



Deep Learning for Automated Adequacy Assessment of Cervical Cytology Samples

VLADYSLAV MOSIICHUK

julho de 2022

POLITÉCNICO DO PORTO
INSTITUTO SUPERIOR DE ENGENHARIA DO PORTO

Deep Learning for Automated Adequacy Assessment of Cervical Cytology Samples

Vladyslav Mosiichuk

Master in Electrical and Computer Engineering
Specialization Area of Autonomous Systems



DEPARTAMENTO DE ENGENHARIA ELETROTÉCNICA
Instituto Superior de Engenharia do Porto

July, 2022

This dissertation partially satisfies the requirements of the Thesis/Dissertation course of the program Master in Electrical and Computer Engineering, Specialization Area of Autonomous Systems.

Candidate: Vladyslav Mosiichuk, No. 1160805, 1160805@isep.ipp.pt

Scientific Guidance: Paula Maria Marques Moura Gomes Viana, pmv@isep.ipp.pt, and Luís Filipe Caeiro Margalho Guerra Rosado, luis.rosado@fraunhofer.pt

Host Entity: Associação Fraunhofer Portugal Research



DEPARTAMENTO DE ENGENHARIA ELETROTÉCNICA
Instituto Superior de Engenharia do Porto
Rua Dr. António Bernardino de Almeida, 431, 4200-072 Porto

July, 2022

Acknowledgements

First and foremost, I would like to express my sincere gratitude to my thesis supervisors: to Doctor Luís Rosado for giving insightful information regarding key aspects of the project, for guiding me in the right direction and giving suggestions and hints whenever advice was needed; to Professor Paula Viana for helping managing the stress and being able to give me very useful feedback about writing; and to both for being extremely comprehensive and supportive through the course of writing of this dissertation.

I'm also extremely grateful to my girlfriend, who was always there for me, to encourage me and give me a boost to finish this dissertation. To my parents and my brother for the patience and support along this journey.

This work was done under the scope of the project Transparent Artificial Medical Intelligence (TAMI), co-funded by Portugal 2020 framed under the Operational Programme for Competitiveness and Internationalization (COMPETE 2020), Fundação para a Ciência and Technology (FCT), Carnegie Mellon University, and European Regional Development Fund under Grant 45905. And, under the scope of “CLARE: Computer-Aided Cervical Cancer Screening”, project with reference POCI-01-0145-FEDER028857 and also financially supported by FEDER through Operational Competitiveness Program – COMPETE 2020 and by National Funds through Foundation for Science and Technology FCT/MCTES.

Abstract

Cervical cancer has been among the most common causes of cancer death in women. Screening tests such as liquid-based cytology (LBC) were responsible for a substantial decrease in mortality rates. Still, visual examination of cervical cells on microscopic slides is a time-consuming, ambiguous and challenging task, aggravated by inadequate sample quality (e.g. low cellularity or the presence of obscuring factors like blood or inflammation). While most works in the literature are focused on the automated detection of cervical lesions to support diagnosis, to the best of our knowledge, none of them address the automated assessment of sample adequacy, as established by The Bethesda System (TBS) guidelines. This work proposes a new methodology for automated adequacy assessment of cervical cytology samples. Since the most common reason for rejecting samples is the low count of the squamous cell nuclei, our approach relies on a deep learning object detection model for the detection and counting of different types of nuclei present in LBC samples. A dataset of 41 samples with a total of 42387 nuclei manually annotated by experienced specialists was used, and after extensive system parameters tuning, the best solution proposed achieved promising results for the automated detection of squamous nuclei (AP of 82.4%, Accuracy of 79.8%, Recall of 73.8% and F1 score of 81.5%). Additionally, by merging the developed automated cell counting approach with the adequacy criteria stated by the TBS guidelines, we validated our approach by correctly classifying an entire subset of 12 samples as adequate or inadequate.

Keywords: Cervical Cancer , Cervical Cytology , Machine Learning , Deep Learning , Adequacy Assessment , Nuclei Detection.

Resumo

Cancro cervical é uma das causas mais comuns de morte por cancro entre as mulheres. Testes de triagem, como citologia em meio líquido, foram responsáveis por uma diminuição substancial nas taxas de mortalidade. Porém, o exame visual das células cervicais em lâminas microscópicas é uma tarefa demorada, ambígua e desafiadora, que ainda poderá ser agravada pela qualidade inadequada da amostra (por exemplo, baixa celularidade ou presença de fatores obscurecedores como sangue ou inflamação). Enquanto a maioria dos trabalhos na literatura estão focados na detecção automática de lesões cervicais para apoiar o diagnóstico, até onde sabemos, nenhum deles aborda a avaliação automática da adequabilidade da amostra, conforme estabelecido pelas diretrizes do The Bethesda System (TBS). Este trabalho propõe uma nova metodologia para avaliação automática de adequabilidade de amostras de citologia cervical. Como o motivo mais comum para a rejeição de amostras é a baixa celularidade de núcleos de células escamosas, a nossa abordagem irá basear-se num modelo *deep learning* de detecção de objetos para a detecção e contagem de diferentes tipos de núcleos presentes em amostras. Foi utilizado um conjunto de dados de 41 amostras com um total de 42387 núcleos anotados manualmente por especialistas experientes. Após o ajuste extensivo de parâmetros do sistema a melhor solução proposta alcançou resultados promissores para a detecção automática de núcleos escamosos (AP de 82,4 %, Precisão de 79,8 %, Recall de 73,8% e F1 Score de 81,5%).

Palavras-Chave: Cancro Cervical, Citologia Cervical, Machine Learning, Deep Learning, Avaliação de Adequabilidade, Detecção de Núcleos

Contents

List of Figures	ix
List of Tables	xiii
List of Acronyms	xv
1 Introduction	1
1.1 Context and Motivation	1
1.2 Research Goals	2
1.3 Research Contributions	3
1.4 Document Structure	3
2 Cervical Cancer	5
2.1 Cervical Cancer	5
2.2 Cervical Cancer Screening	6
2.3 The Bethesda System (TBS) guidelines	8
3 Computer Vision and Deep Learning: Background Concepts	13
3.1 Computer Vision	13
3.2 Machine Learning	15
3.2.1 Classification Task, Segmentation and Object Detection . . .	15
3.2.2 Neural Networks	16
3.2.3 Training Neural Network	19
3.2.4 Deep Learning	19
3.2.5 Convolutional Neural Networks	20
3.3 Computer Vision and Cervical Cytology: State-of-the-art Review . .	22
4 Previous work and CLARE dataset	25
4.1 μ SmartScope for cervical cytology	25
4.2 CLARE dataset	27
4.2.1 Missing annotations	33
4.2.2 Data particularities	36
Overlapping cells	36
Out-of-focus	36

Chromatic aberration	37
5 Methodology	39
5.1 Overview	39
5.2 Proposed approach	41
5.2.1 Tools and model selection	42
5.2.2 Pre-processing	43
Dataset preparation	43
Slicing	45
5.2.3 Hyperparameter tuning	47
5.2.4 Data manipulation	48
Semi-supervised training	48
Merging and eliminating classes	48
Correction of chromatic aberrations	49
5.3 Evaluation	49
5.3.1 Evaluation metrics	50
IoU	50
Precision	51
Recall	51
Accuracy	52
Specificity	52
F1 Score	52
Youden Index	52
5.3.2 Evaluation pipeline	52
6 Experiments and Results	55
6.1 Hardware details	55
6.2 Baseline parameters search	57
6.3 Hyperparameter tuning	60
6.4 Chromatic aberration correction test	62
6.5 Data-centric experiments	64
6.5.1 Merging and eliminating classes	64
6.5.2 Semi-supervised approach	65
6.6 Final script evaluation	66
6.6.1 Final system parametrization	66
Image level results	67
Sample level results:	68
7 Conclusion and Future Work	71
References	74

A	Detailed description of dataset split	83
B	Tensorboard charts: mAP for each class during training	85
C	Hyperparameters tuning tests and results on test set	89
D	Results for merging and eliminating classes approach	93
E	Results of testing Overlap, NMS threshold and IOU threshold for the final system	95
	E.1 Filtered table	95
	E.2 Complete table	101

List of Figures

2.1	Satisfactory 2.1a and Unsatisfactory 2.1b LBC ThinPrep preparations at 10x magnification [1].	10
2.2	Examples of transformation zones in Conventional preparation (CP) and Liquid-based preparation samples [1].	11
2.3	Unsatisfactory specimen due to obscuring white blood cells. It should be noted that when assessing the adequacy of the slide considering obscuring factors and cellularity, the cellularity refer only to well visualized cells [1].	12
3.1	Comparing a biological neuron to an artificial neuron [2].	16
3.2	Graph of a step function [2].	17
3.3	Sigmoid function graph 3.3a and Graph of the ReLU activation function 3.3b [2].	17
3.4	Example basic fully-connected neural network [2].	18
3.5	Visual depiction of passing image data through a neural network, getting a classification [2].	18
3.6	Typical CNN architecture [3].	20
3.7	Example of convolution operation and pooling operation (max pooling). 21	
4.1	The μ SmartScope prototype, with smartphone attached and microscopic slide inserted [4].	26
4.2	The μ SmartScope application pipeline. The ML4CervicalAdequacy module correspond to the module developed in the scope of this thesis. 27	
4.3	Microscope slide with LBC preparation (blue) with area of sampling (orange). Images of the sample are taken along the orange area with 40 \times amplification.	28
4.4	Examples of raw images that compose the dataset with one image per smartphone.	29
4.5	Visual comparison between two samples. Sample 4.5a shows an inadequate sample with a scarce amount of nuclei. Only two or three can be identified in all three images. On the other hand, the adequate sample 4.5b is almost fully covered with squamous cells.	31
4.6	Examples of cropped objects of each class in relative scale to each other. 32	

4.7	Examples of images with a dense arrangement of glandular nuclei.	33
4.8	Example of the distribution of two predominant types of nuclei in samples. There are images where mostly only one type of nuclei is present, such as Figure 4.8a and Figure 4.8b. Figure 4.8c shows a case where the distribution is even.	34
4.9	Examples of images with a dense arrangement of nuclei. Some nuclei appear without annotations.	35
4.10	Example of cell overlapping problem, 4.10a, and comparison to a clear image, 4.10b.	36
4.11	Example of chromatic aberration. The image was cropped to its actual size to facilitate observation.	37
5.1	The ML4CervicalAdequacy module pipeline.	40
5.2	Model-centric and data-centric development pipelines.	41
5.3	Visual representation of dataset split into train and test sets. Also, cross-validation split into three folds can be visualized.	44
5.4	Percentage of annotations in each set, train and test, for each class.	44
5.5	Percentage of annotations of the fold 1 in each set: train and validation; for each class.	45
5.6	Example of cut cells.	45
5.7	The green boxes indicate the annotations that were kept on the patch. The red boxes indicate that these patches were eliminated from the patch. It is possible to verify that most of the eliminated annotations are not containing any meaningful information and, therefore, should not be kept.	46
5.8	Cropping image to the size of the optical disk: (a) Original image; (b) Segmentation mask of optical circle with cropping; (c) Original image cropped to the size of segmentation mask.	46
5.9	Semi-supervised pipeline to improve model performance by increasing training data.	48
5.10	Tensorboard interface.	50
5.11	IoU visual explanation [5].	51
6.1	Examples of predictions in form of bounding boxes with respective classification and confidence (on the left side of each image) compared to the ground truth annotations (on the right side of each image).	58
6.2	Overall mAP (on vertical axis) progression of two models: MobileNet (Orange) and ResNet50 (Blue); during the training on fold 1.	60

6.3	Results of hyperparameters tuning evaluated on the test set. Image (a) outlines relative performance increase of bigger batch size values when compared to batch size of 4. Image (b) depicts mAP of squamous nucleus class across different values of LR with batch size of 16. Image (c) shows mAP across different classes with LR = 5.29×10^{-4} and batch size = 16.	62
6.4	Example of chromatic aberration correction. Upper part of 6.4a and 6.4b is the original image and bottom part is the same image after correction. The images were cropped to its actual size to facilitate observation. As can be noticed, correction did an excellent job in eliminating blue and orange fringing around nuclei and other objects.	63
6.5	Confusion Matrix of the model after semi-supervised training.	66
6.6	Detections and classifications made by the final system (a) and (b). Image (a) depicts mostly correct detection and classifications (green and blue bounding boxes means correct classifications of two different classes), while on the image (b) yellow and red bounding boxes correspond to misclassification and false positives detections respectively. Image (c) shows the confusion matrix of the test set. "Nothing" corresponds to false positives and false negatives.	68
6.7	Results for the proposed solution: (a) average number of detected squamous nuclei per sample (sample 1 to 5 are inadequate and 6 to 12 adequate); (b) and (c) illustrative images of automated nuclei detection at the image level.	69
B.1	Artifact	85
B.2	Glandular nucleus	86
B.3	Inflammatory cell	86
B.4	Squamous nucleus	86
B.5	Undefined all	87
B.6	Undefined nucleus	87

List of Tables

4.1	Dataset distribution by clusters/abnormality levels.	30
4.2	Dataset distribution by nucleus type.	31
5.1	Models considered for training.	42
5.2	Learning Rate values considered for the tests.	47
6.1	Virtual Desktops used for the development	56
6.2	VD Host Server	56
6.3	Dataset distribution by nucleus type.	65
6.4	Performance metrics of the model after semi-supervised training, taken on the final system.	66
6.5	Performance metrics of the final system.	67
A.1	Dataset division: Stratification of nuclei annotations	84
C.1	Results for model MobileNet using cross-validation folds.	90
C.2	Results for model ResNet50 and EfficientDet D1 using cross-validation folds.	91
C.3	Results for model ResNet50 and MobileNet for best parameters found during hyperparameter tuning. Results obtained using TRAIN and TEST sets.	92
D.1	Results for model ResNet50 with some classes merged and eliminated.	94
E.1	Part 1 of the filtered table.	97
E.2	Part 2 of the filtered table.	98
E.3	Part 3 of the filtered table.	99
E.4	Part 4 of the filtered table.	100
E.5	Part 1 of the complete table.	102
E.6	Part 2 of the complete table.	103
E.7	Part 3 of the complete table.	104
E.8	Part 4 of the complete table.	105

List of Acronyms

3D	Three Dimensional
AP	Average Precision
API	Application Programming Interface
AR	Average Recall
ASC-H	Atypical Squamous Cells
ASC-US	Atypical Squamous Cells of Undetermined Significance
CADx	Computer-aided Diagnosis
CNN	Convolutional Neural Network
CP	Conventional Preparation
CPS	Conventional Papanicolaou Smear
CPU	Central Processing Unit
CV	Computer Vision
DNA	Deoxyribonucleic Acid
DNN	Deep Neural Networks
DSS	Decision Support System
FN	False Negative
FOV	Field of View
FP	False Positive
FPN	Feature Pyramid Network
FhP	Fraunhofer Portugal
GPU	Graphics Processing Unit
HFF	Hospital Fernando Fonseca
HPV	Human Papillomavirus
HSIL	High Grade Squamous Intraepithelial Lesion
IPO	Instituto Português de Oncologia do Porto
IoU	Intersection Over Union
LBC	Liquid-Based Cytology
LBP	Liquid-Based Preparation
LR	Learning Rate
LSIL	Low-grade Squamous Intraepithelial Lesion
ML	Machine Learning
NMS	Non-maximum Suppression

NN	Neural Network
PAA	Pattern Analysis and Applications
PCR	Polymerase Chain Reaction
RAM	Random Access Memory
RNA	Ribonucleic Acid
SCC	Squamous Cell Carcinoma
SSD	Single Shot Detector
TBS	The Bethesda System
TF	TensorFlow
TNR	True Negative Rate
TP	True Positive
TPR	True Positive Rate
TZ	Transformation Zone
USB	Universal Serial Bus
VD	Virtual Desktop
VDI	Virtual Desktop Infrastructure
VIA	Visual Inspection with Acetic Acid
VOC	Visual Object Classes
VRAM	Video Random Access Memory
WHO	World Health Organization
mAP	mean Average Precision
vSAN	Virtual Storage Area Network

Chapter 1

Introduction

1.1 Context and Motivation

Cervical cancer is the fourth most frequently diagnosed cancer and the fourth leading cause of cancer death in women. Nearly 605 000 new cases were registered in 2020 and about 342 000 died from cervical cancer worldwide [6][7][8]. Most cases of cervical cancer are caused by infection with human papillomavirus (HPV), which is preventable with a vaccine. However, in the low-to-middle-income countries, where a lack of awareness and limited access to health services are the main factors responsible for nearly 84% of new cases and 87% of deaths, this practice fails to be successful[9][10][11][12].

Nevertheless, screening tests such as cytology have been responsible for a strong decrease in cervical cancer deaths over the past years. Screening cytology reduced the incidence of cervical cancer by 60% to 90%, and the death rate by 90% [13]. The test consists of collecting cells from the patient's uterine cervix and submitting to expert cytologist for analyzing the changes under the microscope of HPV effects [14]. However manual screening can be difficult, tedious, time-consuming, expensive and subjected to error induced by inaccurate sample preparation. Several additional factors can worsen a visual inspection by obscuring factors like excessive presence of blood, inflammations, bacteria, lubricant or by the simple fact of not having a minimum amount of cells.

The increasing interest in the development of computer-aided diagnosis (CADx)

systems for cervical screening is related to difficulties experienced in these under-resourced health facilities, such as the shortage of specialized staff and equipment. CADx systems often use machine learning and deep learning techniques to reduce the dependence on manual input and respectively increase their autonomy. This technique does not require a highly experienced cytotechnologist and can be deployed on relatively cheap hardware and, at the same time, be faster than human in inspecting the sample. During the last years, various computational methods have been proposed and implemented to support cervical cancer screening [15][16][17].

In a recent review article [18], the authors analysed and discussed focus and adequacy assessment, segmentation and classification computational approaches used for the analysis of microscopic images from cervical cytology smears. Among the general conclusion about the technical state of computer vision methods, the authors outlined smear adequacy assessment as a topic scarcely addressed in the literature. Most of the works simply ignore it, while others implement some techniques to detect and remove unwanted objects such as inflammatory cells, dirt, blood or other artifacts.

