labelers based on their performance, with the
 614 presence of an intermediate reliability map allowing the model to learn from the
 615 most reliable labelers and ignore the noisy labels. This approach differs from
 616 existing CNN-based segmentation models, as it does not require explicit
 617 supervision of the labeler performance, making it more generalizable and
 618 adaptable to different datasets and labelers.

Table 1-1 Summary of state-of-the-art approaches for handling multiple annotators in medical image segmentation

Approach	Key Characteristics			Advantages	Limitations
Majority Voting (MV)	Simple aggregation of multiple annotators' labels			<ul style="list-style-type: none"> - Simple to implement - No complex algorithms required - Proven effective in many cases 	<ul style="list-style-type: none"> - Assumes equal expertise - Sensitive to correlated errors - May introduce bias in expert tasks
STAPLE	Probabilistic framework estimating true segmentation and rater reliability			<ul style="list-style-type: none"> - Handles binary and multiclass segmentation - Estimates rater sensitivity/specificity - Robust to correlated errors 	<ul style="list-style-type: none"> - Computationally expensive - Sensitive to low-quality annotations - Tends to over-smooth results - Limited interpretability
U-Net based approaches	Deep learning architecture with encoder-decoder structure			<ul style="list-style-type: none"> - Captures both global and local information - Effective for complex structures - Works well with limited data 	<ul style="list-style-type: none"> - Requires large amounts of training data - Complex optimization process - May need specialized loss functions
Bayesian Models (CCGPMA)	Gaussian Processes modeling annotator reliability			<ul style="list-style-type: none"> - Preserves annotator interdependencies - Handles label noise effectively - Works for classification and regression 	<ul style="list-style-type: none"> - Computationally expensive - Not directly applicable to segmentation - Complex optimization process
Generalized Entropy (GCE)	Cross-Entropy (GCE)	Loss function balancing MAE and CE		<ul style="list-style-type: none"> - Robust to noisy labels - Accounts for non-stationary behavior - Handles annotator dependencies 	<ul style="list-style-type: none"> - Limited to classification tasks - Requires parameter tuning - May need adaptation for segmentation

619 1.4.1 General Aim

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

626 1.4.2 Specific Aims

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

636 1.5 Outline and Contributions

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

- 639 • A novel loss function for ISS tasks in medical images capable of inferring the
640 best possible segmentation and codifying the reliability of the labelers
641 without needing separate inputs about the labeler performance or multiple
642 networks to be trained.
- 643 • A Chained Deep Learning model for ISS tasks in medical images, which can
644 learn from inconsistent labels and improve the segmentation performance
645 outperforming state of the art approaches on a Triple Negative Breast Cancer
646 histopathology dataset.
- 647 • A mechanism for noisy annotations emulations by using networks encoding
648 layers weight disturbance for building crowdsourced version of the well
649 known Oxford-IIIT Pet Dataset.
- 650 • A python package for using the proposed loss function and architecture in
651 CNN models for ISS tasks in medical images.¹¹
- 652 • Triple negative breast cancer and Oxford-IIIT Pet Datasets mapping as lazy
653 loaders for the proposed loss function ready to go in the package.¹²
- 654 • 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

655

CHAPTER

656

657

TWO

658

659

CONCEPTUAL PRELIMINARIES

660

2.1 Modern concept of digital image

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

666

2.1.1 Types of digital images

667

Grayscale images

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

673 Color images

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

680 Multispectral images

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

687 3D images and volumetric data

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

693 2.1.2 Mathematical representations

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

- 695 • **Sampling:** The process of converting a continuous image into a discrete
696 representation. According to the Nyquist-Shannon sampling theorem, the
697 sampling frequency must be at least twice the highest frequency present in
698 the image to avoid aliasing.
- 699 • **Quantization:** The process of converting continuous intensity values into
700 discrete levels. The number of quantization levels determines the image's
701 bit depth and affects its quality and storage requirements.
- 702 • **Resolution:** The number of pixels per unit length in an image, typically
703 measured in pixels per inch (PPI) or dots per inch (DPI).
- 704 • **Dynamic range:** The ratio between the maximum and minimum measurable
705 light intensities in an image, often expressed in decibels (dB).

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

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

707 where $I(x, y)$ is the digital image, $f(i, j)$ represents the intensity values, and $\delta(x -$
708 $i, y - j)$ is the Kronecker delta function.

709 For color images, the representation extends to:

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

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

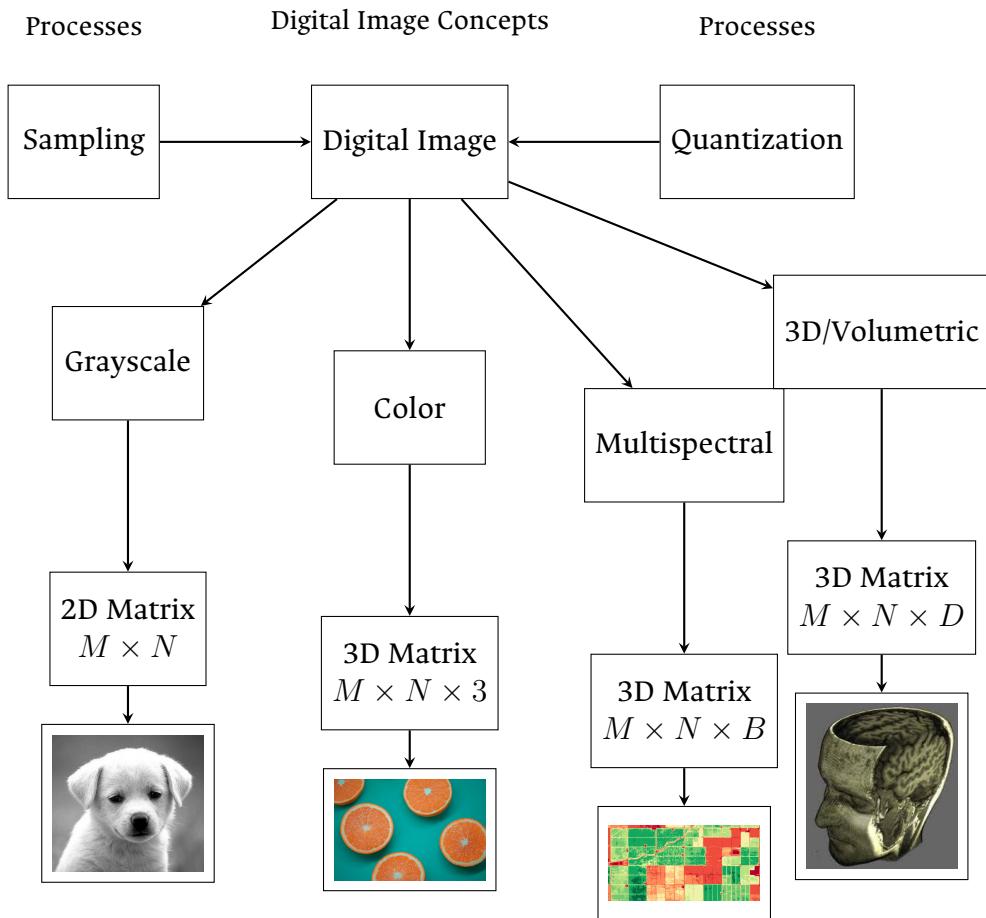


Figure 2-1 Overview of digital image concepts and their mathematical representations. The figure shows the main types of digital images (grayscale, color, multispectral, and volumetric), their mathematical representations, and the fundamental processes of sampling and quantization. Example images are included to illustrate each type.

711 2.2 Digital histopathological images

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

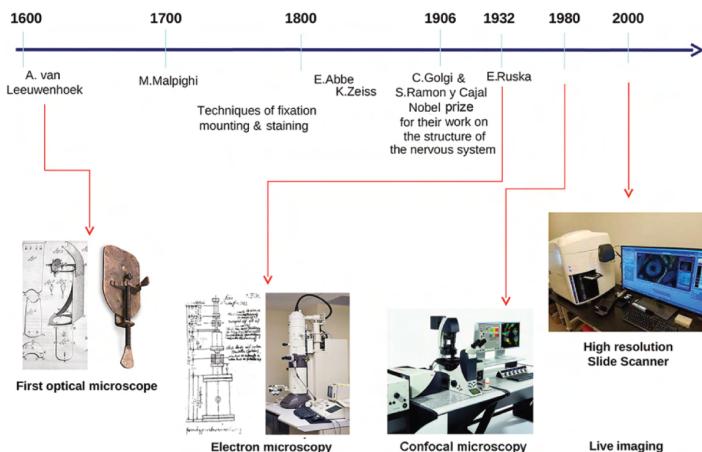


Figure 2-2 Histology evolution timeline. (Image from [Mazzarini et al., 2021]).

718 2.2.1 Whole Slide Imaging (WSI)

719 Whole Slide Imaging (WSI) is the process of digitizing entire glass slides at high
720 resolution, creating a digital representation that can be viewed, analyzed, and
721 shared electronically. Modern WSI scanners use sophisticated optical systems that
722 capture multiple fields of view at high magnification, which are then stitched
723 together to create a seamless digital image [Hu et al., 2025]. These systems
724 incorporate high-resolution objectives with magnifications ranging from 20x to
725 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 2-3 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].

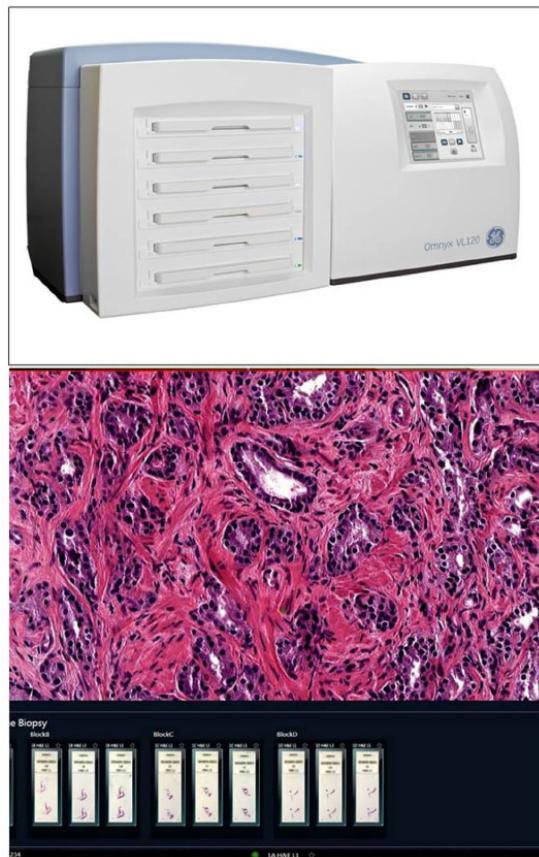


Figure 2-3 (Above) Whole slide imaging system by Omnyx for slide digitization. (Below) Comprehensive digital pathology interface from Omnyx designed to streamline pathologists' diagnostic workflow. (From [Farahani et al., 2015]).

732 2.2.2 Regions of Interest (ROI)

733 In digital histopathology, Region of Interests (ROIs) are specific areas within a
734 whole slide image that contain diagnostically relevant information. These regions
735 can be manually annotated by pathologists, automatically detected using
736 computer vision algorithms, or defined based on specific tissue characteristics or
737 abnormalities. The importance of ROIs lies in their ability to focus computational
738 analysis on relevant areas, reduce computational complexity in automated
739 systems, facilitate targeted diagnosis and research, and enable efficient storage and
740 transmission of critical information.

741 2.2.3 Staining Techniques

742 Histopathological analysis relies heavily on various staining techniques to enhance
743 the visibility of different tissue components and cellular structures. The choice of
744 staining method depends on the specific diagnostic requirements and the type of
745 tissue being examined.

746 Hematoxylin and Eosin (H&E)

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

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

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

760 **Special Stains**

761 In addition to H&E, various special stains are used for specific diagnostic purposes.
762 Immunohistochemistry (IHC) uses antibodies to detect specific proteins, playing a
763 crucial role in subtyping breast cancers. Key IHC stains include Estrogen Receptor
764 (ER) staining for detecting estrogen receptors, Progesterone Receptor (PGR)
765 staining for assessing progesterone receptor status, Human Epidermal Growth
766 Factor Receptor 2 (HER2) staining for evaluating HER2 protein expression, and
767 Ki67 staining for measuring cellular proliferation rates. These markers are
768 particularly crucial in breast cancer diagnosis and treatment planning, as they help
769 determine the molecular subtype of the cancer and guide personalized therapeutic
770 approaches. Other specialized stains include Periodic Acid-Schiff (PAS) for
771 highlighting carbohydrates and basement membranes, Masson's Trichrome for
772 distinguishing between collagen and muscle fibers, and silver stains for detecting
773 microorganisms and nerve fibers. These specialized staining techniques
774 complement H&E by providing additional diagnostic information that is crucial for
775 accurate diagnosis and treatment planning [Weitz et al., 2023]. Examples of these
776 staining techniques are shown in Figure 2-9.

777 **2.3 Deep learning fundamentals**

778 Deep learning has emerged as a powerful subset of machine learning,
779 revolutionizing the field of artificial intelligence. Its roots can be traced back to the
780 early development of artificial neural networks in the 1940s and 1950s, with
781 significant milestones including the perceptron in 1958 and the backpropagation

782 algorithm in the 1980s. However, it wasn't until the early 21st century, with the
783 advent of more powerful computational resources and the availability of large
784 datasets, that deep learning truly began to flourish.

785 2.3.1 Learning Paradigms

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

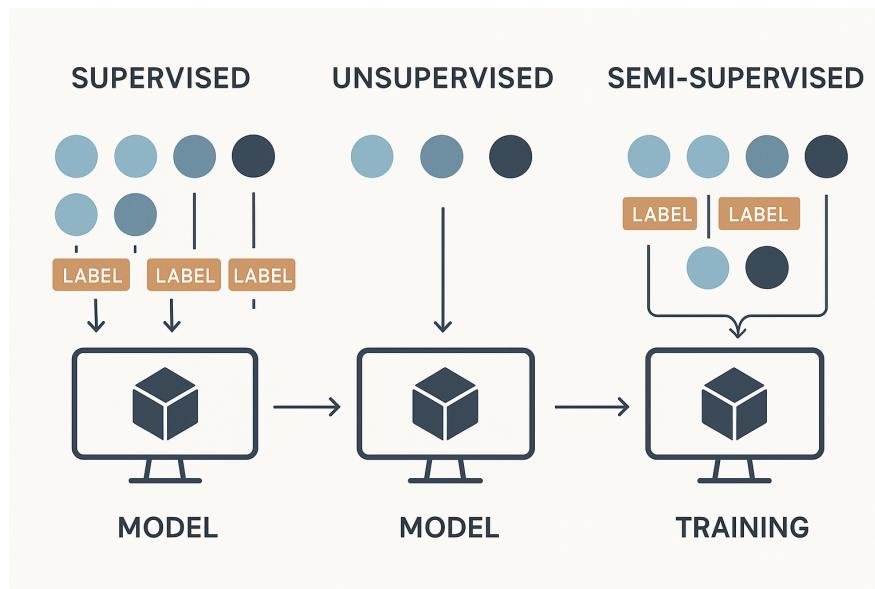


Figure 2-4 Common learning paradigms.

795 2.3.2 Architecture and Training

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

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

809 2.3.3 Challenges and Solutions

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

821 2.3.4 Deep Learning Frameworks

822 The development of powerful open-source frameworks has significantly
823 accelerated deep learning research and applications. TensorFlow, developed by
824 Google, provides a comprehensive ecosystem for building and deploying machine
825 learning models [Abadi et al., 2016]. PyTorch, created by Facebook’s AI Research
826 lab, offers dynamic computation graphs and has become particularly popular in
827 research settings. Caffe, known for its speed and modularity, is widely used in
828 computer vision applications.

829 These frameworks have democratized deep learning by providing efficient
830 implementations of common operations, automatic differentiation for gradient
831 computation, and GPU acceleration for faster training. They also offer pre-trained
832 models and transfer learning capabilities, along with active communities for
833 support and knowledge sharing. The combination of these frameworks with
834 modern hardware has enabled researchers and practitioners to develop
835 increasingly sophisticated models, pushing the boundaries of what’s possible in
836 artificial intelligence. As shown in Figure 2-5, which presents data from Google
837 Trends over the last five years (as of April 2025), TensorFlow and PyTorch have
838 emerged as the two most prominent frameworks in the deep learning landscape.

839 2.4 Datasets and data sources

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

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

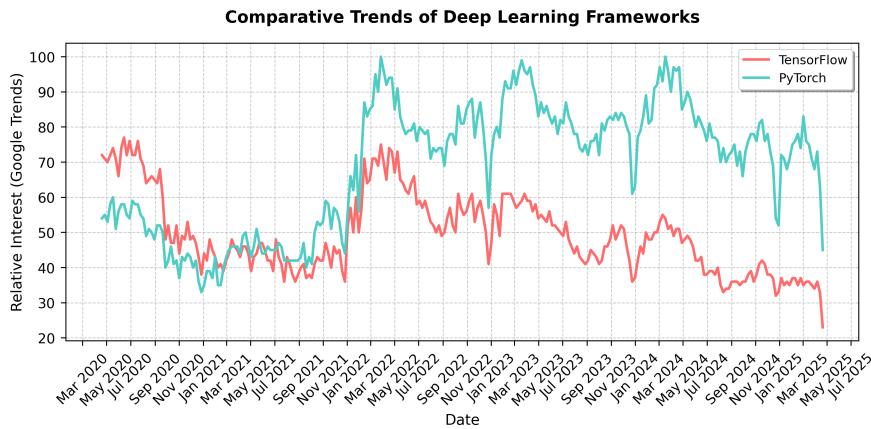


Figure 2-5 Comparative Trends of the top two most popular Deep Learning Frameworks, apparently, tendency was switched to PyTorch since 2022

847 • Distribution of segmentation labels is not uniform across the image, since
848 some tissues and structures are more common than others.

849 • Visualization of performance of the models in debug time (like per epoch
850 analysis) is not simple for non experts in the subject, which makes it hard to
851 evaluate whether the model is overfitting or not at a glance.

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

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

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

861 2.4.1 Datasets with emulated noisy annotations

- 862 A challenge arises for the creation of emulated noisy annotations datasets, since it
 863 is expected for images annotations to have some degree of expertise variability,
 864 similar to what is expected in real crowdsourced datasets. Simply introducing
 865 random noise into the annotated masks does not work, since the original
 866 morphological structures from the expected ground truth are far from being
 867 preserved. Instead, a “noisy” labeler is expected to produce an annotation which
 868 has at least some degree of morphological consistency on itself, even if it shows
 869 discrepancies when compared with some metric (like DICE score) against the
 870 ground truth annotation.
- 871 This has been proven experimentally when introducing random noise into any
 872 popular segmentation dataset, in which the resulting mask is just a non coherent
 873 map of noise across the image, as shown in Figure 2-6.

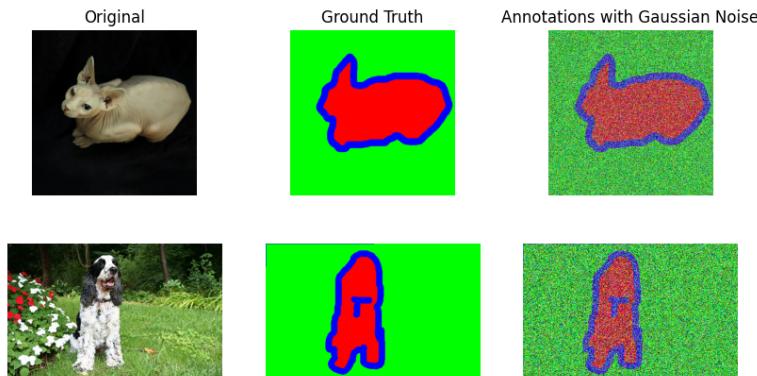


Figure 2-6 Example of a noisy mask generated by naively introducing random noise into a ground truth mask. Morphological consistency is lost.

- 874 For this reason, a more sophisticated approach was needed to create datasets with
 875 emulated noisy annotations. The approach used here was to use a pre-trained
 876 U-Net model to produce a good enough segmentation of the image, which was
 877 then used to produce a “noisy” annotation by introducing random noise into the
 878 last encoder layers of the model, thus preserving the morphological consistency of

the original annotation. This strategy works since slight modifications introduced into the encoder layers weights, somehow resemble conceptual disturbances with respect to the original ground truth label, which kind of emulates the cultural behavior of having a different interpretation of the image, either for having a different level of expertise or for having a different point of view or school of thought. The first encoding layers were not modified, since it is expected human labelers would agree on the most fundamental structures (analog to extracted features from initial convolutional layers) in a similar way, thus preserving the morphological consistency of the original annotation.

In this way, the level of “disturbance” with respect to the ground truth for an emulated annotator can be controlled by the level of noise introduced into the weights of the last encoder layers, thus:

$$\mathbf{W}_{noisy} = \mathbf{W}_{original} + \mathcal{N}(0, \sigma^2) \quad (2-3)$$

where $\mathbf{W}_{original}$ represents the original weights of the last encoder layers, $\mathcal{N}(0, \sigma^2)$ is a Gaussian distribution with mean 0 and variance σ^2 , and \mathbf{W}_{noisy} are the resulting noisy weights. The variance σ^2 controls the level of noise introduced, and thus the degree of disturbance in the resulting segmentation masks.