1.2 Research Goals

This work will mainly address automated approaches for the assessment of cytological specimen adequacy. The major focus will be given to cellularity assessment, being low squamous cellularity the most common cause for the unsatisfactory specimens, as well as following The Bethesda System (TBS) guideline. The starting point will be exploring a private image dataset collected by Fraunhofer Portugal (FhP) in the ambit of the CLARE project. The CLARE dataset is being collected with the μ SmartScope prototype at IPO-Porto and Hospital Fernando Fonseca, and contains annotations of cellular structures highly relevant for our work. The initial work will be focused in the implementation and testing of the detection and the classification of cervical cells (major focus on nuclei detection). Having detected and classified relevant objects, we will work on their counting and extraction of metrics in order to measure objectively the performance of our system, and, finally, we will explore data-centric and model-centric approaches to improve the model's performance.

Considering that there are not many mobile-based and lightweight systems to assess cervical specimen adequacy available, during this work we will try to give answers to the following questions:

- Can lightweight Deep Learning object detection approaches be used for the detection of cervical cells, giving major focus on nuclei detection?

- Can lightweight Deep Learning approaches be used for automated classification and cell counting of different types of cervical cells, and are the respective performances suitable for sample adequacy assessment purposes?
- Can the developed object detection models be improved via optimization strategies such as semi-supervised or data-centric approaches?

1.3 Research Contributions

During writing of this dissertation, we were able to contribute to the Fraunhofer AICOS and also to the scientific society with the work we developed. The major contributions are:

- Development of lightweight object detection deep learning models for detecting and classifying cellular structures in microscopic images of LBC samples relevant for sample adequacy assessment.
- Explored and evaluated the impact of applying optimization strategies based on semi-supervised and data-centric approaches to improve the performance of the developed object detection models.
- Leverage the suitability of the developed models for usage scenarios with constraint computational and memory resources, through the use and optimization of pre- and post-processing image slicing techniques.
- Contribution to the scientific society by presenting a scientific paper at IbPRIA 2022: 10th Iberian Conference on Pattern Recognition and Image Analysis. The paper was already presented and published:

Mosiichuk, V., Viana, P., Oliveira, T., Rosado, L. (2022). Automated Adequacy Assessment of Cervical Cytology Samples Using Deep Learning. In: Pinho, A.J., Georgieva, P., Teixeira, L.F., Sánchez, J.A. (eds) Pattern Recognition and Image Analysis. IbPRIA 2022. Lecture Notes in Computer Science, vol 13256. Springer, Cham. https://doi.org/10.1007/978-3-031-04881-4_13

1.4 Document Structure

This document is structured as follows: In Chapter 1 an introduction, along with the quick contextualization and motivation are made. Moreover, research goals and contributions are present as well. Chapter 2 introduces cervical cancer and cervical cancer screening, and makes an overview of TBS guideline and specimen adequacy criteria. Chapter 3 introduces computer vision and machine learning techniques. Existing computational methods related to detection and classification of objects in

cytological images are overviewed. Chapter 4 contextualizes previous and current research work carried out at FhP that are relevant for this thesis, and makes an overview of the dataset used in the development of the deep learning model. In Chapter 5 the proposed methodology is discussed. In Chapter 6 we will evaluate and discuss the developed work. Finally, Chapter 7 summarizes the results obtained and gives a critical appreciation of this dissertation.

Chapter 2

Cervical Cancer

The main purpose of this chapter is to make a more in-depth introduction to cervical cancer, cervical screening types and TBS guidelines. In the first subchapter more about cervical cancer will be presented, followed by an overview of various screening types, their implications, advantages and some disadvantages present in the second subchapter. In the final subchapter, a quick review of a few TBS guidelines will be done. Moreover, it is important to note that the TBS guidelines play a key factor in further assessment of samples quality since they specify the threshold between satisfactory and unsatisfactory results.

2.1 Cervical Cancer

Cervical cancer is a disease in which cancer, or malignant cells, form in the tissues of the cervix. Cervical cancer usually develops slowly over time. It can take 15 to 20 years for cervical cancer to develop in women with normal immune systems and only 5 to 10 years in women with weakened immune systems. Before cancer appears in the cervix, the cervical cells undergo changes known as dysplasia, in which abnormal cells begin to appear in the cervical tissue. Over time, the abnormal cells may become cancer cells and start to grow and spread more deeply into the cervix and surrounding areas [19].

Long-lasting infection with certain types of HPV is the main cause of cervical cancer [19]. HPV is a group of viruses that are extremely common globally. There are more than 100 types of HPV, 14 of which are known to be cancer-causing or

high-risk type. But only two HPV types are responsible for 70% of cervical cancers and pre-cancerous cervical lesions. Other types of HPV infection, in most cases, clear up without any intervention within a few months after the acquisition, and about 90% clear within 2 years [20].

Generally, the peak time for acquiring infection is shortly after becoming sexually active. HPV is sexually transmitted, but penetrative sex is not required for transmission, being the skin-to-skin genital contact a well-recognized mode of transmission, as well [20].

The World Health Organization (WHO) perceives three categories of epithelial tumours of the cervix: squamous, glandular and other epithelial tumours.

Squamous carcinomas are composed of cells that are, as its name suggests, distinctly squamous. They are skin-like cells that cover the outer surface of the cervix and vary in either growth pattern or cytological morphology [21][22].

Squamous cell carcinomas are accounted for 70% to 80% of cervical cancers, being the HPV-16 the major causal agent [13].

Adenocarcinomas are characterized by neoplasia of epithelial tissue that has a glandular origin, glandular characteristics, or both. Usually, adenocarcinoma cells start developing along the inside of the passage that runs from the cervix to the womb, designated by the endocervical canal, where glandular are scattered. Even though squamous carcinomas and adenocarcinomas generally are associated with being caused by HPV, adenocarcinomas are particularly linked to HPV-18. Compared to squamous carcinomas, adenocarcinomas are much less common, being accounted for 20% of cervical cancers.

Adenosquamous carcinomas are characterized by being composed of malignant squamous and malignant glandular cells and represent around 5% to 6% of women diagnosed with cervical cancer [13][22]. Adenosquamous and other rare cancers such as small cell cancer, lymphomas and sarcomas fall into other types category since they represent a small portion of the global number of cases and may require a different approach for screening and treatment when compared to squamous carcinomas and adenocarcinomas [22].

2.2 Cervical Cancer Screening

When screening tests are conducted, we might be focusing on both pre-cancerous and cancerous lesions. Thus, screening tests can be performed in women without any reported symptoms. With effective screening strategies, pre-cancerous lesions can be detected in the early stages and treated before progressing to heavier stages, which is proven to prevent up to 80% of cervical cancers [20].

Currently there are three types of screening tests that are recommended in the latest guidelines from World Health Organization (WHO): i) HPV testing for high-risk HPV types; ii) Visual Inspection with acetic Acid (VIA); iii) and cervical cytology: conventional (Pap) test and Liquid-Based Cytology (LBC)[20].

HPV Test: is a laboratory test used to check DNA or RNA for certain types of HPV infection. Cells are collected from the cervix and checked using polymerase chain reaction (PCR) or hybrid capture, to find out if there is any cervical-cancer related HPV infection. HPV test can be performed sole or as a consequence of finding abnormal cervical cells in a Pap test. For women age 30 or older, Pap and HPV cotesting or HPV testing alone are more sensitive than Pap testing alone [23]. Particularly, testing for HPV has demonstrated greater sensitivity for high-grade cervical intraepithelial neoplasia and provides 60% to 70% greater protection against invasive cervical cancer [13][18].

VIA: or Visual Inspection with Acetic Acid, is a method that enables physicians to directly see cervical lesions. It requires the use of a speculum and light source, and must be performed by a trained health-care provider. A 3%-5% acetic acid solution is liberally applied to the cervix with a large cotton swab. After removing the cotton swab, the provider waits for 1-2 minutes and visually inspects for white areas. Usually, in the first minute, areas that become faintly white simply due to inflammation physiological cell changes, such as metaplasia, will recede, while damaged tissues will continue white and more likely will be associated with cervical pre-cancer or cancer. Generally, cervical lesions appear near the squamocolumnar junction. Once white areas are detected, the physician can remove the damaged tissues using cryotherapy or other techniques. Despite the fact that VIA yields high sensitivity, it's not recommended for postmenopausal women due to the fact that it is more difficult visualize the squamocolumnar junction in older women, which affects negatively the result of the test [18][19]. VIA is quite inexpensive, utilizes locally sourced supplies and does not rely on specialized laboratory services. However, due to the subjective nature of the test, quality control and quality assurance for VIA is particularly important [19].

Cervical Cytology: has been, for many years, the standard method for cervical cancer screening, responsible in many cases for drastic reduction of mortality rate. In many countries of the world, the Papanicolaou (Pap) test helped to reduce the incidence rate by 60% to 90% and decreased death rate by 90% [13]. At present, there are two types of the test that are performed: Conventional Papanicolaou Smear (CPS) and Liquid-Based Cytology (LBC). While LBC is more common in developed countries, in low resource settings, CPS is the mainstay screening system [24]. In the last years many researches were done in regards to evaluating and comparing efficacy of LBC and CPS as the screening tool for cervical cancer. Results have demonstrated that LBC test provides higher quality samples, as it offers

an uniform spread of cells and better clarity, which affects positively the efficacy of the test by reducing the number of unsatisfactory samples (cases) and consuming less time [25][26][27][28]. Particularly, LBC offers better results in detecting cervical pathology than CPS in cases of postmenopausal women, where menopause-induced anatomical and hormonal changes make screening a rather challenging procedure [26]. However, even LBC being considered a higher quality test than CPS, some works concluded that, aside from the case mentioned earlier, it doesn't offer a significant difference in diagnostic performance, thus making CPS more viable in contexts where LBC is not affordable [29][30]. Cervical cytology has undeniably provided substantial contribution in cervical screening and continues to do so, and even though it has some inherited limitations such as being a morphological method requiring subjective interpretation by well-trained cytologists. Despite continuous efforts to improve the performance of cervical cytology, its sensitivity isn't optimal yet and the method yields high number of borderline results, which leads to an excessive number of referrals for colposcopy and overtreatment [13][31][32].

2.3 The Bethesda System (TBS) guidelines

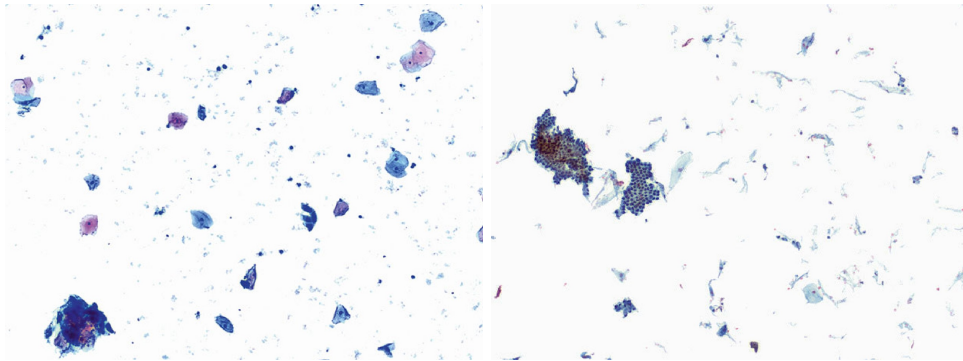
There are some classification systems for classifying cancerous or pre-cancerous lesions of the cervix. They may be based either on cytology or histology and can have different clinical and reporting purposes. In case of screening tests, the most applicable classification systems is The Bethesda System. TBS provides a standard set of terms to describe abnormal findings from CPS and LBC which reports sample equality and cytology results [33]. Since the main purpose of our work is to explore and implement automated approaches for assessment of cytological specimen adequacy, the TBS guidelines will be particularly crucial to define classification parameters of our implementation to separate satisfactory samples from unsatisfactory ones.

Several factors can worsen sample quality and, therefore, hinder the interpretation of certain specimens. These factors usually are: technical problems in the slide preparation, low cellularity or obscuring factors like excessive presence of blood, bacteria, lubricant or inflammations. False negatives (FN) are often reported during manual screening due to these reasons, thus, specimen quality assessment and its reasoning are an important step in decreasing the number of FN. In order to achieve better sample quality, by decreasing the subjectivness and uniformizing the analysis, TBS defined a set of guidelines that are recommended for laboratories to follow [1]:

- **Cellularity:** In 2001 The Bethesda System stated that an adequate LBC preparation should have an estimated minimum of at least 5000 well-preserved squamous cells. It is important to note that endocervical cells and completely obscured cells should be excluded from the estimate. When referring to CPS, a specimen should contain an estimated minimum of approximately 8,000–12,000

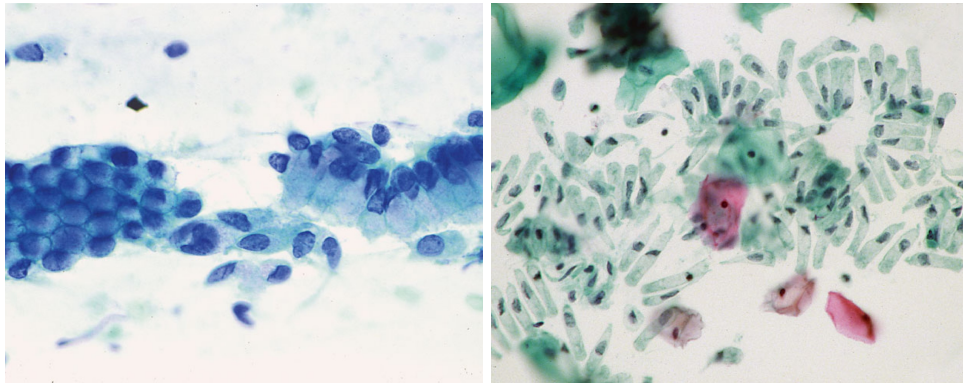
well-preserved and well-visualized squamous epithelial cells. These are the minimum values for a specimen. The difference between threshold for LBC and CPS is mainly due to the nature of the LBC test, which presents a more random, and presumably more representative, sample of the collected cervical material. Also, it should be noted that thresholds for LBC and CPS specimens, should not be rigidly applied in vaginal and post-therapy specimens [1]. Some examples of cellularity assessment are given in Figure 2.1.

- **Transformation Zone:** Although transformation zone sampling is not necessary for an adequate specimen, only squamous cellularity is necessary. Presence or absence should be present in reports as it may be a useful quality assurance measure. In the case of accounting for Transformation Zone (TZ), for CPS and LBC, an adequate transformation zone sample requires at least 10 well-preserved endocervical cells or squamous metaplastic cells, in a cluster or singly, as can be seen in a Figure 2.2. One group is often enough to indicate that the TZ has been sampled. In the past, there were concerns that the samples of the squamocolumnar junction were not adequate when the cytology specimens lacked TZ component. However, some studies showed that there was almost no difference in detection of high grade squamous intraepithelial lesion (HSIL) between patients who were TZ negative in the initial screening evaluation and those who were TZ positive [1].
- **Obscuring factors:** are another component which can lead to an unsatisfactory specimen classification. Specimens in which more than 75% of squamous cells are obscured, and no abnormal cells are identified, should be designated as unsatisfactory. When the percentage is between 50% and 75%, the report should describe a specimen as partially obscured, but satisfactory result. It is important that cells' nuclei preserve and could be well visualized since partial obscuring of cytoplasmic detail may interfere with specimen evaluation. Also, it should be noted that the percentage of obscured cells should be evaluated and not the slide area obscured, even though a minimal cellularity criteria should also be applied. These criteria are similar to CPS and LBC samples [1].



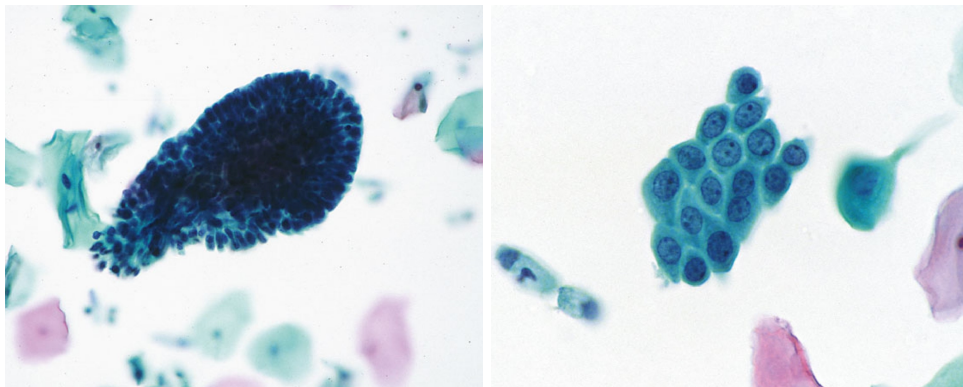
(a) Satisfactory, but borderline level for (b) Unsatisfactory specimen due to scant specimen. Based on four cell average per field squamous cellularity. Should be noted that at 40x slightly over 5000 cells were estimated. endocervical cells, seen in a honeycomb arrangement, are not taken into account

Figure 2.1: Satisfactory 2.1a and Unsatisfactory 2.1b LBC ThinPrep preparations at 10x magnification [1].



(a) Endocervical cells (Conventional preparation (CP)). Distinct cytoplasmic borders are seen in the cluster of cells on the left, and the cell cluster on the right is seen from a side view, giving the “picket fence” appearance.

(b) Endocervical cells (Liquid-based preparation (LBP), SurePath). Cellular dissociation is more frequent in liquid-based preparations than in conventional smear preparations.



(c) Endocervical cells (LBP, SurePath). Normal endocervical cells may appear in large hyperchromatic fragments, often in the center of some LBP. The thickness of the fragment may give the appearance of architectural disarray; however, note normal appearing cells at the periphery of the fragment. Additionally, focusing up and down through the fragment reveals normal spacing of cells, distinct cytoplasmic borders, and bland nuclear chromatin. Normal endocervical cell groups with this appearance should not be confused with neoplastic clusters that show more crowding (even within a single layer of cells), nuclear enlargement, nuclear membrane irregularity, and abnormal chromatin pattern.

(d) Transformation zone component (LBP, SurePath). Normal endocervical cells from the upper region of the endocervical canal can closely mimic squamous metaplastic cells.

Figure 2.2: Examples of transformation zones in Conventional preparation (CP) and Liquid-based preparation samples [1].