895 Oxford-IIIT Pet Dataset

Using the techniques described above, the Oxford-IIIT Pet Dataset [Parkhi et al., 2012] was used to create a dataset with emulated noisy annotations. The almost perfectly uniform distribution of the dataset classes makes it an ideal playground dataset for testing segmentation models with a high degree of confidence in the ground truth annotations, at the same time that cats and dogs are a common sight in the daily life of most people, making it easier to find a labeler that is able to

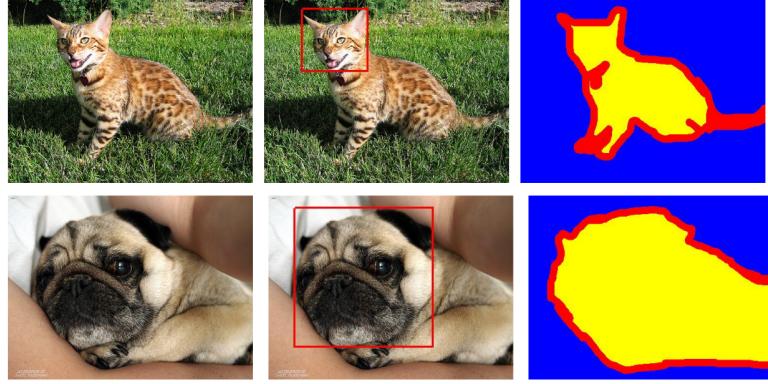


Figure 2-7 Annotations in the Oxford-IIIT Pet data. From left to right: pet image, head bounding box, and trimap segmentation (blue: background region; red: ambiguous region; yellow: foreground region).

902 annotate the images with a high degree of accuracy, which facilitates model initial
 903 debugging. Figure 2-7 shows an example of the annotations in the original dataset.
 904 With the application of the encoder layer weight perturbation technique, the
 905 resulting noisy masks are shown in Figure 2-8. It can be seen that the
 906 morphological consistency is preserved, even though the resulting masks are far
 907 from the ground truth annotations, which goes perfectly well for testing the
 908 robustness of the models against noisy annotations in crowdsourcing-like
 scenarios.

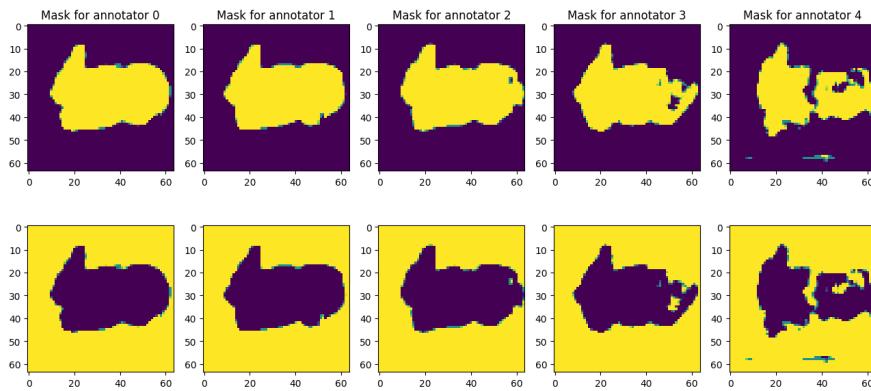


Figure 2-8 Noisy mask generated by enhancing the disturbances in the encoder layers weights for the Oxford-IIIT Pet Dataset. Morphological consistency is preserved. From left to right, SNR levels of noise in the encoder layer are 10, 5, 2, 0, -5 dB.

910 2.4.2 Real histopathology datasets

911 Multi-Stain Breast Cancer Histological Dataset

912 The Multi-Stain Breast Cancer Histological Dataset [Weitz et al., 2023] represents
913 one of the largest publicly available collections of whole slide images (WSIs) from
914 surgical resection specimens of primary breast cancer patients. This dataset is
915 particularly valuable for our work because it contains matched pairs of H&E and
916 IHC-stained tissue sections from the same tumor, with a total of 4,212 WSIs from
917 1,153 patients. The IHC stains include important biomarkers such as ER, PGR,
918 HER2, and KI67, which are routinely used in breast cancer diagnosis and
919 treatment planning (more on staining techniques in Section ??).

920 The dataset’s relevance to our work stems from several key aspects. The matched
921 H&E and IHC stains allow for studying the consistency of segmentation across
922 different staining modalities, which is crucial for understanding how different
923 visualization methods affect annotation quality. With 1,153 patients, the dataset
924 provides a robust foundation for training and evaluating segmentation models in a
925 real-world clinical setting. The inclusion of routine diagnostic cases makes the
926 dataset representative of actual clinical practice, where variations in staining
927 quality and tissue preparation are common. Furthermore, the multiple biomarker
928 stains (ER, PGR, HER2, KI67) enable the study of how different tissue
929 characteristics affect segmentation performance and annotator agreement.

930 This dataset serves as an ideal testbed for our crowdsourced segmentation
931 approach. It allows us to evaluate how different staining modalities affect
932 annotator performance and agreement, while also providing insights into the
933 relationship between tissue characteristics and segmentation difficulty. The
934 dataset’s comprehensive nature enables validation of our models’ performance
935 across different biomarker expressions and assessment of the generalizability of
936 our approach to real-world clinical data.

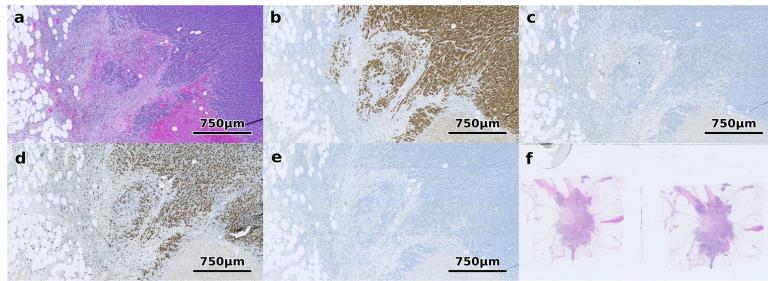


Figure 2-9 Different staining techniques obtained from multi-stain breast cancer dataset [Weitz et al., 2023]. (a) shows H&E, (b) ER, (c) HER2, (d) Ki67 and (e) PGR. (f) shows an example of a WSI that was excluded since it contains multiple tissue sections.

937 Structured Crowdsourcing Dataset for Histology Images

938 The dataset presented by [Amgad et al., 2019] is particularly relevant to our work
939 as it represents one of the first systematic studies of crowdsourced annotations in
940 histopathology. The authors recruited 25 participants with varying levels of
941 expertise (from senior pathologists to medical students) to delineate tissue regions
942 in 151 breast cancer slides using the Digital Slide Archive platform, resulting in
943 over 20,000 annotated tissue regions.

944 Key aspects of this dataset make it valuable for our work. The systematic
945 evaluation of inter-participant discordance revealed varying levels of agreement
946 across different tissue classes, with low discordance for tumor and stroma, and
947 higher discordance for more subjectively defined or rare tissue classes. The
948 inclusion of feedback from senior participants helped in curating high-quality
949 annotations, demonstrating that fully convolutional networks trained on these
950 crowdsourced annotations can achieve high accuracy (mean AUC=0.945). The
951 dataset also provides evidence that the scale of annotation data significantly
952 improves image classification accuracy.

953 This dataset is particularly valuable for our work because it provides crucial
954 insights into how annotator expertise affects segmentation quality. It
955 demonstrates the feasibility of using crowdsourced annotations for training
956 accurate segmentation models, showing that even with varying levels of expertise,

957 aggregated annotations can produce reliable ground truth. The dataset includes a
 958 systematic analysis of inter-annotator agreement, which is crucial for
 959 understanding the challenges in crowdsourced histopathology segmentation and
 960 informing the development of more robust segmentation approaches.

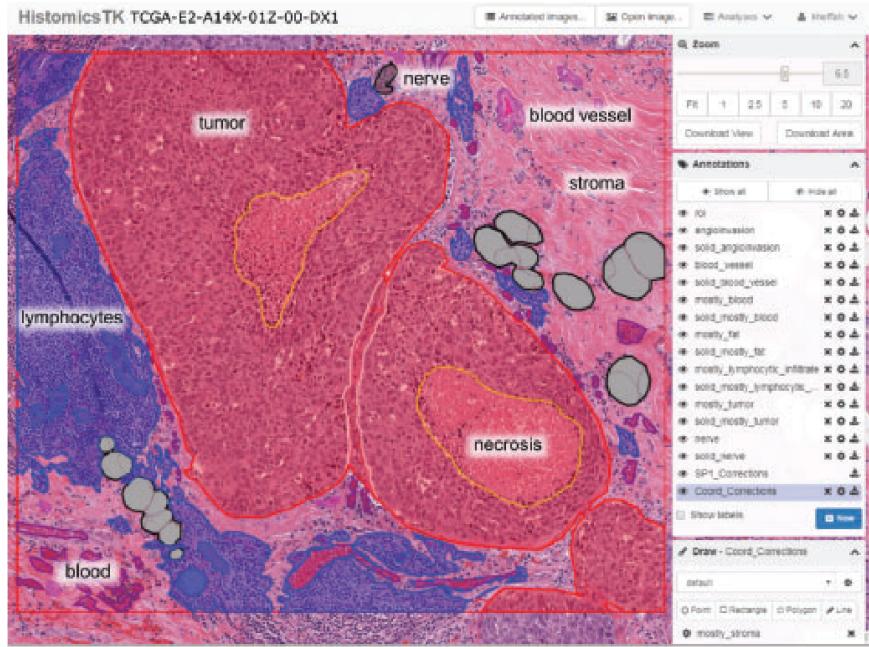


Figure 2-10 Screenshot of the DSA and HistomicsTK web interface while creating the crowdsourced annotations for the dataset presented by [Amgad et al., 2019].

961

962

CHAPTER

963

964

THREE

965

CHAINED GAUSSIAN PROCESSES

966

3.1 Background and Related Methods

967 In recent years, several approaches have been proposed to handle supervised
968 learning problems in the context of multiple annotators. This section reviews key
969 methods that form the foundation for understanding chained Gaussian processes
970 and their application to multiple annotator scenarios in regression and
971 classification tasks.

972

3.1.1 Kernel Alignment-Based Annotator Relevance 973 Analysis (KAAR)

974

The Kernel Alignment-Based Annotator Relevance Analysis (KAAR) method,
975 introduced in [Julián and Álvarez Meza Andrés Marino, 2023], addresses the
976 challenge of estimating annotator expertise in scenarios where the ground truth is

977 unavailable. The key innovation of KAAR lies in its use of Centered Kernel
 978 Alignment (CKA) to measure the similarity between input features and annotator
 979 labels.

980 Given a dataset with input features \mathbf{X} and labels from multiple annotators \mathbf{Y} , KAAR
 981 computes a kernel matrix \mathbf{K}_ν as a convex combination of R basis kernels:

$$\mathbf{K}_\nu = \sum_{r=1}^R \nu_r \mathbf{K}_r, \quad (3-1)$$

982 where \mathbf{K}_r is computed using a kernel function over the input features, and ν_r
 983 represents the relevance weight for the r -th annotator. The weights are optimized
 984 by maximizing the CKA between \mathbf{K}_ν and a target kernel \mathbf{F} computed from the
 985 input features:

$$\rho(\mathbf{K}_\nu, \mathbf{F}) = \frac{\langle \bar{\mathbf{K}}_\nu, \bar{\mathbf{F}} \rangle_F}{\|\bar{\mathbf{K}}_\nu\|_F \|\bar{\mathbf{F}}\|_F}, \quad (3-2)$$

986 A key advantage of KAAR is its ability to capture dependencies between annotators
 987 through the kernel framework, without requiring explicit parametric modeling of
 988 these relationships.

989 3.1.2 Localized Kernel Alignment-Based Annotator 990 Relevance Analysis (LKAAR)

991 LKAAR extends KAAR by introducing locality in the kernel alignment framework.
 992 Rather than assuming constant annotator performance across the input space,

993 LKAAR models the annotator relevance as a function of the input features. The
 994 kernel combination in LKAAR takes the form:

$$\mathbf{K}_q = \sum_{r=1}^R \mathbf{Q}_r \mathbf{K}_r \mathbf{Q}_r, \quad (3-3)$$

995 where \mathbf{Q}_r is a diagonal matrix whose elements represent the local relevance of
 996 annotator r . These elements are computed as:

$$q_r(\mathbf{x}_n) = \begin{cases} \beta_0^{(r)} + \sum_{n'=1}^N \beta_{n'}^{(r)} \kappa_\beta(\mathbf{x}_n, \mathbf{x}_{n'}), & \text{if } n \in N_r \\ 0, & \text{Otherwise} \end{cases} \quad (3-4)$$

997 This localized approach allows LKAAR to better model inconsistent annotators
 998 whose performance varies across different regions of the input space.

999 3.1.3 Regularized Chained Deep Neural Network 1000 (RCDNN)

1001 RCDNN represents a deep learning approach to multiple annotator learning,
 1002 inspired by the chained Gaussian processes model. The method models both the
 1003 ground truth estimation and annotator reliability through a deep neural network
 1004 architecture. The likelihood function in RCDNN takes the form:

$$p(\mathbf{Y}|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{r \in R_n} \left(\prod_{k=1}^K \zeta_{n,k}^{\delta(y_n^{(r)} - k)} \right)^{\lambda_n^{(r)}} \left(\frac{1}{K} \right)^{1-\lambda_n^{(r)}}, \quad (3-5)$$

1005 where $\lambda_n^{(r)}$ represents the reliability of annotator r for instance n , and $\zeta_{n,k}$ is the
 1006 estimated probability of class k for instance n .

1007 The RCDNN architecture is built upon a sophisticated multi-layer structure that
 1008 combines several key innovations in deep learning. At its core, it employs multiple
 1009 dense layers with carefully chosen activation functions, each designed to capture
 1010 different aspects of the annotator-data relationship. To prevent overfitting, a
 1011 comprehensive regularization strategy is implemented, combining both l1 and l2
 1012 norms with dropout mechanisms. The network architecture splits into specialized
 1013 branches: one dedicated to estimating the ground truth labels, and another
 1014 focused on modeling annotator reliability. This dual-branch approach allows the
 1015 network to simultaneously learn both the true underlying patterns in the data and
 1016 the individual characteristics of each annotator. Furthermore, the implementation
 1017 of Monte Carlo dropout during prediction provides a mechanism for uncertainty
 1018 estimation, offering not just predictions but also confidence measures for those
 1019 predictions.

1020 3.1.4 Chained Gaussian Processes Model

1021 The Chained Gaussian Processes (CGP) model provides a probabilistic framework for
 1022 modeling multiple outputs through a chain of Gaussian processes. In the context of
 1023 multiple annotators, CGP can be used to model both the ground truth and annotator
 1024 behavior. The key idea is to model a likelihood function with J parameters as in [?]:

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n|\theta_1(\mathbf{x}_n), \dots, \theta_J(\mathbf{x}_n)), \quad (3-6)$$

1025 where each parameter $\theta_j(\mathbf{x})$ is modeled using a separate Gaussian process. This
 1026 framework provides remarkable flexibility in modeling complex relationships

1027 within the data. The use of Gaussian processes enables natural uncertainty
1028 quantification through posterior distributions, allowing the model to express
1029 varying degrees of confidence in its predictions. The framework inherently
1030 handles missing data through its probabilistic nature, and the automatic relevance
1031 determination property of GPs helps identify the most important features for
1032 prediction.

1033 The CGP framework becomes particularly powerful when integrated with concepts
1034 from KAAR and LKAAR. By incorporating kernel alignment principles, the model
1035 can make more informed choices about kernel functions, better capturing the
1036 underlying structure of the data. The framework can be extended to include
1037 input-dependent lengthscales, allowing for varying degrees of smoothness across
1038 the input space. Furthermore, through multi-output GP formulations, the model
1039 can capture complex dependencies between different annotators, providing a rich
1040 representation of how different experts' opinions relate to each other.

1041 However, implementing CGP in practice presents several significant challenges.
1042 The computational complexity of the model scales cubically with the dataset size,
1043 making it potentially prohibitive for large-scale applications. The need for
1044 approximate inference techniques in non-Gaussian likelihood scenarios can
1045 introduce additional complexity and potential approximation errors. When
1046 dealing with high-dimensional inputs, the model may struggle with the curse of
1047 dimensionality, requiring careful feature selection or dimensionality reduction
1048 preprocessing. Additionally, the selection of appropriate kernel functions remains
1049 a critical challenge, as it significantly impacts the model's performance and its
1050 ability to capture relevant patterns in the data.

1051 3.2 Tooling and Implementation

1052 The implementation of Chained Gaussian Processes for multiple annotator
1053 learning requires robust and flexible software tools that can handle both the

1054 probabilistic nature of Gaussian processes and the complexity of deep
1055 architectures. In our implementation, we leverage two key frameworks: GPflow
1056 and GPflux. These tools provide the foundation for building scalable and efficient
1057 implementations of our models while maintaining the theoretical rigor necessary
1058 for proper uncertainty quantification.

1059 The choice of these frameworks is motivated by several factors. First, both are
1060 built on top of TensorFlow, providing access to automatic differentiation and GPU
1061 acceleration capabilities. Second, they offer a modular design that allows for easy
1062 extension and modification, which is crucial when implementing novel model
1063 architectures. Third, they provide robust implementations of various Gaussian
1064 process models and inference methods, allowing us to focus on the specific
1065 challenges of multiple annotator learning rather than reimplementing basic GP
1066 functionality.

1067 3.2.1 GPflow Overview

1068 GPflow is a Python package for building Gaussian process models in TensorFlow.
1069 Developed as a modern and scalable alternative to GPy, it leverages TensorFlow's
1070 computational capabilities to provide efficient implementation of Gaussian process
1071 models. The framework is designed with four key principles in mind: ease of use,
1072 computational efficiency, extensibility, and automatic differentiation. This work
1073 was originally published in [[de G. Matthews et al., 2016](#)].

1074 At its core, GPflow provides implementations of various Gaussian process models
1075 including GPR (Gaussian Process Regression), GPMC (Gaussian Process Monte
1076 Carlo), SGPR (Sparse Gaussian Process Regression), and SGPMC (Sparse Gaussian
1077 Process Monte Carlo). These implementations are built on top of TensorFlow's
1078 computational graph framework, allowing for automatic differentiation and GPU
1079 acceleration where available.

1080 The framework’s architecture is modular, with clear separation between model
1081 components such as kernels, likelihoods, and optimization routines. This
1082 modularity makes it straightforward to extend existing models or create new ones.
1083 GPflow includes a comprehensive suite of kernel functions, ranging from the basic
1084 RBF (Radial Basis Function) kernel to more sophisticated ones like the Spectral
1085 Mixture kernel.

1086 One of GPflow’s distinguishing features is its robust handling of model parameters.
1087 It provides facilities for parameter transformation (ensuring parameters remain in
1088 valid ranges during optimization), prior specification, and flexible optimization
1089 strategies. The framework also includes utilities for model evaluation and
1090 prediction, with built-in support for uncertainty estimation.

1091 3.2.2 GPflux Framework

1092 GPflux extends GPflow to enable the construction and training of Deep Gaussian
1093 Process (DGP) models using modern deep learning practices. Built on top of both
1094 GPflow and TensorFlow, GPflux provides a flexible framework for implementing
1095 deep probabilistic models that combine the expressiveness of deep neural
1096 networks with the uncertainty quantification capabilities of Gaussian processes.
1097 This work was originally published in [Dutordoir et al., 2021].

1098 The framework implements Deep Gaussian Processes as a composition of Gaussian
1099 process layers, where each layer transforms its inputs through a Gaussian process
1100 mapping. This hierarchical structure allows the model to capture complex,
1101 non-stationary patterns in data while maintaining uncertainty estimates
1102 throughout the entire processing pipeline. GPflux handles this complexity through
1103 a carefully designed architecture that separates the concerns of model
1104 specification, inference, and prediction.

1105 A key innovation in GPflux is its integration with Keras-style layers, allowing
1106 practitioners to combine deterministic neural network layers with Gaussian

process layers in a single model. This hybrid approach enables the construction of models that can leverage both the structural assumptions encoded in neural network architectures and the flexibility of Gaussian processes.

For inference, GPflux implements various approximate methods, including stochastic variational inference with inducing points, which makes it possible to train deep Gaussian process models on large datasets. The framework provides built-in support for minibatch training, making it practical to work with datasets that don't fit in memory.

3.2.3 Implementation Considerations

When implementing Chained Gaussian Processes for multiple annotator learning using these tools, several key considerations must be addressed. First, the model architecture must be carefully designed to capture both the individual annotator characteristics and their interdependencies. This requires extending the basic GP layers provided by GPflux to handle multiple output streams - one for each annotator's reliability model.