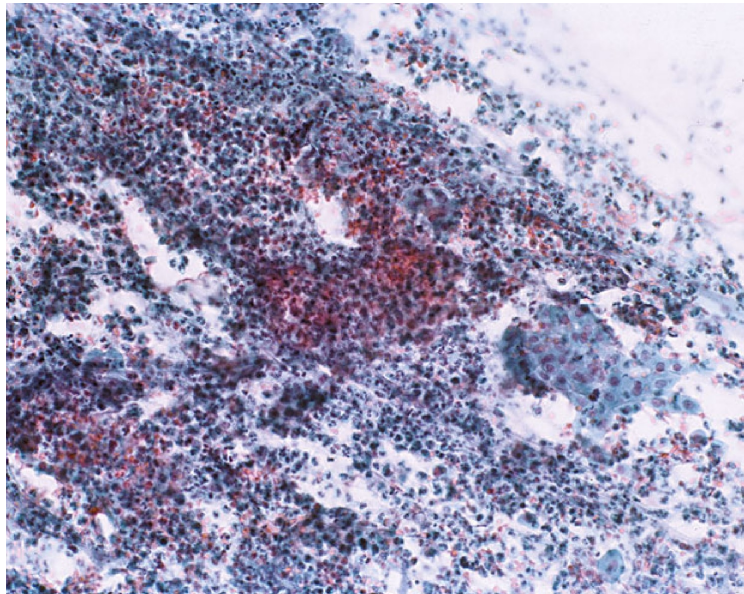


Figure 2.3: Unsatisfactory specimen due to obscuring white blood cells. It should be noted that when assessing the adequacy of the slide considering obscuring factors and cellularity, the cellularity refer only to well visualized cells [1].

Chapter 3

Computer Vision and Deep Learning: Background concepts and State-of-the-art Review

This Chapter introduces the concepts of Computer Vision and Deep Learning. Techniques such as Machine Learning and Deep Learning are also introduced alongside Convolutional Neural Networks, Neural Networks and how they are trained. And finally, the overview of the state-of-the-art techniques for object detection and classification in cytological images is made.

Existing computational methods related to detection and classification of objects in cytological images are overviewed

3.1 Computer Vision

Computer Vision (CV) is a field of computer science that aims to provide computers with capabilities to process and make a human-like interpretation of visual scenes.

As humans, when we are looking at some scene, we perform a lot of implicit calculations in our heads, being left at the end with the visual representation of what we see and the information about it. For example, we can roughly estimate how far away the objects are from us, what the object is, its position within the scene, and other information we may or may not fill, considering our memory and knowledge

we already have acquired through the course of our lives. Unfortunately, we can not just teach computers to interpret images/videos as we do because computers are just very complicated calculators, and they have too little in common with us. For instance, an image for the computer is just a matrix with values that represents the intensity of given color in a given spot. Whether for humans, it's a set of photons bouncing from the surfaces and coming through the eye to our brain, which ends up processing it.

On a lower and more technical level, we can decompose the problem of image interpretation to image analysis and pattern recognition. The goal is to use some algorithm that, given an image, can retrieve meaningful information about what is in the image, such as the object present in it, forms, and other properties. While speaking of image analysis, it comprises tasks such as pre-processing, detection, and segmentation. Pre-processing can be done on the level of the whole dataset structure or the image level. When referring to the entire dataset structure, the pre-processing task may include adding or eliminating individual images, clustering them by their classes/groups, and dividing in training and test sets. On the other hand, the pre-processing on the image level focuses on modifying each image's characteristics such as brightness, contrast, and geometry through techniques like: histogram stretching and equalization, color correction, geometry transformations or some noise filtering. Usually, the goal of these operations is to create a uniform dataset and help distinguish essential structures from non-essential and noise. In other words, pre-processing aims to help the algorithm achieve better results.

The detection's task is to find structures of interest and localize them within the limits of the image. Once objects are localized, this information may be used for purposes such as counting the number of objects in the scene or localizing them, and extracting relevant information about positions in order to each other or in reference to some key point.

Pattern recognition is strictly related to detection. While simple detection tasks may rely on detection or recognition of elementary geometry figures such as circles, squares, or lines, which indeed can be defined as patterns, most of the real-world problems require recognition of more complex patterns, for example, shape of cats, cars, or people. As stated in [34]: "The field of pattern recognition is concerned with the automatic discovery of regularities in data through the use of computer algorithms and with the use of these regularities to take actions such as classifying the data into different categories." On the technical level and in more cases where patterns aren't simple objects, pattern recognition uses feature extraction to extract information about patterns and then uses matching techniques to determine the presence and localization of the desired pattern in the image. In the past years, more sophisticated Machine Learning (ML) algorithms, such as Deep Neural Networks (DNN), started to be used to automate the process of extraction and matching, as

well as drawing some conclusions about found objects.

3.2 Machine Learning

Machine Learning is mainly characterized as a field of computer science that enables computers to learn without being explicitly programmed. In contrast to the conventional programming where the function is needed to be explicitly defined, ML allows to approximate the desired function by giving it input and output data, and with means of some sort of cost function reduce the distance to the correct output. In order to "program" ML algorithm, the key insights from the raw data are needed to be obtained. ML leverages this process by determining the insights using models that encode structural descriptions of the analysed data, which later can be used to make predictions. The structural description is obtained through the process of optimisation, or by another words, through the training of the algorithm. To facilitate this process, usually a feature space is defined and relevant features are extracted, instead of feeding the raw data.

There are three main types of ML: Supervised Learning encompasses problems where available data for training the model includes output ground truth values, whereas Unsupervised Learning performs the training without consistent true output values; The third type is Reinforcement Learning (RL) which is characterized by not needing nor input nor output labels, it tries to approximate the desired function by allowing it performing any type of action, but also rewarding it for the correct action, or for one that drives the algorithm in the right direction.

Machine Learning algorithms can also be divided in two categories depending on its output: Discrete and Continuous. Discrete outputs usually can be classifications or clustering where a set of possible answers known and unknown respectively, and it is limited. For the Continuous we can have Regression or Reduction of Dimensionality [35].

3.2.1 Classification Task, Segmentation and Object Detection

In the context of ML, classification can have two different meanings: (1) establishing the existence of classes or clusters in a given set of observation or (2) establishing a set of rules where one can classify a new observation into one of pre-established existing possible classes. The first definition refers to the Unsupervised Learning, whilst the second definition, which can be considered the standard definition of classification, refers to Supervised Learning. In the context of CV, classification has a similar definition to the one given in ML, as it is defined as the process of categorizing a vector of stimuli or features into a finite set of classes, meaning that classification in CV usually involves recognizing the dominant content in a scene. It is important to note that in CV, detection differs from classification, since

the former requires knowing the location of target objects in a scene, essentially making detection a process composed of both classification and localization. The term segmentation (or image segmentation) consists of the partitioning of an image into a set of regions that cover it, with the goal of obtaining regions that represent meaningful areas of the image, such as tumours or polyps in a medical image, or crops, cars, urban areas and forests in satellite images. These regions can either be a border of pixels surrounding the area of interest or a shape, such as a circle, ellipse, or polygon. The goal of segmentation is simultaneously to decompose the image into parts for further analysis and to perform a change of representation of the regions into more meaningful concepts and efficient units for further processing and analysis. Image segmentation encompasses a wide range of fields, such as object recognition and detection, which consist in the detection and recognition of one or more objects in a specific scene, depending on the type of problem. Semantic image segmentation is regarded as the more fine-detailed case of image segmentation, where each pixel is classified into belonging to a specific class, essentially clustering the pixels in accordance to their classes. Although object detection is encompassed by image segmentation, it differs from the former in that object detection has the goal of detecting the bounding shape (usually a box) of the target object, unlike segmentation, which includes the region shape of the detected object.

3.2.2 Neural Networks

Inspired by the human brain, see Figure 3.1, "artificial" Neural Networks (NN) translated to the computer feature neurons, activation and interconnectivity [2].

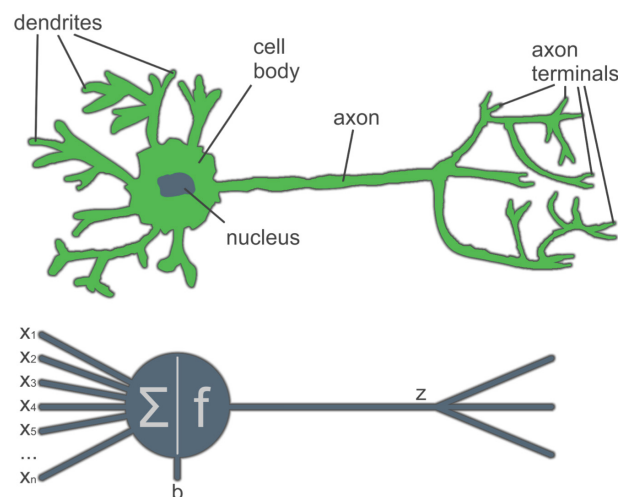


Figure 3.1: Comparing a biological neuron to an artificial neuron [2].

Similarly to brain neurons, "artificial" neurons, also called perceptrons, interconnect with each other forming dense layer. In a dense layer, each neuron of a

given layer is connected to every neuron of the next layer, therefore, its output becomes input of the next neurons. Each connection between neurons has a weight associated with it, which gets modified through the training process and essentially defines how much the connection is relevant to the network. In practice, the value that passes through the neuron gets multiplied by the weight. Once all inputs and weights of the neuron get multiplied, they are summed and a bias value, which is another trainable parameter, is added. Then, the application of activation function generate the perceptron output, or lack of it, depending on the function used. As a very general overview, the step function meant to mimic a neuron in the brain, either "firing" or not. In our case, the function that depicts this behaviour is called Heaviside function, or a step function, as can be seen in the Figure 3.2

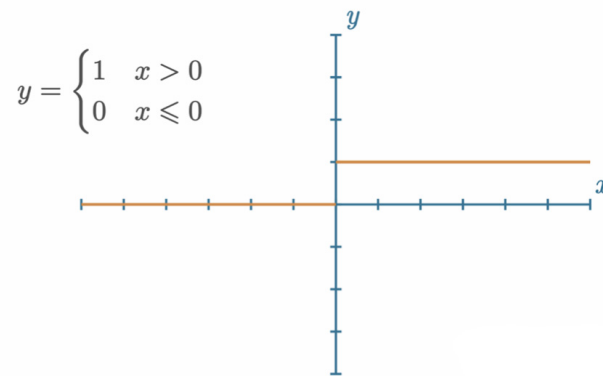


Figure 3.2: Graph of a step function [2].

While a step function could still be used in models, there are better performing functions such as already mentioned ReLU or Sigmoid functions, depicted in Figure 3.3

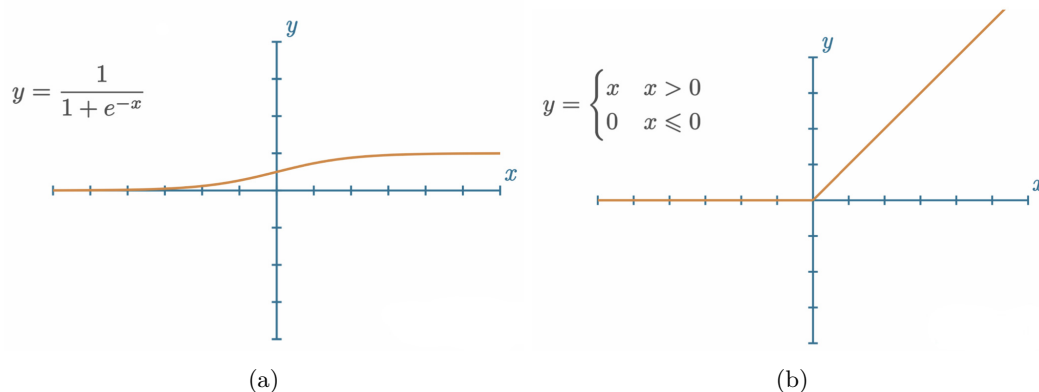


Figure 3.3: Sigmoid function graph 3.3a and Graph of the ReLU activation function 3.3b [2].

Figure 3.4 shows an example of a basic network composed of one input layer, two dense hidden layers, and lastly, one output layer.

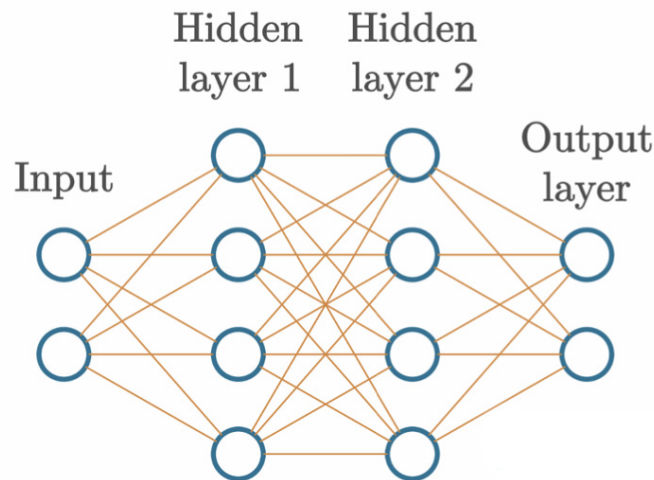


Figure 3.4: Example basic fully-connected neural network [2].

While input layers already could take raw data, some pre-processing is actually needed, it is recommended that the data is normalized and equally scaled. Usually normalization consists of scaling down or up values to the range from 0 to 1 or from -1 to 1, while maintaining all the features. The output layer would be the actual output of what the network will return. For the classification task, the number of output perceptrons will actually match the number of classes, but for binary problem, we can have only one perceptron. For instance, for a problem where we try to classify a "dog" and a "cat", we will have two output perceptrons, but we can also have only one, where output can be interpreted as "dog" and "not dog". On Figure 3.5 an input image can be seen to the left of the neuronal network. On the NN itself, dots or perceptrons can be observed with different color intensity, where the darker the color, the stronger the connection, and therefore, that was a path for the features to get and "dog" output from the network.

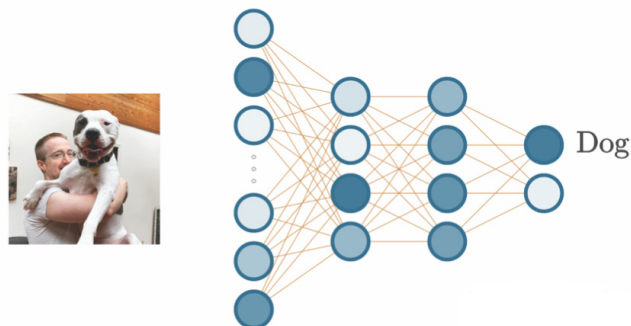


Figure 3.5: Visual depiction of passing image data through a neural network, getting a classification [2].

3.2.3 Training Neural Network

In order to train the network, the connection weights and bias are initially initialised with , usually, random values and then are updated through an iterative process where the network is exposed to new data and tries to make predictions. For the single data instance, the network computes output values and compares them with the ground truth using a loss function. A loss function essentially quantifies the error between prediction and the real value. The main goal in the training process is to reduce the error. In each iteration the networks weights are updated accordingly to the output of the loss function, in order to minimize it. However, due to the, usually, large amount of weights and biases that are to be updated in the NN, the relation between the loss function's first derivative and the weights of each layer have to be established. To achieve this, the backpropagation algorithm propagates the cost function's gradient in a reverse direction, from the output layers to the input, taking into account partial derivatives at each perceptron and the derivation chain rule [36].

It is important to note that a single perceptron can also be seen as a classifier. However, it would not be capable of tackling complex non linear problems. For this we have to use multi layer networks, which mainly are similar to the network present in the Figure 3.4, but usually offer even more layers.

3.2.4 Deep Learning

Deep Learning is a sub-field of Machine Learning that extends it by allowing automatic extraction of features through usage of neural networks [36]. The neural network is fundamentally formed of layers, which process raw data and extract features from it. The number of these layers represent the depth of the network and if the number of hidden layers is greater than two, then we have a deep neural network. In each layer linear and non-linear operations are performed, allowing extraction of more complex and abstract features, as more layers are added and network become deeper. In comparison to the conventional machine learning approaches, the abstract representations learned by DL networks allows considerable leap forward in performance. Which mainly can be addressed to the remarkable ability of bypassing the manual feature extraction step [37].

The idea of using human brain alike neurons to perform machine learning is not new, but at the time of proposing it, the computers were not capable of executing algorithm quickly enough to be able to show some significant performance [38]. Nevertheless, a revolution in the computer vision field and deep learning took place in 2012 at ILSVRC competition [39]. Alex Krizhevsky and his associates used NN-based, particularly a Convolutional Neural Network (CNN), approach and obtained 16.4% top-5 error, which was a major improvement from the previous year.

The outstanding performance of the network created, AlexNet, was achieved with a combination of several factors: availability of a vast amount of labeled data, new data augmentation strategies, development of dropout regularization, Rectified Linear Unit (ReLU) for the activation function and the rise of fast graphical processing units (GPU). The victory of AlexNet motivated development of several CNN models, as well as the incentive to try those models in fields such as CV.

3.2.5 Convolutional Neural Networks

Despite the fact that Convolutional Neural Networks was not a new technology back in 2012, as stated above, the annotated data availability and fast hardware paved a path for CNN to be the state-of-the-art network when it comes to the CV applications. A traditional CNN, as can be seen in the Figure 3.6, consists of single or multiple convolution and pooling layers, followed by one or more fully connected layers and an output layer. The convolution layer is the core of an CNN. This layer allows the network to learn feature representation of the input. It is composed of several learnable convolution kernels, which are in charge of computing feature maps. Each new feature map is produced by convolving the input with the kernels and applying elementwise non-linear activation function on the convolved result. Then, each unit of feature map is connected to a receptive field of the next layer.

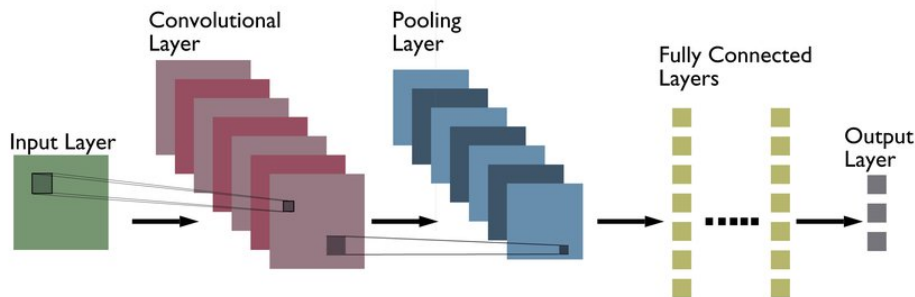


Figure 3.6: Typical CNN architecture [3].

The kernel defines not only the receptive field, but also the coefficient used to multiply the data, or in another words, the weights. Along with the additional bias values, the weights are the values that are optimised during the training process in the way that, when training is completed, they represent the learned features. For the neurons at the same depth level weights are shared as a result of the need of reducing the number of parameters to optimize, and thus, simplify the training process. This leads to the number of weights needed to be small, which also turn out to be memory efficient.