Second, the inference procedure must be adapted to handle the unique challenges of multiple annotator data, such as missing labels and varying levels of annotator expertise across different regions of the input space. The variational inference capabilities of GPflux provide a solid foundation for this, but careful attention must be paid to the design of the variational approximation to ensure it captures the relevant uncertainty in both the ground truth and annotator reliability estimates.

Finally, computational efficiency becomes a critical concern when scaling to large datasets with multiple annotators. The minibatch training capabilities of both frameworks help address this, but additional considerations such as appropriate inducing point selection and batch size tuning become important for achieving good performance in practice.

¹¹³⁴ 3.3 Experimental Setup

¹¹³⁵ Synthetic datasets were used to evaluate the proposed methods in classification and
¹¹³⁶ regression tasks, following very similar procedures to the ones used in [?], which are
¹¹³⁷ as follows:

¹¹³⁸ 3.3.1 Classification

¹¹³⁹ A fully synthetic data for a 1-D ($P = 1$) multiclass classification problem with
¹¹⁴⁰ $K = 3$ classes is generated. The input feature matrix \mathbf{X} is constructed by randomly
¹¹⁴¹ sampling $N = 100$ points from a uniform distribution over the interval $[0, 1]$. For
¹¹⁴² each sample n , the true label is determined by computing three target values:

$$\begin{aligned} t_{n,1} &= \sin(2\pi x_n), \\ t_{n,2} &= -\sin(2\pi x_n), \\ t_{n,3} &= -\sin(2\pi(x_n + 0.25)) + 0.5, \end{aligned}$$

¹¹⁴³ and assigning the class label as the index $i \in \{1, 2, 3\}$ that maximizes $t_{n,i}$. This
¹¹⁴⁴ construction ensures a nontrivial, nonlinear decision boundary between the classes.

¹¹⁴⁵ For evaluation, the test set is generated by extracting 200 equally spaced samples
¹¹⁴⁶ from the interval $[0, 1]$, providing a dense coverage of the input space to assess
¹¹⁴⁷ generalization performance.

¹¹⁴⁸ Note that at this point, the fully synthetic dataset does not contain real annotator
¹¹⁴⁹ labels. Therefore, it is necessary to simulate annotator labels as corrupted versions
¹¹⁵⁰ of the hidden ground truth. In these experiments, the simulation is designed to
¹¹⁵¹ account for: 1) dependencies among annotators, and 2) labeler performance that
¹¹⁵² varies as a function of the input features.

¹¹⁵³ To achieve this, an **SLFM**-based approach (termed **SLFM-C**) is employed, which
¹¹⁵⁴ generates annotator labels as follows:

- 1155 1. **Definition of deterministic functions:** Define Q deterministic functions $\hat{\mu}_q :$
 1156 $\mathcal{X} \rightarrow \mathbb{R}$ and their combination parameters $\hat{w}_{l,r,q} \in \mathbb{R}$, for all $r \in R, n \in N$.
- 1157 2. **Annotator performance computation:** Compute $\hat{f}_{l,r,n} = \sum_{q=1}^Q \hat{w}_{l,r,q} \hat{\mu}_q(\hat{x}_n)$,
 1158 where $\hat{x}_n \in \mathbb{R}$ is the n th component of $\hat{\mathbf{x}} \in \mathbb{R}^N$. Here, $\hat{\mathbf{x}}$ is the 1-D
 1159 representation of the input features in \mathbf{X} , obtained using the t-distributed
 1160 Stochastic Neighbor Embedding (t-SNE) approach.
- 1161 3. **Reliability calculation:** Calculate $\hat{\lambda}_n^r = \varsigma(\hat{f}_{l,r,n})$, where $\varsigma(\cdot) \in [0, 1]$ is the
 1162 sigmoid function.
- 1163 4. **Label assignment:** The r th annotator's label for the n th sample is given by

$$y_n^r = \begin{cases} y_n, & \text{if } \hat{\lambda}_n^r \geq 0.5 \\ \tilde{y}_n, & \text{if } \hat{\lambda}_n^r < 0.5 \end{cases}$$

1164 where y_n is the true label and \tilde{y}_n is a flipped version of the actual label y_n .
 1165 This procedure ensures that the simulated annotator labels exhibit both
 1166 inter-annotator dependencies and input-dependent reliability.

1167 3.3.2 Regression

1168 A fully synthetic dataset is also generated for a 1-D regression problem, where the
 1169 ground truth for the n th sample is defined as:

$$y_n = \sin(2\pi x_n) \sin(6\pi x_n)$$

1170 Here, the input matrix \mathbf{X} is constructed by randomly sampling 100 points uniformly
 1171 from the interval $[0, 1]$. The test set is created by extracting equally spaced samples

1172 from the same interval, ensuring a clear separation between training and testing
1173 data.

1174 To simulate the presence of multiple annotators, the ground truth is corrupted with
1175 annotator-specific noise. For each annotator r , the observed label for the n th sample
1176 is given by:

$$y_n^r = y_n + \epsilon_n^r$$

1177 where $\epsilon_n^r \sim \mathcal{N}(0, v_n^r)$ is Gaussian noise, and v_n^r represents the error variance for
1178 the r th annotator, which may depend on the input features. This setup allows us
1179 to model both the bias and the reliability of each annotator, as well as potential
1180 dependencies between annotators.

1181 To further increase realism, the error variance v_n^r can be modeled as a function of
1182 the input features, capturing scenarios where annotator reliability varies across the
1183 input space, thus:

$$\epsilon_n^r \sim \mathcal{N}(0, v_n^r)$$

1184 being v_n^r , the r th annotator's error variance for the n th sample.

1185 Again, for modeling the error as a function of the input features, an SLFM-based
1186 approach (termed SLFM-R) is employed to build the noisy labels, as follows:

- 1187 **1. Definition of deterministic functions:** Define Q deterministic functions $\hat{\mu}_q : \mathcal{X} \rightarrow \mathbb{R}$ and their combination parameters $\hat{w}_{l,r,q} \in \mathbb{R}$, for all $r \in R, n \in N$.
- 1189 **2. Annotator performance computation:** Compute $\hat{f}_{l,r,n} = \sum_{q=1}^Q \hat{w}_{l,r,q} \hat{\mu}_q(\hat{x}_n)$,
1190 where $\hat{x}_n \in \mathbb{R}$ is the n th component of $\hat{\mathbf{x}} \in \mathbb{R}^N$. Here, $\hat{\mathbf{x}}$ is the 1-D
1191 representation of the input features in \mathbf{X} , obtained using the t-distributed
1192 Stochastic Neighbor Embedding (t-SNE) approach.
- 1193 **3. Compute variance :** Finally, determine $v_n^r = \exp(\hat{f}_{l,r,n})$.

₁₁₉₄ **3.4 Summary**

₁₁₉₅ **3.4.1 Comparative Analysis**

₁₁₉₆ The evolution of multiple annotator learning methods represents a fascinating
₁₁₉₇ progression in how we approach the challenge of learning from diverse expert
₁₁₉₈ opinions. KAAR established a foundation with its non-parametric approach
₁₁₉₉ through kernel alignment, though it operates under the assumption of constant
₁₂₀₀ annotator performance across the input space. LKAAR advanced this concept by
₁₂₀₁ introducing input-dependent annotator behavior, recognizing that expert
₁₂₀₂ performance often varies across different types of inputs. RCDNN brought the
₁₂₀₃ power and flexibility of deep learning to the problem, offering excellent scalability
₁₂₀₄ but requiring careful attention to regularization and architecture design. CGP
₁₂₀₅ represents perhaps the most theoretically elegant approach, providing a principled
₁₂₀₆ probabilistic framework, though it faces significant computational challenges in
₁₂₀₇ practical applications.

₁₂₀₈ The selection of an appropriate method for a given application requires careful
₁₂₀₉ consideration of several factors. The size and dimensionality of the dataset play a
₁₂₁₀ crucial role, as some methods scale better than others with increasing data
₁₂₁₁ volume. The need for uncertainty quantification may favor probabilistic
₁₂₁₂ approaches like CGP, while computational constraints might push towards more
₁₂₁₃ efficient methods like KAAR. The presence of input-dependent annotator behavior
₁₂₁₄ could make LKAAR or RCDNN more appropriate choices. Additionally, when
₁₂₁₅ interpretability is a key requirement, as in many medical applications, methods
₁₂₁₆ like KAAR and LKAAR might be preferred over the more complex deep learning
₁₂₁₇ approaches. Table 3-1 provides a comprehensive comparison of the key
₁₂₁₈ characteristics and capabilities of each method discussed.

Table 3-1 Comparison of Multiple Annotator Learning Methods

Characteristic	KAAR	LKAAR	RCDNN	CGP
Input-dependent annotator behavior	No	Yes	Yes	Yes
Annotator dependencies modeling	Yes	Yes	Yes	Yes
Uncertainty quantification	No	No	Partial	Yes
Computational scalability	High	Medium	High	Low
Missing data handling	Yes	Yes	Yes	Yes
Non-linear relationships	Yes	Yes	Yes	Yes
Training complexity	Low	Medium	High	High
Hyperparameter sensitivity	Low	Medium	High	Medium
Interpretability	High	Medium	Low	Medium
Memory requirements	Low	Medium	High	High

Note: Partial uncertainty quantification in RCDNN refers to the use of Monte Carlo dropout for approximate uncertainty estimation.

1219 3.5 Conclusion

1220 In this chapter, we have presented a comprehensive overview of the Chained
 1221 Gaussian Processes model and its application to multiple annotator learning. We
 1222 have discussed the theoretical underpinnings of the model, its implementation
 1223 using GPflow and GPflux, and its evaluation on synthetic datasets. It was shown in
 1224 [?] that CCGPMA outperforms other method in state of the art methods in terms of
 1225 accuracy and uncertainty quantification (see 1.3.1).

1226

1227

CHAPTER

1228

FOUR

1229

1230

TRUNCATED GENERALIZED CROSS ENTROPY FOR SEGMENTATION

1231

1232

4.1 Loss functions for multiple annotators

1233 As mentioned in Section 2.3.2, a loss function is a key element for defining the
1234 objective function of a deep learning model. The categorical cross-entropy loss is a
1235 common loss function for classification tasks. However, in the case of multiple
1236 annotators, the categorical cross-entropy loss is not able to handle the varying
1237 reliability of the annotators. In this section, we will propose a loss function that is
1238 able to handle multiple annotators' segmentation masks while accounting for their
1239 varying reliability across different regions of the image.

1240 4.1.1 Generalized Cross Entropy

1241 The Generalized Cross Entropy (GCE) loss function was first introduced by [Zhang
1242 and Sabuncu, 2018] as a robust alternative to the standard cross-entropy loss,
1243 particularly effective in handling noisy labels. Let us first consider the Cross
1244 Entropy (CE) and Mean Absolute Error (MAE) loss functions:

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

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

1245 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, $\mathbf{1}^\top \mathbf{y} =$
1246 $\mathbf{1}^\top f(\mathbf{x}) = 1$, $\mathbf{1} \in \{1\}^K$ being an all-ones vector. In addition, the MAE loss can be
1247 rewritten for softmax outputs, yielding:

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

1248 where \odot stands for the element-wise product.

1249 The CE loss function exhibits several distinct characteristics that make it
1250 particularly sensitive to noisy labels. It is unbounded from above and heavily
1251 penalizes confident but wrong predictions, which can lead to instability in the
1252 presence of label noise. In contrast, the MAE loss offers a more robust alternative
1253 with its bounded nature and symmetric properties in softmax-based
1254 representations. While the MAE loss is more resilient to noisy labels, it tends to

1255 train more slowly due to its equal weighting of all mistakes regardless of
1256 confidence.

1257 The GCE loss function, as defined by the authors in [Zhang and Sabuncu, 2018],
1258 provides a flexible framework that bridges these two approaches:

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

1259 with $q \in (0, 1]$. Remarkably, the limiting case for $q \rightarrow 0$ in GCE is equivalent to the
1260 CE expression, and when $q = 1$, it equals the MAE loss. In addition, the GCE holds
1261 the following gradient with regard to θ :

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

1262 The GCE loss combines the best of both worlds, offering robustness to label noise
1263 while maintaining the convexity property necessary for optimization. The
1264 truncation parameter q provides a mechanism to control the sensitivity to outliers,
1265 allowing for fine-tuning of the loss function's behavior based on the specific
1266 characteristics of the dataset.

1267 4.1.2 Extension to Multiple Annotators

1268 In the context of multiple annotators, we need to consider the varying reliability
1269 of each annotator across different regions of the image. Let's consider a k -class
1270 multiple annotators segmentation problem with the following data representation:

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

¹²⁷¹ where the segmentation mask function maps the input to output as:

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

¹²⁷² The segmentation masks \mathbf{Y}_r satisfy the following condition for being a softmax-like
¹²⁷³ representation:

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

¹²⁷⁴ 4.1.3 Reliability Maps and Truncated GCE

¹²⁷⁵ The key innovation in our approach is the introduction of reliability maps Λ_r for
¹²⁷⁶ each annotator:

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

¹²⁷⁷ These reliability maps serve as a sophisticated mechanism to estimate the
¹²⁷⁸ confidence of each annotator at every spatial location (w, h) in the image. By
¹²⁷⁹ learning these maps jointly with the segmentation model, the network gains the
¹²⁸⁰ ability to adapt to the varying levels of expertise across different regions of the
¹²⁸¹ image. This approach allows for dynamic weighting of each annotator's

contribution based on their estimated reliability in specific areas, effectively handling cases where annotators might demonstrate different levels of expertise in different parts of the image.

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

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

where $q \in (0, 1)$ controls the MAE or CE level in the same way as in the GCE loss function. 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:

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

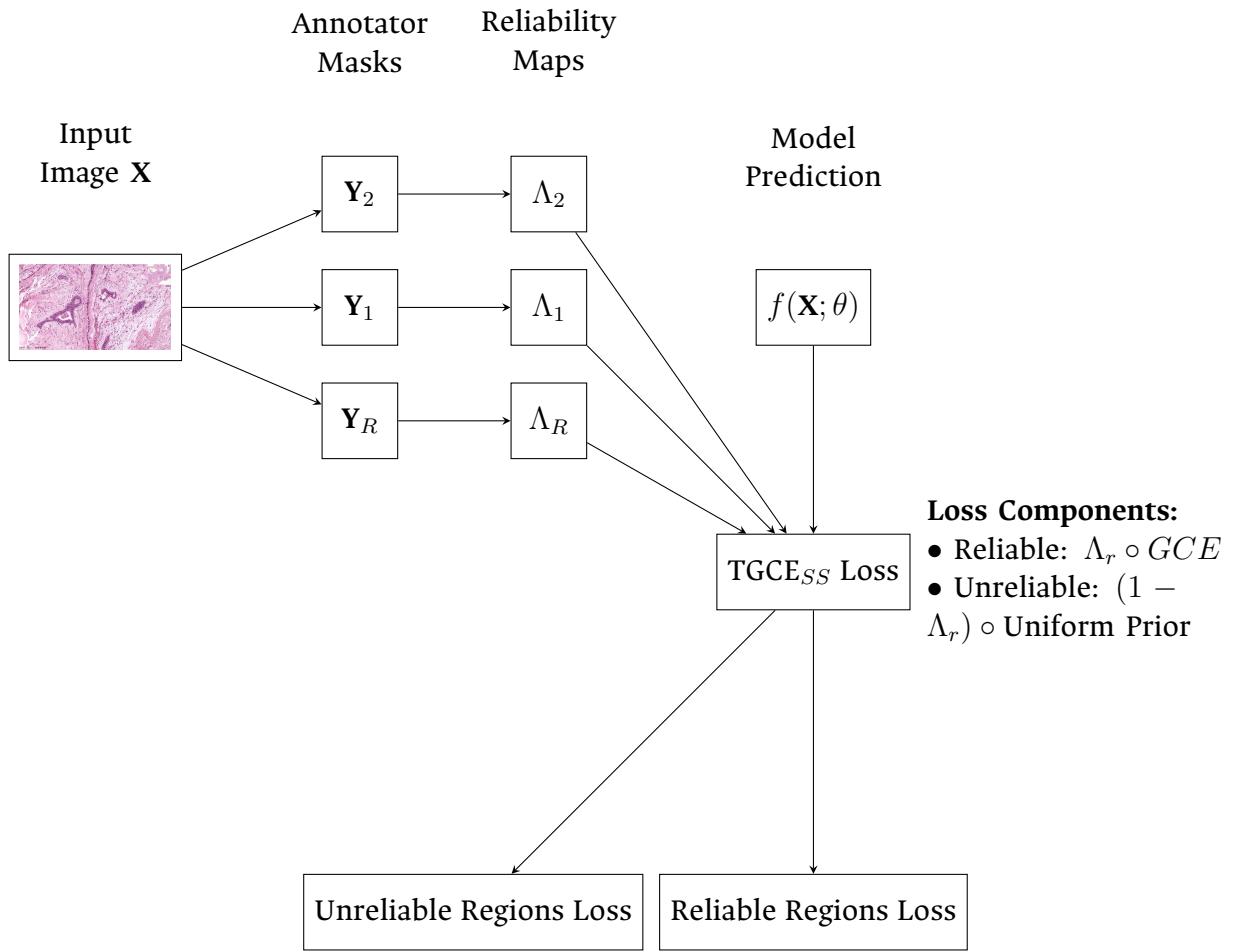


Figure 4-1 Working mechanism of the proposed Truncated Generalized Cross Entropy for Semantic Segmentation (TGCE_{SS}) loss function. The loss combines reliability maps Λ_r with the model's predictions $f(\mathbf{X}; \theta)$ to compute a weighted loss that accounts for annotator reliability across different image regions.

1294

1295 CHAPTER

1296

FIVE

1297

1298

CHAINED DEEP LEARNING FOR IMAGE
SEGMENTATION

1299

1300

5.1 Introduction

1301 As mentioned in Chapter 1, U-shaped CNNs have been proven a good solution for
1302 segmentation tasks in medical images, due to their ability to capture both global
1303 and local information with relatively low datasets.

1304

5.2 Using U-NET as a building block

1305 Our proposed model architecture combines the strengths of UNET with a ResNet-
1306 34 backbone, specifically designed to work with the $TGCE_{SS}$ loss function. The
1307 architecture is illustrated in Figure 5-2.

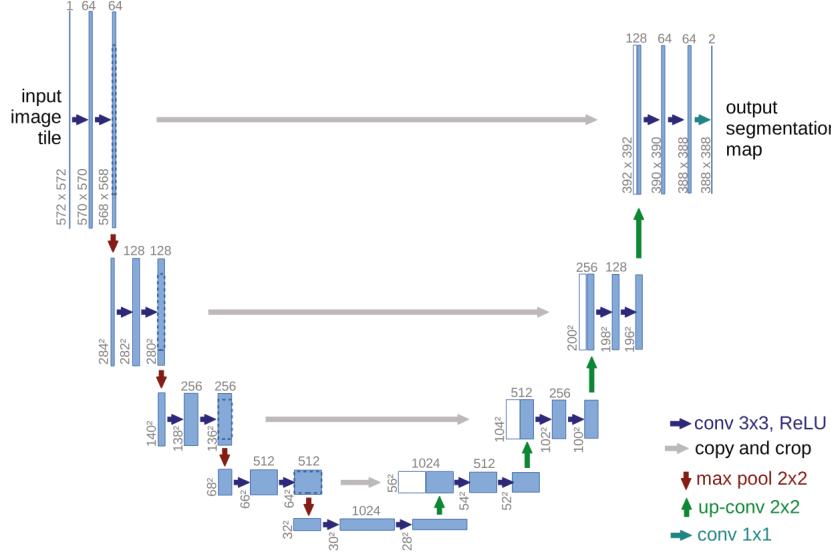


Figure 5-1 Original U-NET architecture.

5.2.1 Backbone Architecture

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

5.2.2 UNET Architecture

The UNET architecture follows a traditional encoder-decoder structure with skip connections, where the encoder path implements the ResNet-34 structure. The decoder path employs transposed convolutions for upsampling, creating a

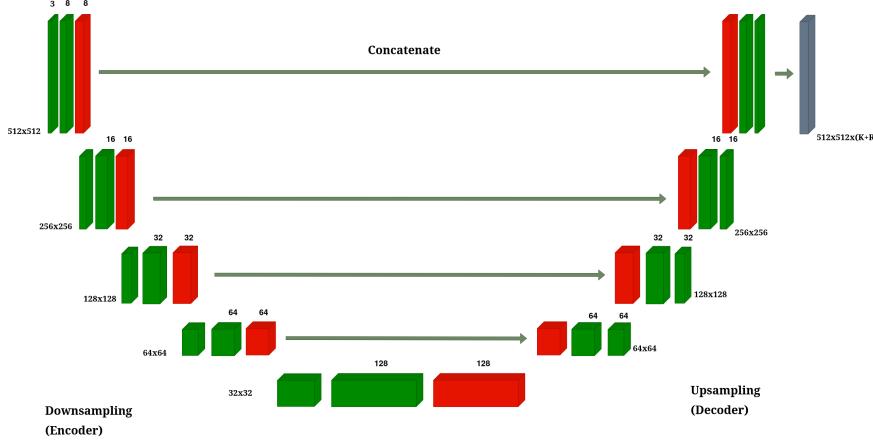


Figure 5-2 Overview of the proposed model architecture.