Convolution kernels can be seen as filters that only allow the transmission of certain features. So, if the input has this features, it passes, if it does not, then the input is blocked. The patterns that are learned by the network can represent low-level

features such as edges, or also a more complex, high-level features which represent certain structures in the data. Usually the high-level features are more abstract and result from the combination of several low-level features. It is important to state that the depth of the network, then, greatly influences the capability of the network to learn certain features of the data and the performance of the network. Besides weights optimisation, the number of layers, amount of filters in each layer, their dimension, stride of the kernel and the introducing of zero-padding in the output can also be seen as parameters to optimise in the CNN.

The following layers, which are pooling layers, are used to sub-sample the feature maps after the convolution operations. Essentially it shrinks the feature maps to a more simpler, while trying to preserve the most dominant feature. The pooling operation is carried out by specifying the pooled region size and the stride, similarly to the convolution operation. The most common pooling type is the max pooling where the maximum value inside the area defined by the kernel is considered as the output value. There are also other pooling types such as average pooling, where an average of the values inside the kernel area becomes the output, or for example min pooling, gated pooling, tree polling, etc. Figure 3.7 depicts well the convolution operation and the pooling.

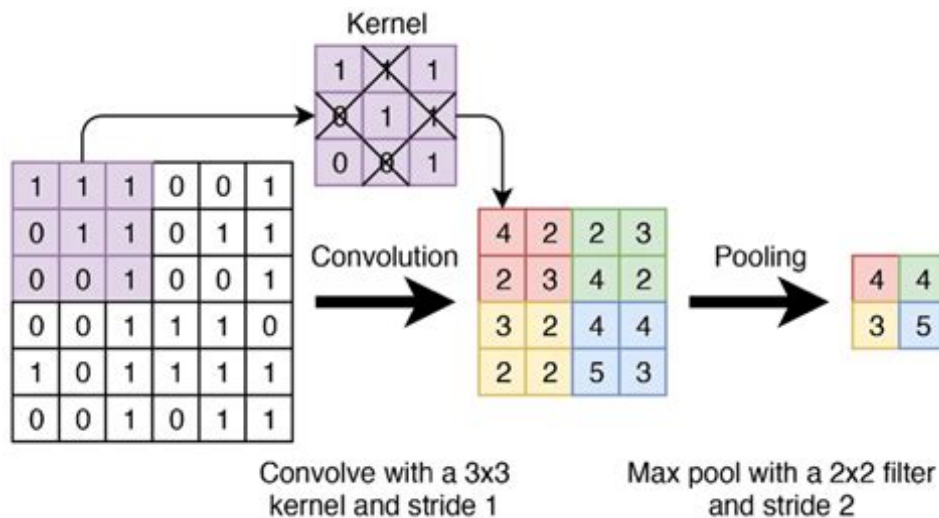


Figure 3.7: Example of convolution operation and pooling operation (max pooling).

Usually pooling layers are placed between consecutive convolutional layers. This allows the reduction of the size of features' representation and also contributes to prevention of the overfitting. The pooling layers do not have parameters that are optimised during the training process, however, the kernel size, the stride and pooling type can be tuned in order to build a better model.

Near the output of the model, the dense layers, or also called fully connected layers, are usually placed. These layers are composed of several neurons to which

all the outputs of the previous layers are connects. The combination with the non-linear functions like softmax produces the networks' output. In the classification type problem, the number of the neurons at the last fully connected layer is equal to the number of possible classes [40][41].

3.3 Computer Vision and Cervical Cytology: State-of-the-art Review

General cell detection, segmentation, and counting are rather well addressed in the literature, and most recent proposals mostly rely on machine learning and deep learning approaches. Particularly, [42, 43, 44] describe and propose using deep learning models and architectures such as U-net and Feature Pyramid Network (FPN) to produce cells detection and segmentation. In [45] the author describes using Single Shot Detector (SSD) in pair with a CNN to localize blood cells and then count them separately. On the other hand, [46] proposes a microscopy cell counting based on density estimation employing fully convolutional regression networks. While all of the approaches output the number of cells, the segmentation also outputs the mask from the detected objects. In decreasing order of output complexity, the detection task normally results only in the localization of the object and a bounding box around it. Lastly, density estimation only gives the final number of objects. In all approaches mentioned above, authors outline results with performance comparable with a human specialist.

A less explored field in the literature is to devise methods that enable implementing efficient smear quality control approaches that could contribute to increase the performance of the detection. In [47], the authors describe an AI assistive diagnostic solution to improve cervical liquid-based thin-layer cell smear diagnosis according to clinical TBS criteria. The developed system consists of five AI models which are employed to detect and classify the lesions. The quality assessment is done on the entire sample and comprises technical image characteristics, such as focus and contrast, and quantitative cell evaluation. In contrast to the detection and classification tasks, the cell counting task relied on simpler methods. The number of the cells was estimated by, firstly, engaging the Otsu thresholding method to separate the cells from the background and then by calculation of the cell to overall area ratio to obtain a rough number of cells in the sample. The authors report a 99.11% average accuracy on the task of samples classification as satisfactory and unsatisfactory on the validation sets. However, it must be noticed that in their approach, the total number of cells was estimated, while the TBS guidelines specify a minimum threshold specifically for squamous nuclei.

Besides the quality of the data and the performance of the detection methods, some effort has also been put into developing low-cost, portable microscopes that

enable supporting microscopy-based diagnosis in areas with limited access or without enough financial resources. As an alternative to benchtop microscopes, and due to the impressive evolution in the quality of the cameras, processing power and memory, smartphone-based solutions are good candidates for implementing a cost-effective platform for microscopic inspection of samples. The work in [48] describes a portable smartphone based brightfield microscope for screening blood smears. A wide range of applications have also been used to test the feasibility of affordable approaches based on smartphones. Examples include the detection of viruses [49], the quantification of immunoassays, the automated classification of parasites [50], etc [51, 52].

Chapter 4

Previous work and CLARE dataset

In order to get a better understanding of the technologies chosen in the methodology and the dataset used, it is important to contextualize the previous and current work in the Fraunhofer project's workflow. Therefore, in this chapter, a brief description of previous work done in the related projects will be made, as well as an overview of the dataset used in this thesis.

4.1 μ SmartScope for cervical cytology

Through the last years, Fraunhofer has been investigating the ways of automatic computer-aided methods of analysis of microscopic samples. The project DEMalariaScope, which aims to develop an automatic detection of malaria in blood smears using smartphones, and the project CLARE, which stands for Computer-Aided Cervical Cancer Screening and targets to create a novel framework that can be used as a Decision Support System (DSS) for the screening of Cervical Cancer, are both examples of a workflow where a sample has to be digitalized using a microscope [53][54]. Since these projects are aimed to solve the screening problems in low incoming countries, where the access to proper equipment and trained personnel can be limited, the key requirement was to develop a cheap and easy-to-use system that could be deployed in the field effortlessly. This requirement led to the development

of the μ SmartScope, an automated 3D-Printed Smartphone microscope with a motorized stage [4]. As can be seen in Figure 4.1, the μ SmartScope consists of a 3D printed chassis with a smartphone on the top. It encapsulates optics, illumination, stage module and the necessary electronics. The stage system is responsible for moving the sample in three different axis, which allows the automatic positioning and focusing of the samples. The smartphone is connected to the μ SmartScope via Universal serial bus (USB) cable, providing power to it and serving the role of a master computer, where it is responsible for the slide positioning, image acquisition and image analysis.



Figure 4.1: The μ SmartScope prototype, with smartphone attached and microscopic slide inserted [4].

The Figure 4.2 shows the overall pipeline of the app installed on the smartphone that is responsible for the operation of the μ SmartScope. The purpose of this flowchart is to contextualize where the outcome of this work will be implemented in the final system. Being aware of it is important because the final implementation will also add constraints to our module, such as maximum processing time and memory limitations. Without mentioning the app details, on the high level, the first thing to do after the start of the application is to position the sample slide and acquire the image. After that, the pre-processing module should take place where the image brightness, contrast, and other properties are to be adjusted. Afterward, our module is ready to take the processed image as input and output the results in parallel with other optional modules. When the processing of the different modules is finished, the information is aggregated, and the slide should be positioned again in order to acquire the image of the next area. If there are no more uncovered regions, the app should output all the aggregated data to the operator and stop. Even though it is not reflected in the pipeline, other modules can run before or even after our module and do not compelled to run always in parallel. Since the development of

the application is out of the scope of this work, it is impossible to predict or impose the final application architecture.

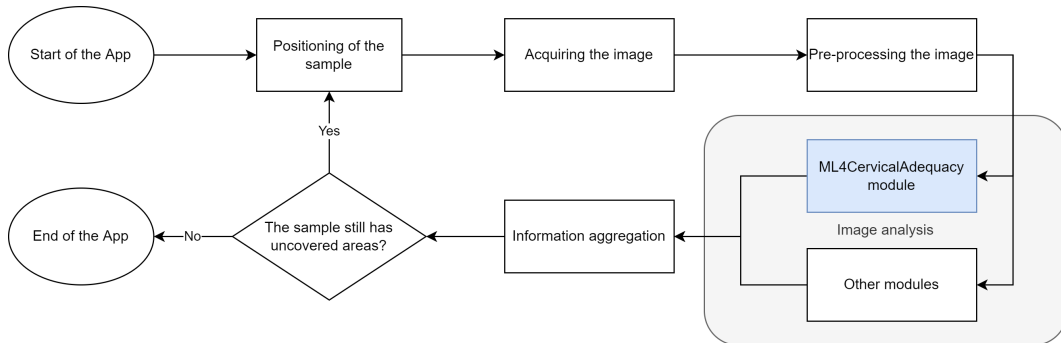


Figure 4.2: The μ SmartScope application pipeline. The ML4CervicalAdequacy module correspond to the module developed in the scope of this thesis.

The expected use case of the μ SmartScope from the operator’s perspective consists of, after obtaining an LBC slide, placing it under the lens, selecting the desired operation, and waiting for the results. Ideally, image analysis modules should not take more time than the module responsible for obtaining images from the sample slide.

At the time of the development of this work, the μ SmartScope already has been through some prototype iterations and is capable of digitalizing samples [55][56]. In the scope of the CLARE project, researchers from Fraunhofer AICOS have already tested the microscope with LBC samples and created a private dataset with cervical cancer samples.

4.2 CLARE dataset

Even though there are some publicly available datasets with nuclei annotations of cervical cells, such as Herlev[57], SIPaKMeD [58], Cervix93[59], and ISBI Challenges[60, 61], they are not adequate for our purpose. In the case of Herlev and SIPaKMeD, they only contain isolated images of cell nuclei, with annotations by abnormality of the cell. While ISBI Challenges and Cervix93 do offer images with the surrounding area, they are not annotated in terms of the type of nuclei. Hence, it would not be possible to use them in our application. Thus, our choice will be a private CLARE dataset.

The CLARE dataset is a private image dataset collected by Fraunhofer Portugal (FhP) in the ambit of the CLARE project. The CLARE dataset is being collected with the μ SmartScope prototype at IPO-Porto and Hospital Fernando Fonseca (HFF), and contains annotations of cellular structures highly relevant for sample adequacy assessment. The dataset is comprised of 40 samples from IPO and

99 from HFF. Each sample contains around one hundred images, which were taken along the sampling area with amplification of $40\times$ (see Figure 4.3).

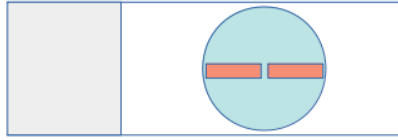


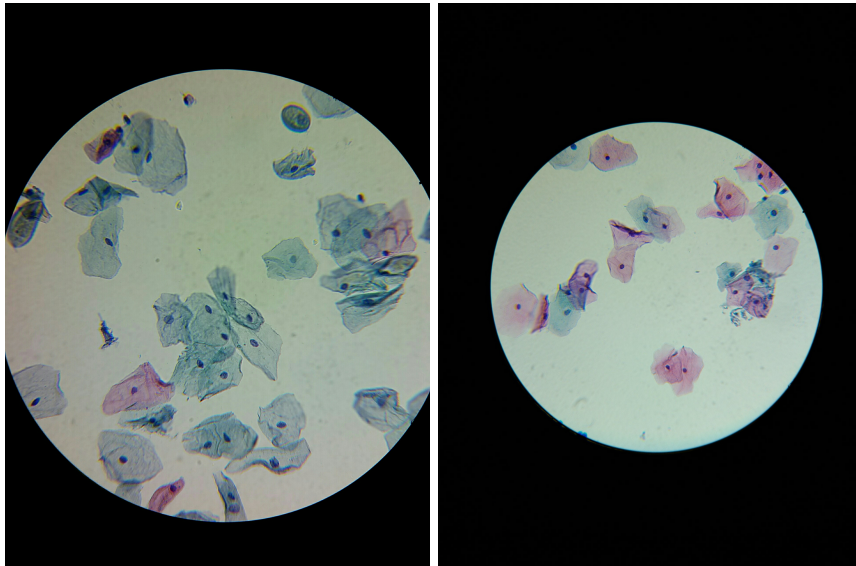
Figure 4.3: Microscope slide with LBC preparation (blue) with area of sampling (orange). Images of the sample are taken along the orange area with $40\times$ amplification.

The images were taken with such smartphones as Samsung Galaxy S5, Google Pixel 2, LG G6, and Google Nexus 5. Images are approximately 5 Mpixels in size, ranging from 1920×2560 pixels to 1944×2592 . Figure 4.4 shows some examples of images present in the dataset. All four devices have different sensors and optical systems and, therefore, have different field of view (FOV). This leads to different circle sizes when compared with the overall image size. From all the four devices, the Samsung Galaxy S5 seems to have the smallest FOV and, so, the biggest circle, which also translates to better object resolution. While circles in the images from LG G6 and Google Nexus 5 appear to be of the same size, the image obtained with Google Pixel 2 has the smallest one.

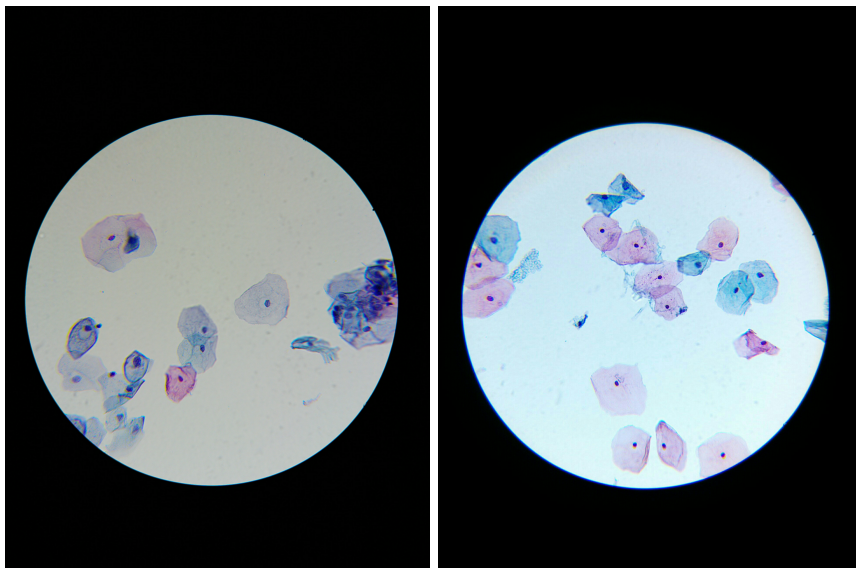
It is important to note that the LG G6 and Samsung Galaxy S6 appear to produce images with the size of 1920×2560 pixels, while smartphones from Google produce images with the size of 1944×2592 pixels. From the specification sheet of each smartphone, it is noticeable that all four smartphones embody image sensors capable of producing images with 8 Mpixels or more and retaining even more information. However, it is possible that in this application, there is not much benefit from increasing image resolution on the current smartphones, and it only leads to slower processing speed.

The dataset is annotated in terms of nuclei type with a bounding box for each nucleus along 6 different classes. It should be noted that this dataset is also partially annotated in terms of abnormal lesions, but this type of annotation are not relevant for our task of cellularity assessment. However, since manual annotation is a very tedious work and requires the time of trained personnel from IPO and HFF, not all images and samples have been annotated yet. Currently, the annotated part of the dataset is constituted of 21 samples from HFF and 20 samples from IPO. It is also worth mentioning that, from these samples, not all images were annotated, but only 15-20 from each sample, totaling 765 annotated images.

The nuclei annotations are divided into four main classes: squamous nucleus, inflammatory cell, glandular nucleus, and artifact. However, there are two additional



(a) Image obtained with Samsung Galaxy S5. (b) Image obtained with Google Pixel 2.



(c) Image obtained with LG G6. (d) Image obtained with Google Nexus 5.

Figure 4.4: Examples of raw images that compose the dataset with one image per smartphone.

classes: undefined nuclei and undefined all. These classes were created to show the uncertainty of the annotators against the object. In the case of undefined nuclei, the annotator meant that it can be either squamous or glandular. In the case of undefined all, it means that it can be squamous, inflammatory or glandular. As per artifact, they just mean that the object probably was not supposed to be there and should not be mistaken with nuclei. The dataset is also clustered in terms of cell abnormality and comprehends possible regions of interest. Samples are divided into seven different clusters: Atypical squamous cells (ASC-H); Atypical squamous cells of undetermined significance (ASC-US); High-grade squamous intraepithelial lesion (HSIL); Low-grade squamous intraepithelial lesion (LSIL); Squamous cell carcinoma (SCC); Negative and Inadequate. The first five mean that there is an abnormality present in the sample, with respective level, while the negative clustered refers to samples with absence abnormalities. The last cluster is inadequate, this means that these samples are not adequate for further analysis. The distribution of clusters can be seen in Table 4.1. Clusters with ASC-H and HSIL are noticeably predominant, but for our use, we can group all clusters aside from inadequate. By doing this, we can account for 38 adequate samples and three inadequate. This would be useful for the final tests, where the whole sample would be submitted for evaluation. Then, the module output, which would be in the form of adequate or inadequate, would be compared against the ground truth. It is unfortunate that in the annotated part of the dataset only three inadequate samples are present. However, in terms of nuclei type, we also had access to two digitalized but not annotated samples that are known to be inadequate, which can not be used in the development but could be used in the final tests.

Table 4.1: Dataset distribution by clusters/abnormality levels.