1321 symmetrical architecture that effectively captures both high-level and low-level
 1322 features. The architecture incorporates four downsampling stages in the encoder,
 1323 corresponding to the ResNet-34 blocks, and four upsampling stages in the decoder.
 1324 These stages are connected through skip connections that bridge corresponding
 1325 encoder and decoder stages, allowing the network to preserve fine-grained details.
 1326 Each convolution operation is followed by batch normalization and ReLU
 1327 activation to ensure stable training and effective feature learning.

1328 5.2.3 Reliability Map Branch

1329 A key innovation in our architecture is the parallel branch dedicated to estimating
 1330 reliability maps. This branch processes the same encoder features as the main
 1331 segmentation path but focuses on learning the confidence of each annotator.
 1332 Through a series of 1×1 convolutions, the branch reduces channel dimensions
 1333 while maintaining spatial information. The final output consists of R reliability
 1334 maps Λ_r , one for each annotator, with values constrained to the $[0, 1]$ range
 1335 through a sigmoid activation function. This design allows the network to learn and
 1336 adapt to the varying reliability of different annotators across different regions of
 1337 the image.

₁₃₃₈ **5.2.4 Integration with TGCE_{SS} Loss**

₁₃₃₉ The model produces two distinct outputs: segmentation masks $\hat{\mathbf{Y}} = f(\mathbf{X}; \theta)$ and
₁₃₄₀ reliability maps $\{\Lambda_r(\mathbf{X}; \theta)\}_{r=1}^R$. These outputs work in tandem with the TGCE_{SS}
₁₃₄₁ loss function described in Section ???. The loss function simultaneously guides the
₁₃₄₂ learning of both the segmentation masks and reliability maps, ensuring that the
₁₃₄₃ model learns to balance the contributions of different annotators based on their
₁₃₄₄ estimated reliability.

₁₃₄₅ **5.2.5 Training Process**

₁₃₄₆ The training process begins with the initialization of the ResNet-34 backbone
₁₃₄₇ using pre-trained weights, providing a strong foundation for feature extraction.
₁₃₄₈ The entire network is then trained end-to-end using the Adam optimizer with a
₁₃₄₉ learning rate of 10^{-4} . The TGCE_{SS} loss function plays a crucial role in updating
₁₃₅₀ both the segmentation and reliability branches, ensuring that the model learns to
₁₃₅₁ effectively handle multiple annotators' inputs while accounting for their varying
reliability.

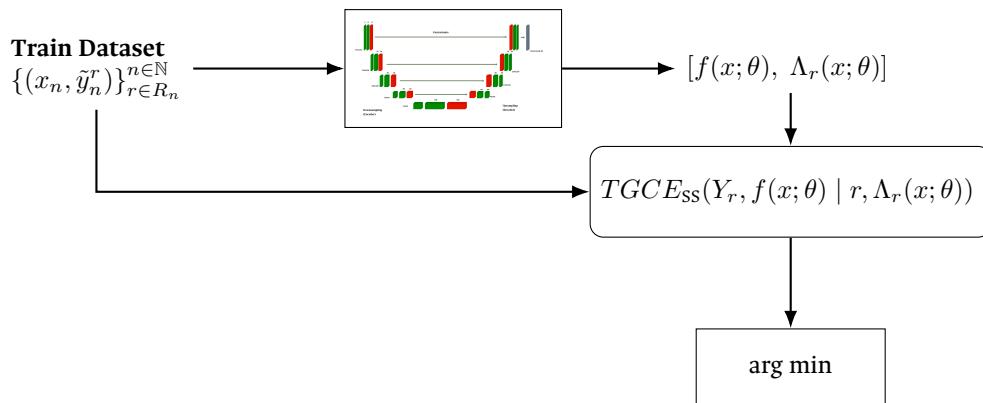


Figure 5-3 Training process of the proposed model overview. Estimated segmentation and reliability maps are computed for each input image, and the loss function is computed for the entire batch.

1352

1353 The model's architecture is specifically designed to address the challenges of
1354 multi-annotator segmentation. Through the ResNet-34 backbone, it learns robust
1355 segmentation features that capture high-level patterns in the data. The UNET's
1356 skip connections enable the preservation of fine-grained details, while the parallel
1357 reliability branch allows the model to adapt to annotator-specific characteristics.
1358 This comprehensive design enables the model to effectively handle multiple
1359 annotators' inputs while maintaining high segmentation accuracy and reliability
1360 estimation.

1361

1362

CHAPTER

1363

1364

SIX

1365

CONCLUSIONS

1366

6.1 Summary

1367 Throughout this work, several techniques have been explored to develop a model
1368 that can face ISS of real world histopathology images in crowdsourcing scenarios.
1369 Different approaches were studied and evaluated to reach an optimal solution
1370 which could model the annotators' performance from the input space. The
1371 following remarks are based on the results obtained from the experiments
1372 conducted and reviewed literature:

1373

1374

1375

1376

1377

- Chained Gaussian Processes approaches are a robust approach for developing classification and regression tasks, however, their computational powerful requirements make them unsuitable for image segmentation tasks, specially since deep learning approaches have shown to be more efficient in the task of extracting features from images.

- 1378 • Existing loss functions in the state of the art were insufficient to model the
1379 annotators' performance across the input space, from which the need for a
1380 new loss function which fulfils the requirements of the task was identified.
- 1381 • A novel loss function was proposed and evaluated, showing to be a good
1382 alternative to the existing loss functions in the state of the art, and providing
1383 an efficient way to optimize parameters for reaching capability to model the
1384 annotators' performance across the input space.
- 1385 • The use of U-Shaped Deep Learning models as a building block for the
1386 proposed solution was a valid approach, as it allows for the model to learn
1387 the spatial relationships between the pixels in the image, and the
1388 annotations, and to use this information to make predictions about the
1389 annotations of new images.
- 1390 • Chained Deep Learning approaches in combination with the proposed loss
1391 function were the best approach for the task, outperforming the state of the
1392 art in the task of **ISS** of histopathology images.

1393 6.2 Future work

1394 As tools and frameworks for bayesian approaches like Chained Gaussian Processes
1395 ¹ evolve, it is expected to see a decrease in computational requirements and hence,
1396 making mixed models more appealing for the task of **ISS** of histopathology images.
1397 In that sense, the use of **CCGPMA** as a building block for the proposed solution could
1398 be a good approach, as it is a powerful model that can take advantage of the inter-
1399 dependencies between the annotators, and the powerfulness of **CNNs** to extract
1400 features from the images.

¹Including GPFlow, and GPFlux.

BIBLIOGRAPHY

- 1402 [Abadi et al., 2016] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J.,
1403 Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga,
1404 R., Moore, S., Murray, D. G., Steiner, B., Tucker, P. A., Vasudevan, V., Warden,
1405 P., Wicke, M., Yu, Y., and Zhang, X. (2016). Tensorflow: A system for large-scale
1406 machine learning. CoRR, abs/1605.08695. (page 35)
- 1407 [Amgad et al., 2019] Amgad, M., Elfandy, H., Hussein, H., Atteya, L. A., Elsebaie,
1408 M. A. T., Abo Elnasr, L. S., Sakr, R. A., Salem, H. S. E., Ismail, A. F., Saad,
1409 A. M., Ahmed, J., Elsebaie, M. A. T., Rahman, M., Ruhban, I. A., Elgazar, N. M.,
1410 Alagha, Y., Osman, M. H., Alhusseiny, A. M., Khalaf, M. M., Younes, A.-A. F.,
1411 Abdulkarim, A., Younes, D. M., Gadallah, A. M., Elkashash, A. M., Fala, S. Y., Zaki,
1412 B. M., Beezley, J., Chittajallu, D. R., Manthey, D., Gutman, D. A., and Cooper, L.
1413 A. D. (2019). Structured crowdsourcing enables convolutional segmentation of
1414 histology images. *Bioinformatics*, 35(18):3461–3467. (pages xviii, 29, 41, and 42)
- 1415 [Avanzo et al., 2024] Avanzo, M., Stancanello, J., Pirrone, G., Drigo, A., and Retico,
1416 A. (2024). The evolution of artificial intelligence in medical imaging: From
1417 computer science to machine and deep learning. *Cancers (Basel)*, 16(21):3702.
1418 Author Joseph Stancanello is employed by Elekta SA. The remaining authors
1419 declare no commercial or financial conflicts of interest. (page 3)

- 1420 [Azad et al., 2024] Azad, R., Aghdam, E. K., Rauland, A., Jia, Y., Avval, A. H.,
1421 Bozorgpour, A., Karimijafarbigloo, S., Cohen, J. P., Adeli, E., and Merhof, D.
1422 (2024). Medical image segmentation review: The success of u-net. *IEEE
1423 Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10076–10095.
1424 (page 2)
- 1425 [Banerjee et al., 2025] Banerjee, A., Shan, H., and Feng, R. (2025). Editorial:
1426 Artificial intelligence applications for cancer diagnosis in radiology. *Frontiers in
1427 Radiology*, 5. (page 8)
- 1428 [Bhalgat et al., 2018] Bhalgat, Y., Shah, M. P., and Awate, S. P. (2018). Annotation-
1429 cost minimization for medical image segmentation using suggestive mixed
1430 supervision fully convolutional networks. *CoRR*, abs/1812.11302. (page 3)
- 1431 [Brito-Pacheco et al., 2025] Brito-Pacheco, D., Giannopoulos, P., and Reyes-
1432 Aldasoro, C. C. (2025). Persistent homology in medical image processing: A
1433 literature review. (page 2)
- 1434 [Carmo et al., 2025] Carmo, D. S., Pezzulo, A. A., Villacreses, R. A., Eisenbeisz,
1435 M. L., Anderson, R. L., Van Dorin, S. E., Rittner, L., Lotufo, R. A., Gerard, S. E.,
1436 Reinhardt, J. M., and Comellas, A. P. (2025). Manual segmentation of opacities
1437 and consolidations on ct of long covid patients from multiple annotators. *Scientific
1438 Data*, 12(1):402. (page 9)
- 1439 [de G. Matthews et al., 2016] de G. Matthews, A. G., van der Wilk, M., Nickson, T.,
1440 Fujii, K., Boukouvalas, A., León-Villagrá, P., Ghahramani, Z., and Hensman, J.
1441 (2016). Gpflow: A gaussian process library using tensorflow. (page 48)
- 1442 [Dutordoir et al., 2021] Dutordoir, V., Salimbeni, H., Hambro, E., McLeod, J.,
1443 Leibfried, F., Artemev, A., van der Wilk, M., Hensman, J., Deisenroth, M. P., and
1444 John, S. (2021). Gpflux: A library for deep gaussian processes. (page 49)
- 1445 [Elhaminia et al., 2025] Elhaminia, B., Alsalemi, A., Nasir, E., Jahanifar, M., Awan,
1446 R., Young, L. S., Rajpoot, N. M., Minhas, F., and Raza, S. E. A. (2025). From
1447 traditional to deep learning approaches in whole slide image registration: A
1448 methodological review. (page 2)

- 1477 [Karthikeyan et al., 2023] Karthikeyan, R., McDonald, A., and Mehta, R. (2023).
1478 What's in a label? annotation differences in forecasting mental fatigue using ecg
1479 data and seq2seq architectures. (page 9)
- 1480 [Kim et al., 2024] Kim, Y., Lee, E., Lee, Y., and Oh, U. (2024). Understanding
1481 novice's annotation process for 3d semantic segmentation task with human-
1482 in-the-loop. In *Proceedings of the 29th International Conference on Intelligent User
1483 Interfaces, IUI '24*, page 444–454, New York, NY, USA. Association for Computing
1484 Machinery. (page 9)
- 1485 [Lam and Suen, 1997] Lam, L. and Suen, S. (1997). Application of majority
1486 voting to pattern recognition: an analysis of its behavior and performance.
1487 *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*,
1488 27(5):553–568. (page 11)
- 1489 [Lin et al., 2024] Lin, Y., Lian, A., Liao, M., and Yuan, S. (2024). Bcdnet: A fast
1490 residual neural network for invasive ductal carcinoma detection. (page 6)
- 1491 [López-Pérez et al., 2023] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A. D.,
1492 Molina, R., and Katsaggelos, A. K. (2023). Crowdsourcing segmentation
1493 of histopathological images using annotations provided by medical students.
1494 In Juarez, J. M., Marcos, M., Stiglic, G., and Tucker, A., editors, *Artificial
1495 Intelligence in Medicine*, pages 245–249, Cham. Springer Nature Switzerland.
1496 (pages 5, 6, 8, and 11)
- 1497 [Lu et al., 2023] Lu, X., Ratcliffe, D., Kao, T.-T., Tikhonov, A., Litchfield, L., Rodger,
1498 C., and Wang, K. (2023). Rethinking quality assurance for crowdsourced multi-
1499 roi image segmentation. *Proceedings of the AAAI Conference on Human Computation
1500 and Crowdsourcing*, 11(1):103–114. (pages 5 and 8)
- 1501 [López-Pérez et al., 2024] López-Pérez, M., Morales-Álvarez, P., Cooper, L. A.,
1502 Felicelli, C., Goldstein, J., Vadasz, B., Molina, R., and Katsaggelos, A. K. (2024).
1503 Learning from crowds for automated histopathological image segmentation.
1504 *Computerized Medical Imaging and Graphics*, 112:102327. (pages xvii, 5, 15, and 17)

- 1505 [Mazzarini et al., 2021] Mazzarini, M., Falchi, M., Bani, D., and Migliaccio, A. R.
1506 (2021). Evolution and new frontiers of histology in bio-medical research.
1507 *Microscopy Research and Technique*, 84(2):217–237. (pages xvii and 29)
- 1508 [Pan et al., 2021] Pan, X., Lu, Y., Lan, R., Liu, Z., Qin, Z., Wang, H., and Liu, Z. (2021).
1509 Mitosis detection techniques in h&e stained breast cancer pathological images:
1510 A comprehensive review. *Computers & Electrical Engineering*, 91:107038. (page 31)
- 1511 [Panayides et al., 2020] Panayides, A. S., Amini, A., Filipovic, N. D., Sharma, A.,
1512 Tsaftaris, S. A., Young, A., Foran, D., Do, N., Golemati, S., Kurc, T., Huang, K.,
1513 Nikita, K. S., Veasey, B. P., Zervakis, M., Saltz, J. H., and Pattichis, C. S. (2020). Ai
1514 in medical imaging informatics: Current challenges and future directions. *IEEE*
1515 *Journal of Biomedical and Health Informatics*, 24(7):1837–1857. (page 2)
- 1516 [Parkhi et al., 2012] Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar, C. V.
1517 (2012). Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*.
1518 (page 38)
- 1519 [Qiu et al., 2022] Qiu, Y., Hu, Y., Kong, P., Xie, H., Zhang, X., Cao, J., Wang, T.,
1520 and Lei, B. (2022). Automatic prostate gleason grading using pyramid semantic
1521 parsing network in digital histopathology. *Frontiers in Oncology*, 12. (page 14)
- 1522 [Rashmi et al., 2021] Rashmi, R., Prasad, K., and Udupa, C. B. K. (2021).
1523 Breast histopathological image analysis using image processing techniques for
1524 diagnostic purposes: A methodological review. *Journal of Medical Systems*, 46(1):7.
1525 (pages 1 and 5)
- 1526 [Ronneberger et al., 2015] Ronneberger, O., Fischer, P., and Brox, T. (2015). U-
1527 net: Convolutional networks for biomedical image segmentation. In Navab, N.,
1528 Hornegger, J., Wells, W. M., and Frangi, A. F., editors, *Medical Image Computing*
1529 and *Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham. Springer
1530 International Publishing. (page 14)

- 1531 [Ryou et al., 2025] Ryou, H., Thomas, E., Wojciechowska, M., Harding, L., Tam,
1532 K. H., Wang, R., Hu, X., Rittscher, J., Cooper, R., and Royston, D. (2025). Reticulin-
1533 free quantitation of bone marrow fibrosis in mpns: Utility and applications.
1534 *eJHaem*, 6(2):e70005. (page 2)
- 1535 [Sarvamangala and Kulkarni, 2022] Sarvamangala, D. R. and Kulkarni, R. V. (2022).
1536 Convolutional neural networks in medical image understanding: a survey.
1537 *Evolutionary Intelligence*, 15(1):1-22. (pages 3 and 6)
- 1538 [Shah et al., 2018] Shah, M. P., Merchant, S. N., and Awate, S. P. (2018).
1539 Ms-net: Mixed-supervision fully-convolutional networks for full-resolution
1540 segmentation. In Frangi, A. F., Schnabel, J. A., Davatzikos, C., Alberola-
1541 López, C., and Fichtinger, G., editors, *Medical Image Computing and Computer
1542 Assisted Intervention - MICCAI 2018*, pages 379–387, Cham. Springer International
1543 Publishing. (page 5)
- 1544 [Shalf, 2020] Shalf, J. (2020). The future of computing beyond moore’s law.
1545 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering
1546 Sciences*, 378(2166):20190061. (page 3)
- 1547 [TIAN and Zhu, 2015] TIAN, T. and Zhu, J. (2015). Max-margin majority voting for
1548 learning from crowds. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and
1549 Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 28.
1550 Curran Associates, Inc. (page 12)
- 1551 [Triana-Martinez et al., 2023] Triana-Martinez, J. C., Gil-González, J., Fernandez-
1552 Gallego, J. A., Álvarez Meza, A. M., and Castellanos-Dominguez, C. G. (2023).
1553 Chained deep learning using generalized cross-entropy for multiple annotators
1554 classification. *Sensors*, 23(7). (page 20)
- 1555 [Warfield et al., 2004] Warfield, S., Zou, K., and Wells, W. (2004). Simultaneous
1556 truth and performance level estimation (staple): an algorithm for the validation
1557 of image segmentation. *IEEE Transactions on Medical Imaging*, 23(7):903–921.
1558 (page 12)

- 1559 [Weitz et al., 2023] Weitz, P., Valkonen, M., Solorzano, L., Carr, C., Kartasalo, K.,
1560 Boissin, C., Koivukoski, S., Kuusela, A., Rasic, D., Feng, Y., Sinius Pouplier,
1561 S., Sharma, A., Ledesma Eriksson, K., Latonen, L., Laenholm, A.-V., Hartman,
1562 J., Ruusuvuori, P., and Rantalainen, M. (2023). A multi-stain breast cancer
1563 histological whole-slide-image data set from routine diagnostics. *Scientific Data*,
1564 10(1):562. (pages xviii, 32, 40, and 41)
- 1565 [Xu et al., 2024] Xu, Y., Quan, R., Xu, W., Huang, Y., Chen, X., and Liu, F. (2024).
1566 Advances in medical image segmentation: A comprehensive review of traditional,
1567 deep learning and hybrid approaches. *Bioengineering*, 11(10). (pages 3 and 8)
- 1568 [Yu et al., 2025] Yu, J., Li, B., Pan, X., Shi, Z., Wang, H., Lan, R., and Luo, X. (2025).
1569 Semi-supervised gland segmentation via feature-enhanced contrastive learning
1570 and dual-consistency strategy. *IEEE Journal of Biomedical and Health Informatics*,
1571 pages 1-11. (page 2)
- 1572 [Zhang et al., 2020] Zhang, L., Tanno, R., Xu, M.-C., Jin, C., Jacob, J., Cicarrelli, O.,
1573 Barkhof, F., and Alexander, D. (2020). Disentangling human error from ground
1574 truth in segmentation of medical images. In Larochelle, H., Ranzato, M., Hadsell,
1575 R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*,
1576 volume 33, pages 15750–15762. Curran Associates, Inc. (page 16)
- 1577 [Zhang and Sabuncu, 2018] Zhang, Z. and Sabuncu, M. R. (2018). Generalized
1578 cross entropy loss for training deep neural networks with noisy labels.
1579 (pages 20, 58, and 59)
- 1580 [Zhao et al., 2020] Zhao, R., Qian, B., Zhang, X., Li, Y., Wei, R., Liu, Y., and Pan,
1581 Y. (2020). Rethinking dice loss for medical image segmentation. In 2020 IEEE
1582 International Conference on Data Mining (ICDM), pages 851–860. (page 20)
- 1583 [Zhou et al., 2021] Zhou, S. K., Greenspan, H., Davatzikos, C., Duncan, J. S.,
1584 Van Ginneken, B., Madabhushi, A., Prince, J. L., Rueckert, D., and Summers, R. M.
1585 (2021). A review of deep learning in medical imaging: Imaging traits, technology
1586 trends, case studies with progress highlights, and future promises. *Proceedings of
the IEEE*, 109(5):820–838. (pages 1 and 2)

- 1588 [Zhou et al., 2024] Zhou, Z., Gong, H., Hsieh, S., McCollough, C. H., and Yu, L.
1589 (2024). Image quality evaluation in deep-learning-based ct noise reduction using
1590 virtual imaging trial methods: Contrast-dependent spatial resolution. *Medical*
1591 *Physics*, 51(8):5399–5413. (page 9)