Cluster name	N ^o of samples
ASC-H	13
ASC-US	4
HSIL	15
LSIL	2
SCC	2
Negative	2
Inadequate	3 (5) ¹
Total	41

A visual comparison between adequate and inadequate samples (see Figure 4.5) reveal a very scant number of nuclei, implying that the reason for being considered inadequate is low cellularity.

¹Two additional samples are not present in the dataset due to lack of nuclei annotations, however, all 5 samples will be considered for the system evaluation of the classification on the sample level.

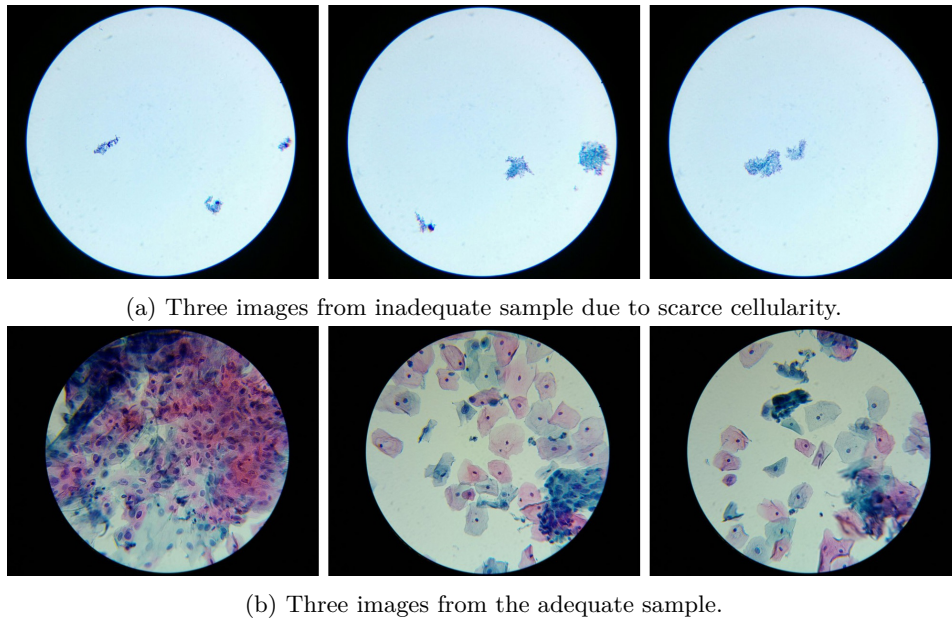


Figure 4.5: Visual comparison between two samples. Sample 4.5a shows an inadequate sample with a scarce amount of nuclei. Only two or three can be identified in all three images. On the other hand, the adequate sample 4.5b is almost fully covered with squamous cells.

Notwithstanding, the part of the dataset that was annotated is still relatively big. As shown in Table 4.2, there are more than forty thousand annotations.

Table 4.2: Dataset distribution by nucleus type.

Class name	N ^o of annotations
Squamous nucleus	21973
Inflammatory cell	17914
Glandular nucleus	684
Undefined nucleus	667
Undefined all	997
Artifact	152
Total	42387

A quick look at the numbers for each class unveils two predominant types: squamous and inflammatory. In the machine learning domain an imbalanced dataset can turn out to be a concern. The lack of examples of glandular, artifacts and undefined types can lead to simply model being unable to learn any features. In this particular dataset, this imbalance can be even more concerning due to a very similar appearance of different nuclei.

Figure 4.6 provides examples of nuclei present in each class. The squamous nuclei, Figure 4.6a, are usually encapsulated in a cytoplasm, distinguishable by blue or pink semi-transparent involucres. Figure 4.6b exemplifies inflammatory nuclei.

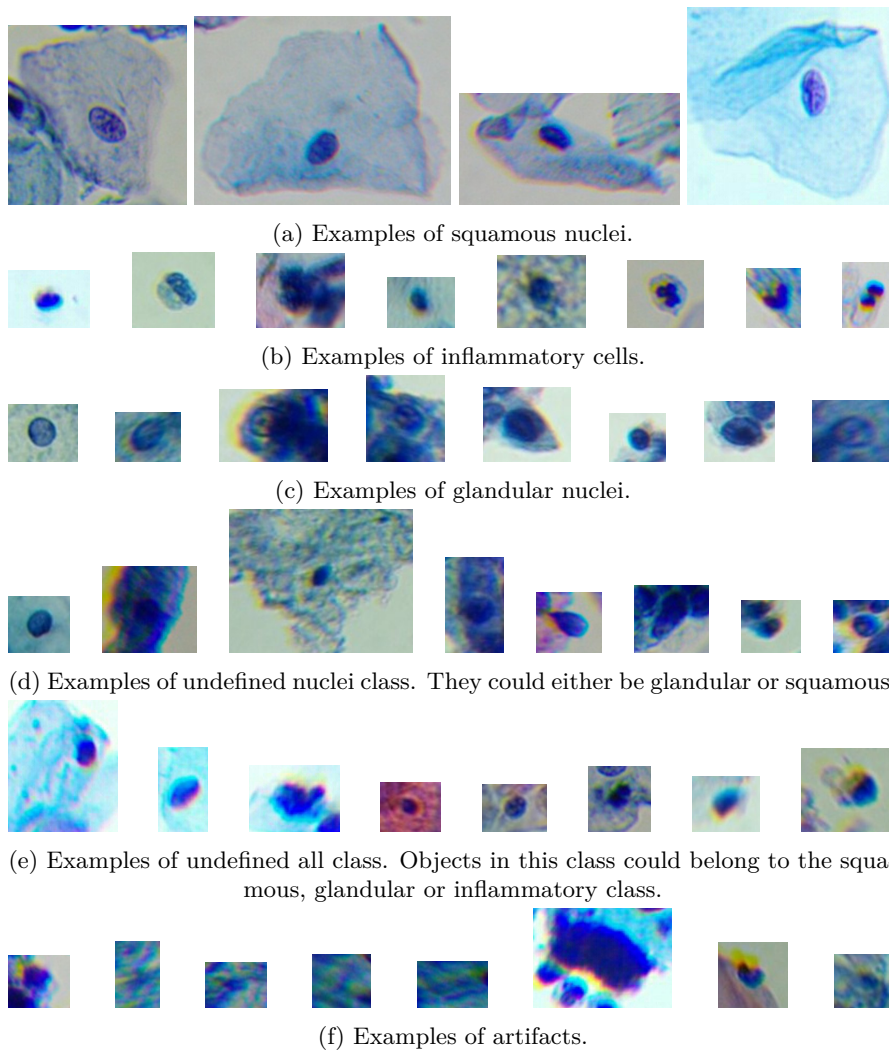


Figure 4.6: Examples of cropped objects of each class in relative scale to each other.

In contrast to squamous, inflammatory nuclei usually are not encapsulated in the cytoplasm. Glandular nuclei, present in Figure 4.6c, are slightly different from squamous and inflammatory and typically could be found in structures exemplified in Figure 4.7. As per undefined classes, Figures 4.6d and 4.6e mimic visually, as expected, the previous three classes. They may belong to one of these classes, but it was unclear for the annotator for which one exactly. The class artifact, Figure 4.6f, contains objects that are not usually found in a sample slide. Their origin may be external, which can occur during slide preparation, or also can be an object that is not supposed to be in the slide and, therefore, not confused with an actual cell or nucleus.

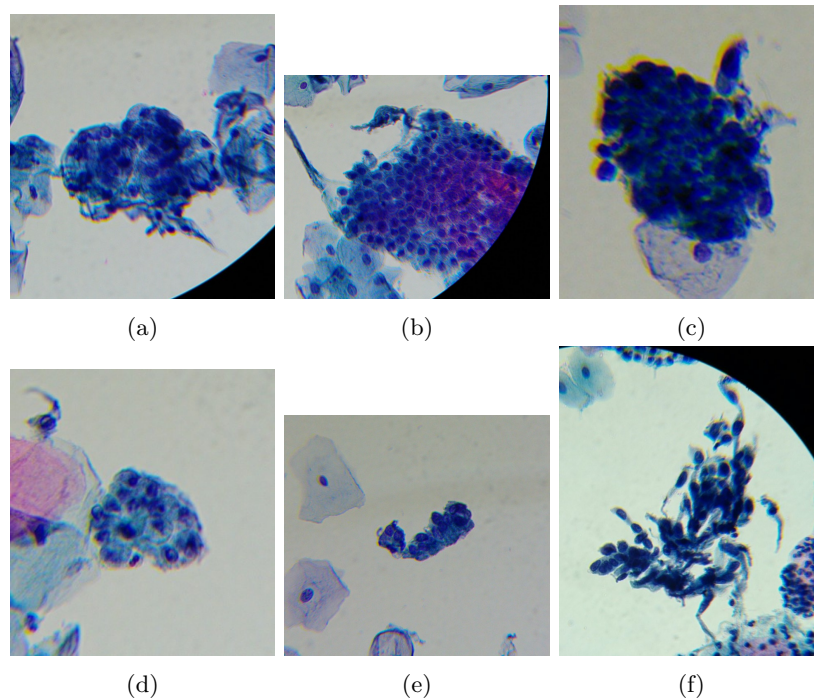


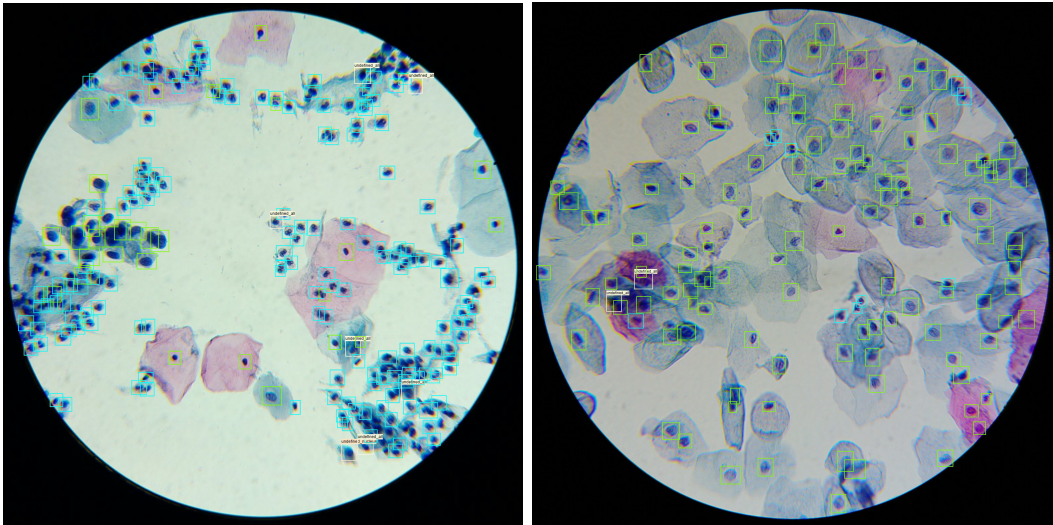
Figure 4.7: Examples of images with a dense arrangement of glandular nuclei.

4.2.1 Missing annotations

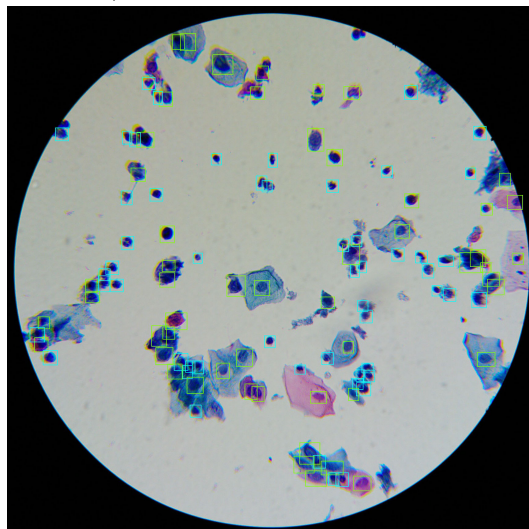
The visual examination of raw images with annotations on top of it, see Figure 4.8 and Figure 4.9, discloses a little more information about the dataset. At a first glance, almost all cells appear annotated. However, in some images, one can notice objects that could be a cell but were not annotated as it. These objects, presumably cells, are relatively easy to spot on Figure 4.9a, 4.9b, 4.9c, and 4.9d. In a meeting with experienced medical specialist from IPO-Porto, it was clarified that annotations are lacking due to the low confidence of the annotator in the type of the cell. The low confidence mainly comes from overlapping nuclei or appearing out of focus, therefore, containing too little information to enable any meaningful classification. In these figures, it is also possible to spot the main difference between squamous and inflammatory nuclei. While squamous cells are enclosed in the cytoplasm, inflammatory nuclei lack it. Different colors present in the cytoplasm are due to dye, which purpose is to help visually separate cells.

It is also important to note that, in a meeting with medical personnel, it was mentioned that the annotator's experience and factors such as fatigue or mood are highly influential in the process of classification of cells. This can lead, for example, to the same cell being classified differently by two different annotators, or even being annotated differently by the same annotator at different times. These findings should be taken into account in the evaluation and discussion steps of the developed

algorithm.

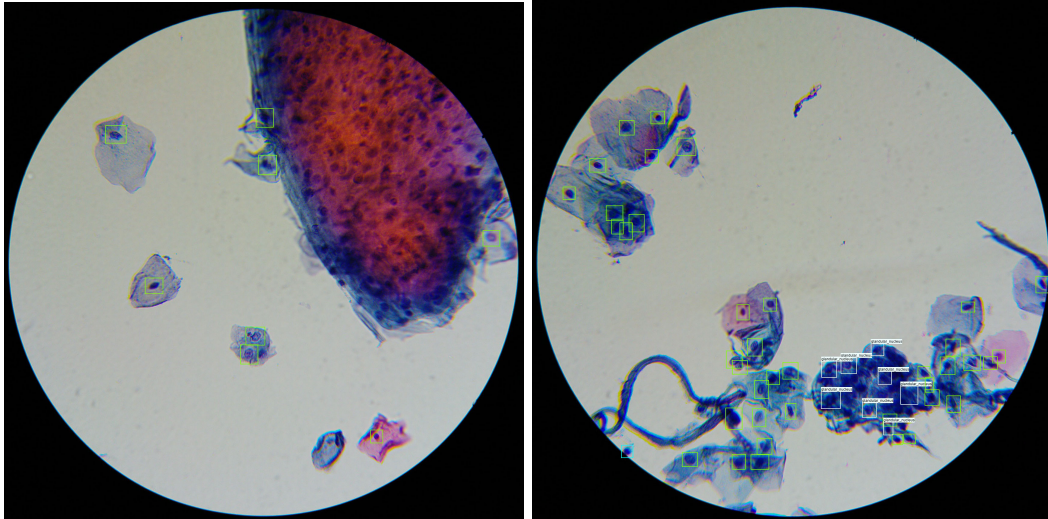


(a) Image with predominant inflammatory nuclei (blue bounding boxes). (b) Image with predominant squamous (green bounding boxes).

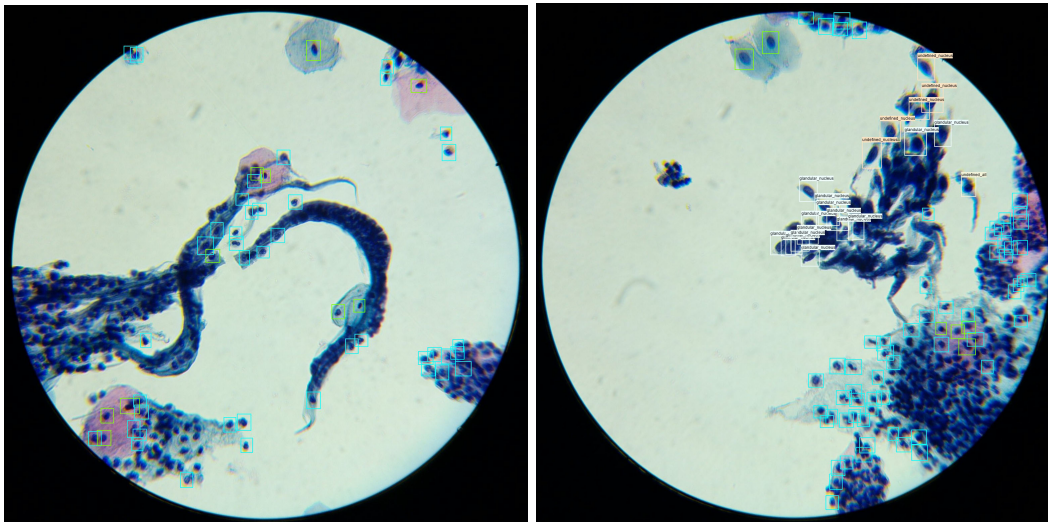


(c) The distribution between squamous and inflammatory is even.

Figure 4.8: Example of the distribution of two predominant types of nuclei in samples. There are images where mostly only one type of nuclei is present, such as Figure 4.8a and Figure 4.8b. Figure 4.8c shows a case where the distribution is even.



(a) A few squamous cells, marked with green bounding boxes, close to the agglomerate of a strange C-shaped form object, and agglomerate of unannotated cells. (b) A few squamous cells (green bounding boxes), and agglomerate of glandular cells (white bounding boxes).



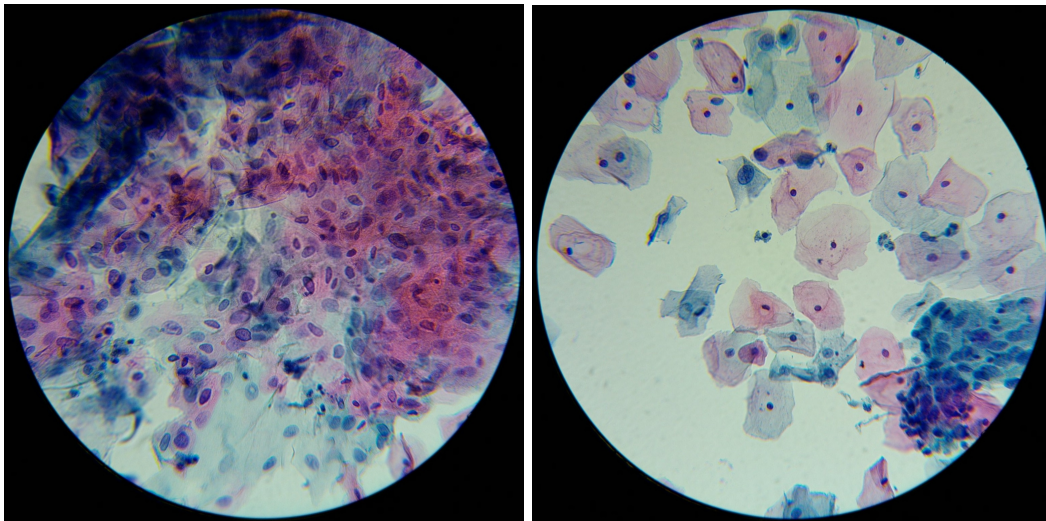
(c) An example of a dense arrangement of cells. (d) Also an example of a dense arrangement of cells, with some glandular (white boxes) and undefined (pink boxes) nuclei.

Figure 4.9: Examples of images with a dense arrangement of nuclei. Some nuclei appear without annotations.

4.2.2 Data particularities

Overlapping cells

Even though LBC provides better-looking samples compared to CPS, it still has its flaws or particularities. When preparing an LBC sample, one of the goals is to equally distribute cells in a uniform fashion across the slide surface and thus eliminating the overlapping cells [62]. Nonetheless, none of the existing processes is perfect, and even LBC provides samples with overlapping cells. In Figures 4.4, 4.8 and 4.10b single cells can be visualized separated from others. This case is considered relatively easy to assess. However, cases such as the ones present in Figures 4.9 and 4.10a only add difficulty to the assessment and annotation process since it is unclear the type of the cell and its characteristics, such as the presence or lack of cytoplasm. Particularly, Figure 4.10a, is a clear example of heavily overlapping cells. As can be expected, in this case, it can be extremely difficult to classify present cells.



(a) Image fully covered with overlapping squamous cells. (b) Almost no overlapping is present in this image.

Figure 4.10: Example of cell overlapping problem, 4.10a, and comparison to a clear image, 4.10b.

Out-of-focus

Although for a human eye, the slide preparation may seem like an almost perfect two-dimensional layer of cells on the white background, the cellular structures are at different depths, which allows cell overlapping problems to occur. Due to this aspect of LBC preparation, it is necessary to focus on the plane of the cells to obtain a sharp image of them. Anyhow, not all cells are positioned on the same plane, which leads to some of the cells appearing blurred, hence, out of focus. This can be observed in Figure 4.10a, where cells positioned on the same plane at which the image was

focused at, are sharp and detailed, whereas out-of-focus cells appear very blurry and almost blend with the background.

Chromatic aberration

In some pictures present in Figure 4.6, it is noticeable blue and orange fringing around objects. This effect is called chromatic aberration and means a failure of a lens to focus on all colors to the same point. The cause of the chromatic aberration is the phenomenon called dispersion. In optics, where the refractive index describes how fast the light travels through the material, its value also tells how the light will bend when traveling through different materials. It happens that the refractive index varies with the wavelength of light and so causes the chromatic aberration [63]. In Figure 4.11, it is possible to observe an example of the chromatic aberration effect. On the right and left parts of the image, this effect is more pronounced since it is where the light suffers more bending, as it goes through the optical system, in opposition to the objects in the center. While the objects are still distinguishable when encountered singly. It turns to be more difficult to classify them due to uncertain appearance.

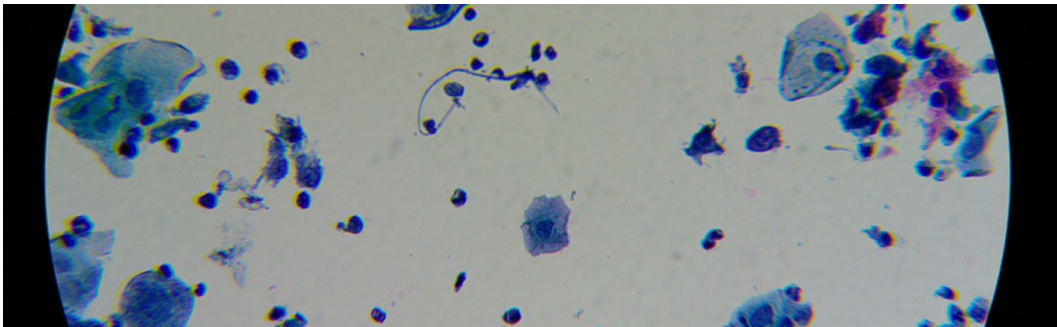


Figure 4.11: Example of chromatic aberration. The image was cropped to its actual size to facilitate observation.

There are a few solutions to reduce this effect. However, usually, it requires more complex optical systems, which costs more. Nonetheless, there are also post-processing algorithms that aim to minimize the chromatic aberration, which will be briefly referred in the next chapters.

Chapter 5

Methodology

With already gathered knowledge from the previous chapters regarding state-of-the-art techniques and previous works, this chapter aims to present our approach along with the methodology used in the research. In the beginning, the overall module pipeline will be described in order to give an understanding of the direction of the development. The second section will target the actual methodology, starting by presenting our research approach and going to its details. Then, the last section will finish with the evaluation methodology.

5.1 Overview

As stated in the introduction, the focus of this work is the assessment of the quality of cervical cytology smears. We have already seen that the most common cause of unsatisfactory cervical cytology smears is low cellularity. There are guidelines given by The Bethesda Systems that specify criteria for cellularity: At least 5000 cells on liquid-based cytology (LBC) slide prepared with ThinPrep, which corresponds to 3.8 cells per high-powered ($40\times$) field. Usually, a cellularity assessment on the LBC slide is done either by comparison with reference images or by counting well-preserved squamous cells in a defined number of fields at high or low power ($10\times$) [64].

Since LBC slides provide relatively clean images where a human eye can distinguish most of the cells individually, as can be observed in Figures 4.4, 4.8, and 4.10, for example, our approach will rely on counting objects separately. In Chapter 3,

we provided studies in a similar context that uses deep learning object detection models and yield satisfactory results.

The flowchart for the module pipeline can be observed in Figure 5.1. This flowchart corresponds to the module of ML4CervicalAdequacy presented in Figure 4.2. Taking into account that the future use case implies running the developed model on the smartphone, as well as the computational power required by current state-of-the-art object detection architecture and the image dimensions of our dataset, it has been decided to process images in patches in order to decrease the required time for training and inference, as well as the required computational resources and memory. Thus, the first operation slices the images into patches of fixed dimensions, as shown in the pre-processing module in Figure 5.1. To overcome the downside where we may lose some information around the objects of interest when those happen to be at the edges of the patches, we allow slicing with a predefined percentage of overlap. Then, patches are fed to the object detection model, and outputs are collected. Since each patch will pass individually through the model, the outputs, which will be in the form of bounding boxes and respective classifications, will be in reference to the patch. In order to count them properly, they must be transferred to a global reference. Also, in case of using patches with overlap, a Non-maximum Suppression (NMS) algorithm is used to eliminate duplicates. After this post-processing, all the detected objects are counted and stored. The process will start again as the new image become available for the module. When all images from the slide/sample are processed, a total number of counted objects is presented along with the classification result of the sample: Adequate or Inadequate.

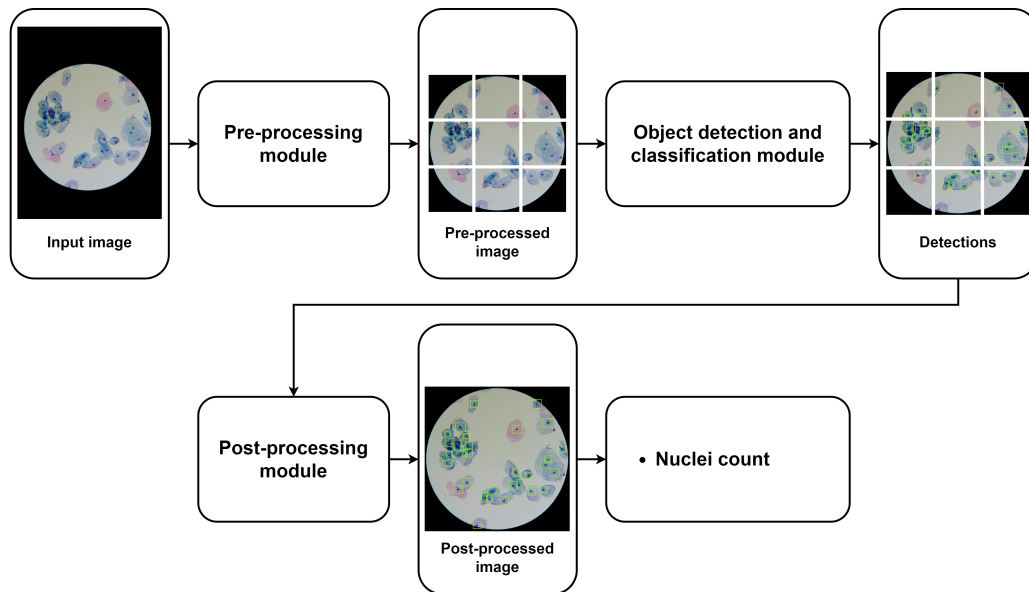


Figure 5.1: The ML4CervicalAdequacy module pipeline.

5.2 Proposed approach

While there are not many questions on how to slice an image and then reconstruct it, developing a machine learning model can be a bit more complex task.

There are two main approaches in the development of a machine learning model: Model-centric and Data-centric. The model-centric method consists of developing experimental research in order to improve the model performance, which involves selecting the best model architecture and its hyperparameters. In this method, the dataset may suffer some general pre-processing at the beginning of the development or at the beginning of each training process. However, it is not subjected to changes through the development process. On the other hand, the data-centric approach consists of fixing the model and its parameters and only manipulating the dataset that is subjected to the training process. These two approaches can be visualized in Figure 5.2.

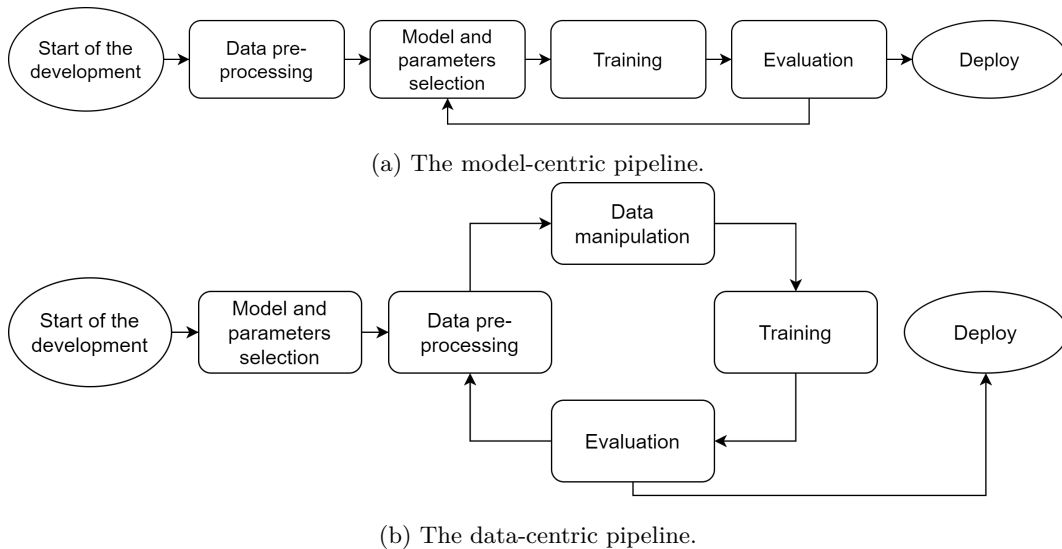


Figure 5.2: Model-centric and data-centric development pipelines.

Data is a critical part of machine learning, and its quality directly reflects on the model performance. Until recent years, most machine learning competitions were focused on optimizing the model to its best. However, it became more and more clear that it was not the only way to optimize the model performance. A recent competition held by DeepLearning.AI demonstrated a twenty percent improvement compared to the baseline, using only a data-centric approach [65, 66].

Our model development will combine some aspects of both model-centric and data-centric approaches. Firstly, a model selection will be made and then, a hyperparameters search performed. After encountering the best model and parameters, they will be fixed, and the dataset will be manipulated in order to improve further the performance of the model.

5.2.1 Tools and model selection

Considering state-of-the-art object detection techniques, we decided to base our approach on TensorFlow (TF) Object Detection API. TensorFlow has been in an active developing state during the last years. This led to a relatively mature base of developers, well-documented API, a community of developers willing to help, and lots of information already on the internet. TensorFlow community provides TensorFlow 2 Detection Model Zoo, a collection of pre-trained detection models. These models can be used for out-of-the-box inference or for initializing models when training on a novel dataset. Since COCO 2017 dataset, the dataset on which Zoo models were trained, does not contain any categories that can be useful for us, we will only be interested in using pre-trained checkpoints for starting training on our dataset.

Model Zoo contains more than forty models, and, due to time and hardware constraints, we would not be able to test every model with a few variations of other parameters. That being said, we selected three models, see Table 5.1, to perform experiments with different hyperparameters.

Table 5.1: Models considered for training.

Model name	Input size	Speed (ms)	COCO mAP
SSD MobileNet V2 FPNLite	320×320	22	22.2
SSD ResNet50 V1 FPN	640×640	46	34.3
EfficientDet D0	512×512	39	33.6

The selection of these models was based mainly on the speed of inference and the performance of the COCO mAP metric. Since this work aims to contribute for the development of a end-to-end lightweight mobile-based and accessible CADx system for cervical cancer screening, the further possibility of porting the model to mobile devices has been taken into account.

The prioritization of speed also affected the input size of the models. To be able to perform any experiments with different batch sizes other than 1, we chose the models with 320×320 pixels of input size. It is important to note that the TensorFlow Model Zoo does not provide pre-trained ResNet50 and EffientDet models with the input size specified above. However, the configuration file allows changing the input size.

5.2.2 Pre-processing

Dataset preparation

Currently, the dataset is composed of images of .jpeg extensions and adjacent .json files, which contain information about annotations in each image. For TensorFlow API to be able to ingest information and process it in an efficient way, there is the need to transform the dataset into a .tfrecord format. This format is a proprietary format of TensorFlow framework. The main purpose of this format is to make the process of the training and online data manipulation easier, i.e., during the training process.

Jumping a bit to the end of the research pipeline, which is evaluation, to perform it correctly, the split between train and test sets is required. Usually, this split is done with an 80% to 20% ratio, where 80% of the data goes to the train set and the other 20% to the test set. The reason for splitting data this way is to guarantee that when evaluating the model, it is used the data that have never been shown during the training phase, thus avoiding evaluation bias. And, therefore, the evaluation will try to reflect real-world performance since it will always see the new data after the deployment.

Also, for the correct search for the hyperparameters, the cross-validation technique is recommended. This technique consists of splitting the train set into K group of train and validation subsets, where K is a number of these groups. It is suggested that K is chosen so that it splits the data evenly. Taking as an example a three-fold cross-validation, the first group will, for example, split the first 67% percent of the train set into train subset and another 33% into validation subset. The second group will maintain the split ratio but, the training and validation subsets are shifted, so it will contain different data in the training subset and validation subset when compared to the first group. The last group will shift the data even more, so it also contains a bit different set of data. It should be noted that before each split, the data is also shuffled in order to get a random distribution across all folds. The search for the hyperparameters will, then, be performed on each fold for each set of parameters. This way, it can be ensured that the results are not biased, as the same set of parameters is evaluated on three different train data sets and also on three different validation sets. In this technique, the validation sets serve the purpose of the test set. Figure 5.3 illustrates the split performed on our dataset, which can also be visualized with greater detail in the Appendix A. The Train set accounts for 80% of the data and the Test set for the remaining 20%. The number of folds was set to three, and each fold is split with 67% and 33% ratio for train and validation subsets, respectively.

It is important to note that the first train and test split is done on the sample level. Since not all the samples contain an equal number of different nuclei, in order

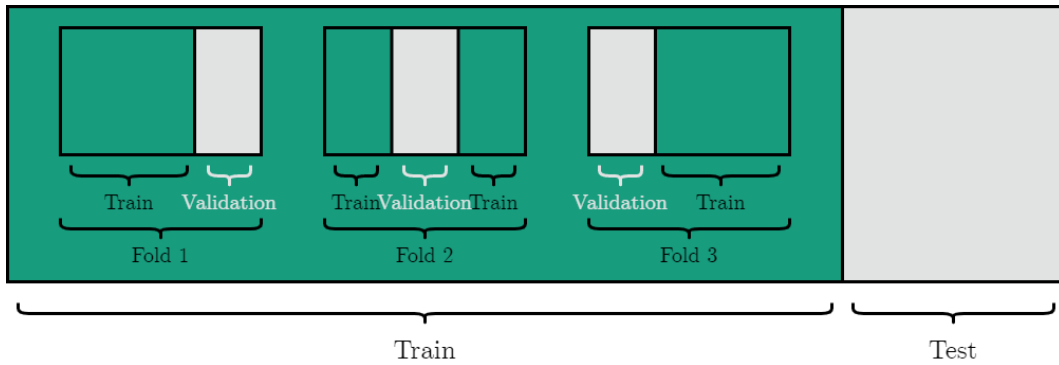


Figure 5.3: Visual representation of dataset split into train and test sets. Also, cross-validation split into three folds can be visualized.

to account for the presence of them in the more or less 80 to 20 ratio, the split was performed manually and the ratio of each class in train and tests sets can be seen in Figure 5.4. It was also ensured that two of the inadequate, from a total of three, samples are present in the test set.

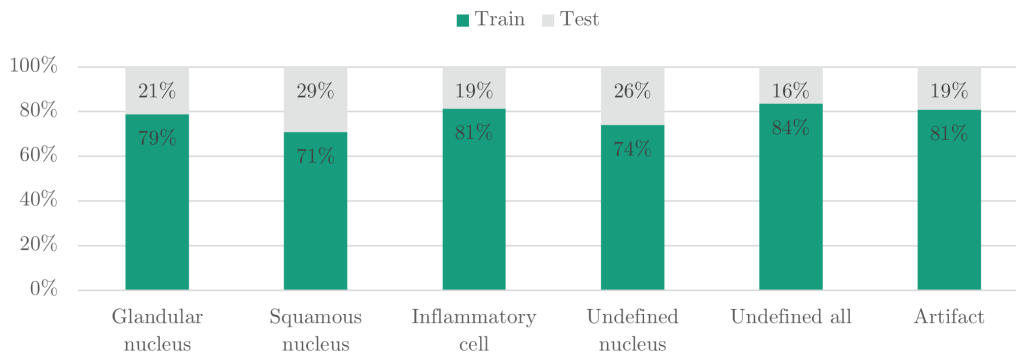


Figure 5.4: Percentage of annotations in each set, train and test, for each class.

As per cross-validation splits, Figure 5.5 demonstrates annotation distribution per subset of Train and Validation of the first fold. At this stage, the splits were done randomly on the image level. As the distribution of the cells is not homogeneous across all the images, it cannot be ensured that the distribution of annotations will also follow 2/3 for the test set and 1/3 for the validation ratios. This explains the slightly different distribution of undefined nuclei.

The other two folds follow the same distribution with minor deviations in classes with a small number of annotations.

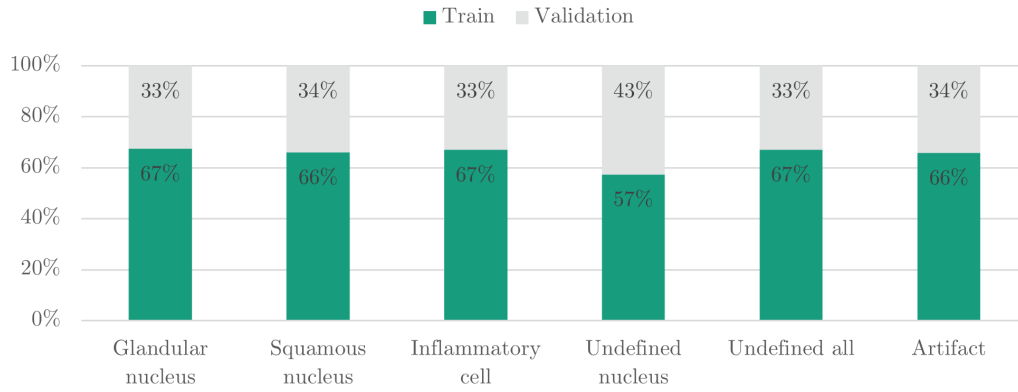


Figure 5.5: Percentage of annotations of the fold 1 in each set: train and validation; for each class.

Slicing

As mentioned before, the images will be sliced into smaller patches. The size of the patches has to be equal to the input size of the selected models, i.e., 320 by 320 pixels. It has been decided to perform slicing on every 320 pixels horizontally and vertically. However, as later discovered, it presented some problems: sometimes it results in slicing on the middle of the object. Pairs 6.1a/5.6b and 5.6c/6.1b in Figure 5.6 display the problem when slicing results in cutting in half annotated objects. Keeping the annotation of the object on two patches would not be the correct solution since it would double the number of annotations for some objects.

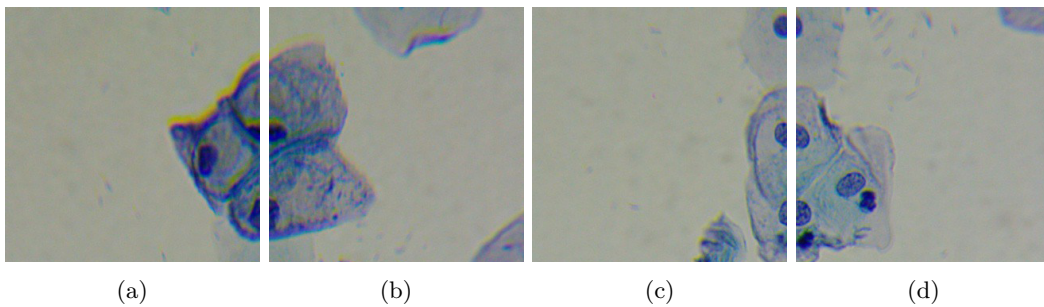


Figure 5.6: Example of cut cells.

While a solution could be to perform any type of vertical or horizontal shift in order to preserve the object entirely on the patch, there are areas with a high density of objects where this solution would not work. To mitigate this issue, it was decided to introduce a threshold for a minimum percentage of an object that has to appear in the image to be considered and still to carry an annotation. Through some tests, it was determined to set the threshold to 50%. Thus, the annotation will only be associated with the patch if it contains at least 50% of its bounding box on the respective patch. When the bounding box is divided exactly in half

between two patches, we decided to exclude these annotations, which was done to prevent double annotation for the same object. Figure 5.7 shows some examples of eliminated annotations and ones that were kept. Also, it is worth mentioning that a simple slicing would also lead to the presence of black patches. To address this issue, prior to slicing, the image is cropped to the size of the optical disc, as can be observed in the Figure 5.8.

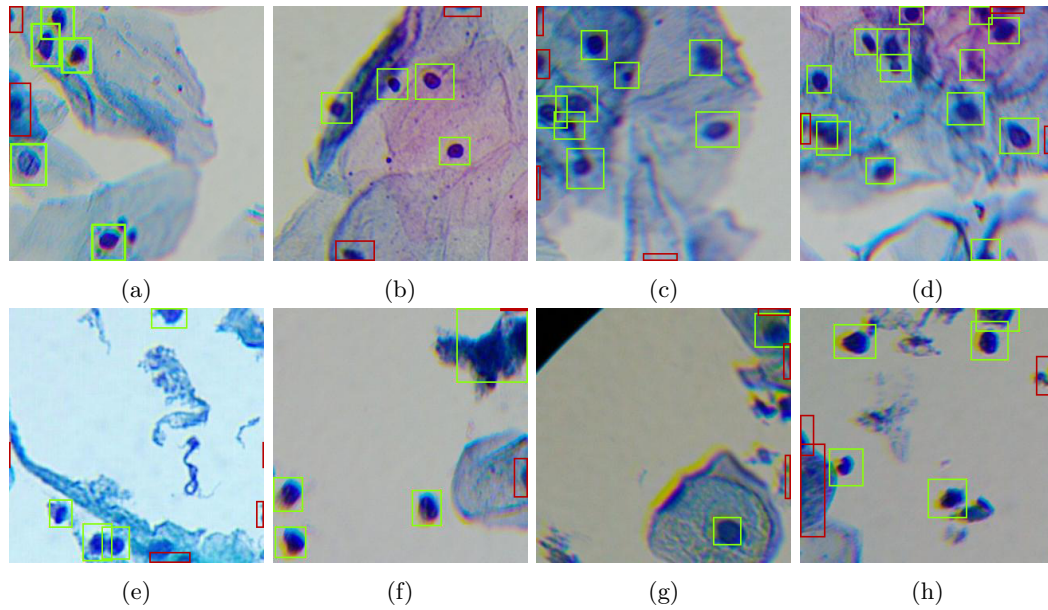


Figure 5.7: The green boxes indicate the annotations that were kept on the patch. The red boxes indicate that these patches were eliminated from the patch. It is possible to verify that most of the eliminated annotations are not containing any meaningful information and, therefore, should not be kept.

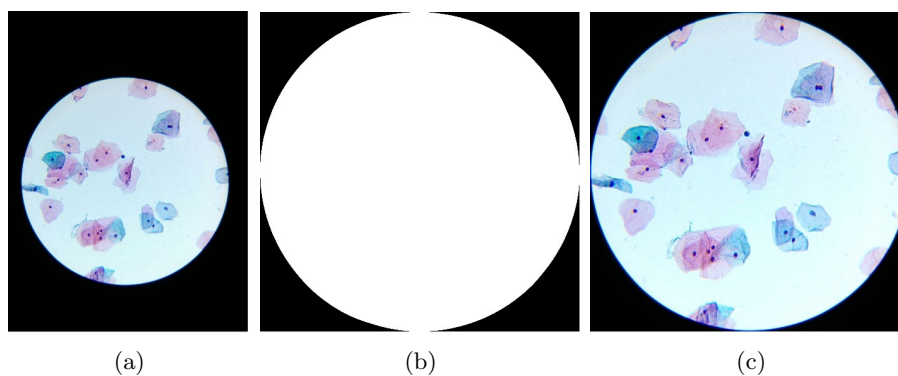


Figure 5.8: Cropping image to the size of the optical disk: (a) Original image; (b) Segmentation mask of optical circle with cropping; (c) Original image cropped to the size of segmentation mask.

5.2.3 Hyperparameter tuning

The next step was defining a set of fixed and variable hyperparameters and determining the range of tests to perform. Since the hardware is a limiting factor, firstly, the capping parameters needed to be determined. Among the set of parameters, the batch size is the most limiting in our case. Through experimentation, we determined that, for ResNet50 and MobileNet, the maximum batch size that allowed successful training was sixteen. For EffientDet, the value was equal to only one. Thus, we established that for ResNet50 and MobileNet, we would try three different values of batch size: 16, 8, and 4; while for EffientDet, there is only one value left: 1.

Another parameter that affects the model training the most is the learning rate. This value sits between 1 and 0. However, not all values in between will produce a positive effect on training. As discussed before, if the LR is too big, the model will diverge; if too small, the model would not learn. Once more, through some trials, we discovered the range of edge values for LR and defined that we will test 5 different values. In order to obtain them, we generated five different values on a logarithmic scale from 1 to 1e-6 and inspected them on a plot for an even distribution. The values can be seen in Table 5.2.

Table 5.2: Learning Rate values considered for the tests.

Index	LR Value
1	3.056e-1
2	2.499e-2
3	5.286e-4
4	1.807e-5
5	1.242e-6

For the optimizer, we went with the default selection on the configuration file of pre-trained models: Momentum. We took into consideration other optimizers, such as Adam, for the fact of being widely recommended among the literature. However, as we discovered later, with some tests where the only variable was the optimizer, it seemed only to affect, as expected, the learning curve, i.e., how early the model will converge. In spite of these findings, for this set of tests, optimizer and other parameters were considered less important for the performance of the final model, or, at least, less promising to help achieve better results.

Once parameters are defined, the test set will consist of a combination of them. Each selected model will be trained with all values of batch size and learning rate defined previously. Also, at this stage, the training process will be done on cross-validation folds. So, the model will be trained and validated on three folds for each parameters tests. In total, the test set will consist of 3(models)*3(folds)*3(batch

sizes)*5(LR values) = 135 tests. Once found the model and parameters that have the best performance on the actual data, they will be fixed for further tests.

5.2.4 Data manipulation

Semi-supervised training

As already indicated in section 4.2 of Chapter 4, there is an unannotated part of the data, which contains 98 samples. It was decided to use a semi-supervised approach toward the use of this part of the dataset in the present work. This approach will consist of the annotation of some unannotated samples by the previously developed model.

Once the model generates annotations, the next step will be to train the model on the train set and the newly annotated samples. In order to ensure the correctness of the comparison, the evaluation has to be made on the test set, which could not be submitted to any changes at this stage. Also, it was defined to use only 15 to 20 images from around one hundred available from each unannotated sample. This is due to the fact that in the annotated samples only 15 to 20 images per sample were used. This way, we will ensure that each sample will have the same amount of representation.

To help understand the technique described above, Figure 5.9 demonstrates visually how this new dataset will enter in the training process.

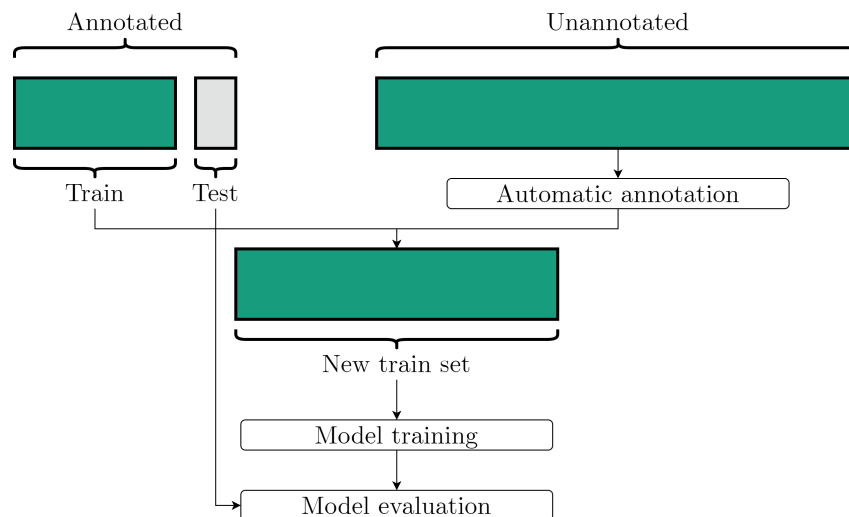


Figure 5.9: Semi-supervised pipeline to improve model performance by increasing training data.

Merging and eliminating classes

Another procedure that could help improve model performance, given the actual dataset distribution per class, is to merge some of the classes that appear to have

similar features and are underrepresented. This would decrease the number of classes and increase the number of annotations in the merged class. Also, underrepresented classes may as well be eliminated if found to be irrelevant for the purpose of the module. Eliminating some annotations can result in decreasing the noise factors for the model.

While the outcome of these procedures is unknown, there is an interest in performing these types of experiments to try to understand better the relevance of each class in the dataset — especially the relevance of the class of undefined nuclei.

Correction of chromatic aberrations

Having mentioned chromatic aberration issues, and being able to use recently developed tool in Fraunhofer to correct the chromatic aberrations, we used that tool to correct all images in our test set. Since this issue was not one of the issues that once solved would bring a huge performance gain, at least in our belief, when also compared to others previously mentioned strategies, we limited the test to only correcting the test set and using already trained model, on uncorrected data, and performing the evaluation.

5.3 Evaluation

Since one of the goals of this work is to set and improve baseline results of object detection of cervical nuclei and classification, it is required to use performance metrics in order to quantify these improvements. It is also important to use standard metrics that are commonly used in the field of object detection. This way, the model and the module performance can be compared to other available solutions. Therefore, the module was evaluated using the main metrics reported in the PASCAL Visual Object Classes (VOC) [67] and COCO [68] challenges, in particular, Average Precision (AP) and Recall.

During the development of the model, there are two main stages where evaluation is performed: during the training process; and, final evaluation of the system. The TensorFlow API allows the creation of checkpoints during the training. These checkpoints enable to resume training from the last checkpoint and also evaluate the model performance. So, during the training stage, the evaluation will be running in parallel in order to evaluate each checkpoint. It will help to prevent unnecessary training of the cases where the model diverges early or has very low performance, also in early stages. For visualizing the evaluation results, TensorFlow provides Tensorboard (see Figure 5.10). This tool helps in visualizing all logged metrics and data during the training, which makes it very convenient to monitor the training process. The main difference between two stages is that during the training process, the images that are evaluated are patches, while on the system evaluation, the whole image

is evaluated, or, even a batch of images of an sample. Also, final system evaluation allows getting a better understanding of the final performance since it accounts for variables that are out of the scope of the object detection model.

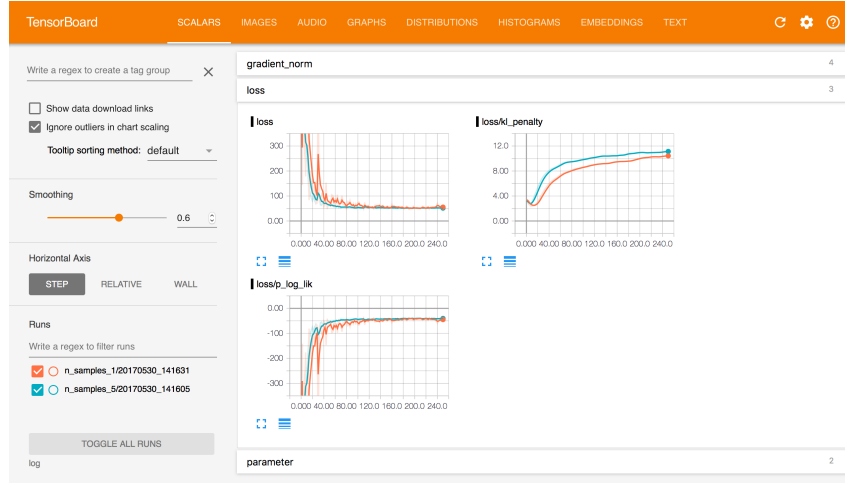


Figure 5.10: Tensorboard interface.

For the final system evaluation, the evaluation script has been developed. The script consists of the image slicing stage (with optional possibility of overlap), performing detection, image reconstruction, and comparing with ground truth annotations. In case of using overlapped patches, some objects were detected more than once due to the fact of being submitted to the model twice (or more times). To cope with duplicate detection we used Non-Maximum Suppression (NMS) algorithm, where detections with high intersection over union (IoU) percentage (threshold value further selected via experimentation) were eliminated until one detection were left for one object. Lastly, by comparing detections made by our model we were able to compare them to ground truth annotations and count the number of true positives (TP), false positives (FP), false negatives (FN) and with these values extracting others metrics such as AP, AR, Accuracy, Specificity, F1 index and Youden index.

5.3.1 Evaluation metrics

IoU

In object detection the intersection over union (IoU) is the most used metric, on which other metrics consequently will rely. IoU, or also known as Jaccard index, aims to reduce the proximity/overlapping of two arbitrary shapes to one value. This way, we can evaluate how similar are two geometric figures (bounding box or object mask) in terms of relative position and their overlap. The IoU can be calculated according to eq. 5.1, where $A, B \subseteq S \in R^n$. Establishing a certain IoU threshold when comparing detection and ground truth annotations allows, then, classify detections as TP, FP or FN.

$$IoU = \frac{|A \cap B|}{|A \cup B|} \quad (5.1)$$

The Figure 5.11 may help to visualize the concept of IoU.

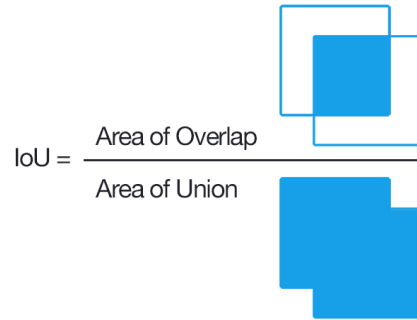


Figure 5.11: IoU visual explanation [5].

Precision

Precision, or confidence, is another fundamental metric which denotes the proportion of predicted positive cases that are correctly Real Positives, or in other form, the ratio of true positive cases over all predicted cases, as can be seen in eq. 5.2.

$$Precision = \frac{TP}{TP + FP} \quad (5.2)$$

There are some extensions to the simple Precision metric, such as **Average Precision (AP)** and **Mean Average Precision (mAP)**. The value for AP, considering Pascal VOC challenge, is computed by obtaining the area under the model's Precision-Recall curve with fixed IoU. In contrast, the COCO's mAP is computed by, firstly, calculating several AP values by averaging precision values over recall values. Each AP is associated with distinct IoU threshold - in this case the range is 0.5-0.95 with 0.05 step. Finally the mean of all AP values is computed, obtaining mAP.

Besides the differences in calculations, both try to evaluate the model's overall performance.

Recall

Recall or True Positive Rate (TPR), or also, Sensitivity, measures the proportion of positives cases that are correctly identified as such, see eq. 5.3.

$$Recall = \frac{TP}{TP + FN} \quad (5.3)$$

Accuracy

Accuracy evaluates how close or far off a given set of measurements are to their ground truth values. Or, in other terms is the number of correct predictions from all predictions made, as can be seen in eq. 5.4.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5.4)$$

Precision, in contrast to accuracy, and to help distinguish both metrics, tries to evaluate how close or disperse the measurements are to each other.

Specificity

Specificity or True Negative Rate (TNR) measures the proportion of negatives that are correctly identified as such (TN), see eq. 5.5

$$Specificity = \frac{TN}{TN + FN} \quad (5.5)$$

F1 Score

The extensive ensemble of tests (see Appendix E) required a more complex measuring metrics, being precision, average precision and recall unable to provide enough information about two models to be able to pick the best. For this we employed F1 Score, the harmonic average of the precision and recall, see eq. 5.6.

$$F_1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (5.6)$$

Youden Index

For the same reason as the F1 Score, we took a look into Youden Index. As can be seen in eq.5.7, derived from sensitivity and specificity, it denotes a linear correspondence balanced accuracy and evaluates the model's ability to avoid failure.

$$YoudenIndex = TPR + TNR - 1 \quad (5.7)$$

The Youden Index has been the main and last metric used in final system parameters search, when various tests regarding different overlap, IoU and NMS threshold values were performed.

5.3.2 Evaluation pipeline

As already mentioned, there are two main stages of system development: model development; and, final system development. During the **model training** process we employed default Object Detection API tools in order to train and evaluate the

model. Having created the .tfrecord files to ingest into model training script and into the evaluation script, we ran both in parallel, saving a training checkpoint every one thousand steps. Whenever a checkpoint was saved, the evaluation script performed an evaluation of the test set, making the predictions on them and comparing with ground truth annotations. Then, we observed the metrics graphs of mAP and Average Recall (AR) and aggregated different test runs in order to pick the best set of parameters used in the tests. Since each checkpoint represents the model at a certain number of training steps, we had to choose one checkpoint to compare with other test runs. For this we looked at each test run and picked the checkpoint that has best mAP and AR values.

For the **final system** evaluation, after the whole algorithm, presented in Figure 5.1, has been developed, we used the same test set as in the previous evaluation step, however, with the difference of using the whole images present in the sample and not only patches of these images.

Having implemented the algorithm of slicing, we added the option of allowing overlapping patches with user selected percentage of overlap. This also brought up the problem of duplicate detections in areas of overlap, where essentially the model will see the object twice. To overcome duplicates, we used NMS algorithm, which aims to eliminate the duplicates by calculating the IoU of two boxes. It has the IoU threshold that defines at which IoU value the two boxes can be assumed duplicates and one of them eliminated. Then the metric evaluation module was added in order to extract metrics of the test set. Here we used an open source toolkit available on GitHub that combines a handful of useful metrics [69]. We then, ran the tests on the whole system and, firstly, the goal of different tests was to pick best overlap percentage, NMS IoU threshold and metrics IoU that defines wheaten the detection corresponds to an object or not. Having fixed these values, we proceeded to setting the baseline of our full system and used it to evaluate performance of data manipulation approaches.

The final stage of the evaluation, when the best model and best parameters were already selected, consisted of running the inference on the test samples and estimate the average number of squamous nuclei per image of the sample. Then, we compared the value of the sample with earlier mentioned threshold set by The Bethesda System, of minimum 3.8 squamous nuclei per image for an adequate sample. Finally, we classified samples as adequate or inadequate, based on the previously described threshold and compared to the samples ground truth classification.

Chapter 6

Experiments and Results

This chapter aims to present results that were obtained through the course of training the model, and different experiments that were performed on the developed system. The chapter will start with hardware details in order to give some idea of the overall performance of the developed system. Then, results from the baseline parameters search will be presented, following with the best hyperparameters search and data-centric tests. Finally, the results of the final system are introduced, along the sample level tests, which will summarize all the developed work through the course of this dissertation.

6.1 Hardware details

Considering the state of the art frameworks and technologies that were used and their computational requirements, it is worth mentioning the hardware that allowed relatively fast development and short model training time.

As we already saw in the Chapter 3, Deep Learning models usually use a predefined set of functions, which are mostly matrix operations. It happens that graphics processing units (GPU), initially designed to accelerate computer graphics, which also and mostly includes matrix operations, are much more suitable for the model training than a central processing unit (CPU).

Having that said, Fraunhofer gave us access to various Virtual Desktops (VD), which were mainly virtual machines with different operational systems. Nevertheless, only two VD were used and their specifications can be seen in the Table 6.1,

present below:

Table 6.1: Virtual Desktops used for the development

	VD 1	VD 2
OS	Ubuntu 20.04 LTS	Ubuntu 20.04 LTS
CPU	4 Cores	4 Cores
GPU	None	1× Nvidia Tesla T4 (8 GB)
RAM	16 GB	16 GB
Storage	3 TB shared	3 TB shared
Networking	1 Gbps	1 Gbps

The VD 1, was mainly used for the tasks when GPU was not needed, as for example code development and debugging. The VD 2, which offered GPU, was used for model training and for final system fast evaluations. The Object Detection API model training pipeline intended to training the model for some steps, save checkpoint, evaluate the checkpoint and repeat the process for the specified number of steps. Unfortunately, we were not able to perform this only using GPU. Due to batch size we used and the fact that TensorFlow allocated almost all the available Video RAM (VRAM) and not freeing it up when the pipeline finished training, the evaluation process could not be started due to the lack of free video memory. To overcome this we had two separate pipelines, or processes, for training and evaluating, and while the training was done on GPU, the evaluation was performed using CPU. While much slower on CPU, the evaluation process, in most cases, was on time to process every checkpoint created by the training pipeline, and when was not catching up, we used VD 1 for evaluation, since the disk space was shared and both desktops had access to the same files. Even though the two VDs had only 4 virtual CPU cores, using VD 1 was a little faster to evaluate checkpoint since it was the main load for it, in contrast to VD 2, which also had CPU load due to the training process.

Since the Virtual Desktops were just virtual machines, to get better understanding of the performance of these VD, the following Table 6.2 shows the specification of the host machine:

Table 6.2: VD Host Server

Server	
OS	VMware Horizon
CPU	2× AMD EPYC™ 7302
GPU	4× Nvidia Tesla T4 (16 GB)
RAM	512 GB
Storage	On another server (vSAN)
Networking	10 Gbps

The AMD EPYC™ 7302 CPU offers 16 cores and 32 threads, and has a base clock of 3.0 GHz. Combining two of them, we can count a total of 64 threads. It is worth noting that virtualization software also creates virtual CPU cores, and the number of total active virtualized cores can be greater than the number of physical ones. While this approach gets the maximum of the hardware, it also brings the disadvantage of lower performance when comparing to using dedicated physical cores. Also, the performance can be inconsistent when too many virtual desktops are used at the same time. As per GPU, there are several Nvidia Tesla T4 with 16 GB of VRAM on board, which were also split between VDs with GPU. Essentially, the GPU were virtualized and presented for the Virtual Desktop as Tesla T4 with 8 GB of dedicated VRAM, but with shared GPU unit. Again, this is great for getting the most of the hardware, but came at the cost of reduced training speed when GPU was hit by two users simultaneously.

Another note is that while network connection to the server seemed to be a 10 Gbps connection, the connection to the actual working space, i.e. under the /home folder, was only up to 1 Gbps. Considering the fact that Virtual Desktop images were located on another server and that the model training and evaluation required a constant flow of images from the dataset and other intermediate files, the 1 Gbps connection, which in the best scenario equals to 125 MBps, in practice up to 80 MBps, was a little to slow and probably was the biggest bottleneck of this system.

By the end of this work, Nvidia Tesla T4 GPUs were joined by two new Nvidia A16, which offer 4 graphical processing units each, equal to $4 \times$ Nvidia A2, on the printed circuit board (PCB) and a total of 64 GB of VRAM, which allowed to use one A16 for 8 Virtual Desktops, instead of only two in the case of Tesla T4, while maintaining the performance.

6.2 Baseline parameters search

In order to compare model improvements, some baseline results had to be set. However, even with some parameters selected for hyperparameters experiments, there were still two variables that needed to be defined prior to the conduction of any experiments: Stop criteria and learning rate scheduler function.

Stop criteria refers to when to stop the model training. It happens that leaving the training process for too many steps can result in model overfitting. And, stopping the training too early will lead to underfitting. To determine the optimal stop criteria, several preliminary model trainings were conducted with different numbers of epochs: 10, 30, 50 and 100. It is important to recall that one epoch is equal to one pass through the model of all data present in the training set, while the pass of one batch of images is defined as a step. The configuration file of object detection API refers to the maximum number of steps to train, thus, epochs must be converted to

the number of steps. Fortunately, the conversion is an easy step and its equation can be seen in the Eq. 6.1.

$$Steps = Epochs \times \frac{images_in_training_set}{Batch_size} \quad (6.1)$$

Through these preliminary tests, it became clear that the performance of the model, namely metrics as loss, mAP and AR were getting better with higher number of epochs until 50. When comparing 50 and 100 epochs, the difference in performance was negligible, but when looking at the confidence of the predictions, they were almost 100% (exemplified in Figure 6.1) for the case with 100 epochs, which may imply the occurrence of overfitting. Given these results, it was decided to use number of epochs of 50 for the next tests. This value seems to be a good starting point to be able to evaluate others parameters.

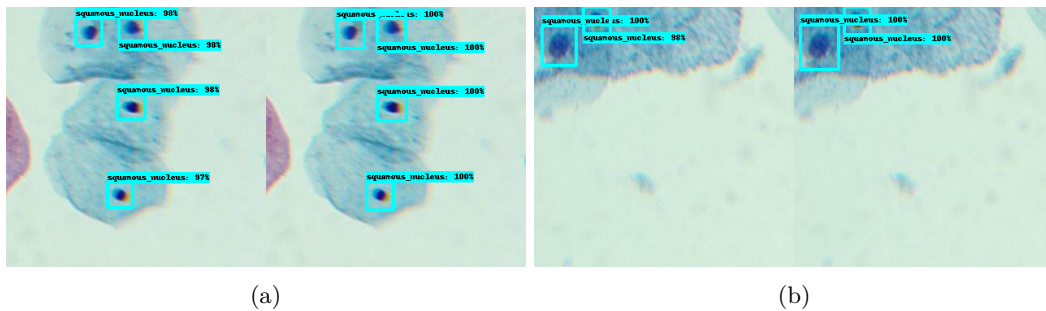


Figure 6.1: Examples of predictions in form of bounding boxes with respective classification and confidence (on the left side of each image) compared to the ground truth annotations (on the right side of each image).

Aside from stop criteria, the learning rate scheduler was another variable to be defined. The learning rate scheduler defines the learning rate value through the training process. It can be either one of the predefined functions or also completely custom or even manual. It happens that when the LR is maintained the same through all training process, the model quickly becomes unstable, and thus, it was decided to implement some sort of decay. As the model weights are modified with the influence of the learning rate through the training process, with each epoch, the magnitude of correction of weights becomes smaller and smaller as the system tries to fine-tune the features learning. In the case of maintaining the LR, this fine-tuning becomes impossible, and the model overshoots the optimal values of the weights.

For this reason, schedulers like Exponential decay, Linear decay, and constant LR were tested. Exponential decay, by default, offered a warm-up of the learning rate. The warm-up consisted of, instead of initializing the LR value with the predefined one, starting with a lot lower value and going to the predefined in some number of steps. It happens that it can be beneficial since, with higher LR values, the model can diverge on the early steps. However, with the warm-up function, the model has

some time with a lower LR value to smooth the initial learning. This way, it is possible to use higher LR without getting early divergence of the model.

The tests showed that there is a difference between constant value and scheduler with decay. In the case of the constant value, the model performed worse than with the decay, and while the first epochs showed almost no difference, at the last ones, the model with constant decay was still unstable, whereas the model with decay was converging. The tests with warm-up and without did not result in substantial differences, aside from being able to smooth a learning curve with higher initial LR values. However, since these variables were not so promising to improve the model quality, no further tests were conducted in regard to the scheduler or a maximum number of steps. As a result, the exponential decay scheduler without warm-up was selected as the default scheduler for the rest of the experiments.

These early tests were performed on the train subset from one of the cross-validation folds. For these tests, the number of steps was equal to 40000 steps. However, for further training on the train set, the number of steps will be equal to 61000 since the number of images in the training set, when compared to the training subset of folds, is bigger.

Another observation made on these tests was that the metric mAP seemed to be very low, even for the first tests. The preliminary tests resulted in mAP of 16-20% (see Figure 6.2). One of the reasons for that was the metric mAP itself, which averages mAP across all the classes. And, as was stated in the previous chapter about the disbalance of the dataset, the classes that have a low count of annotations may be responsible for having very low mAP, and therefore, lowering the average of mAP across all the classes. Switching to mAP per class indeed confirmed that classes with low numbers with annotations were only getting 0-2% of mAP, whereas inflammatory and squamous were ranging from 20-50%. More detailed results can be seen in Appendix B.

It is important to mention again that TBS specifies a minimum number of well preserved squamous nuclei, and therefore, for our work, only squamous nuclei are important. Nevertheless, it would be a bonus if our model could count other objects.

The preliminary tests also showed that the training process takes from 4 to 8 hours, depending on the model chosen.

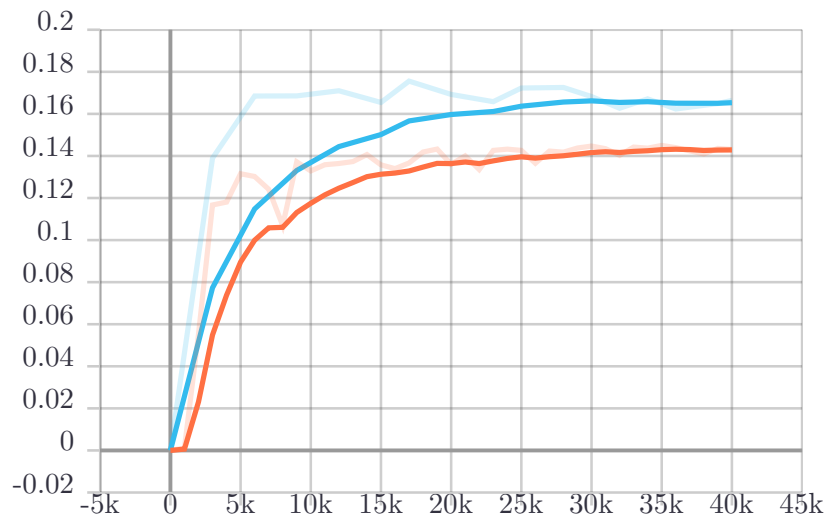


Figure 6.2: Overall mAP (on vertical axis) progression of two models: MobileNet (Orange) and ResNet50 (Blue); during the training on fold 1.

6.3 Hyperparameter tuning

The results, obtained on folds, of the tests are the following: For the model MobileNet, the highest value of LR, with index one on table 5.2, resulted in divergence for all three folds and all three batch sizes. Stepping down to the lower LR value resulted in almost the same behavior, except for the batch size of 16 that performed well, achieving 49.8% mAP on the class of squamous nuclei and 43.3% on the class of inflammatory cells, averaged by three folds. The other two LR values, with index 3 and 4, revealed even lower metrics, and the last one, and the smallest one resulted in, lack of model learning at all, achieving 0.3% mAP on the best run for squamous nucleus.

The variation of the batch size intended to double the number of steps when using half of the previous batch size, which was taken into account. This way, the model was trained for the same number of epochs for all three values of batch size. From the tests, it became clear that the model training was slower with the smaller batch size, i.e., it took more epochs to reach the same performance. A smaller batch size also resulted in slightly lower metrics, accordingly to the Figure 6.3(a).

That being said, the best hyperparameters for the MobileNet on the cross-fold validation were: Batch size of 16 and LR of 0.02499.

As per ResNet50, the highest LR value also resulted in the divergence of the model. The results were almost identical for the rest of the LR values: As LR went lower, so did the mAP and average recall. However, the second LR value performed worse for the ResNet50 when compared to MobileNet, resulting in the best LR value for the ResNet50 being 0.000528 with average mAP along three folds of 50,5% on

the squamous nucleus and 44.4% on inflammatory cells.

The batch size variations produced the same effect as in the MobileNet: the lower the batch size, the lower are the metrics. Thus, the best batch size was also 16.

The last model, EfficientDet D0, was more challenging since the resizing of the model input size to 320 by 320 pixels did not work through the configuration file. The smallest input size that we were able to configure was 512 by 512. The search on how to resize the model input gave the impression that it required more effort to make it work. In addition to the fact that the EfficientDet D0 net also did not have a straight pipeline for using it on the mobile platform, at least compared to MobileNet and ResNet50, it was decided to try the model as it is with input data with 320 by 320 pixels in dimension and adding black padding to the area around the patch. On most of the LR values, the model diverged, only being able to converge with the LR value of 0.0005286. Also, It was only possible to train the model with a batch size of 1 with available hardware. The one test that we were able to conduct resulted in rather poor results when compared to MobileNet and ResNet50: achieving only 10% and 7% of mAP on squamous nucleus and inflammatory cell classes. Appendix C presents Tables C.1, C.2 and C.3 with detailed results of the hyperparameters tuning tests.

Training and evaluating then a smaller set of hyperparameters, on the train and test sets respectively, as depicted in Figure 6.3, we obtained the following results: Even though the MobileNet achieved slightly better results than ResNet50 on cross-fold validations, on the test set the latter was found to be better one, achieving the following performance metrics: squamous nucleus AP@50: 47.3%; inflammatory cell AP@50: 37.8%; AR@100: 43%, which, once again, can be visualised in the Table C.3. It is important to note that the mAP performance can not be compared between tests set and fold validation subsets since the set of the images for the tests is different.

Given the demarked imbalance in the dataset regarding the number of annotations of squamous nucleus and inflammatory cell versus the remaining classes, the poor results regarding these classes were already expected. Since the model appears only to be capable to properly generalize to these two highly represented classes, and being the squamous nuclei the key class for adequacy assessment, most of the following results and respective discussion are focused on this class.

In order to reduce the number of variables in the further tests, and therefore, reduce the time required for each experiment, it was decided to choose one model and parameters that led to the best performance. That being said, it was decided to go further with the ResNet50 with the following hyperparameters: Batch size of 16 and LR of 5.29×10^{-4} .

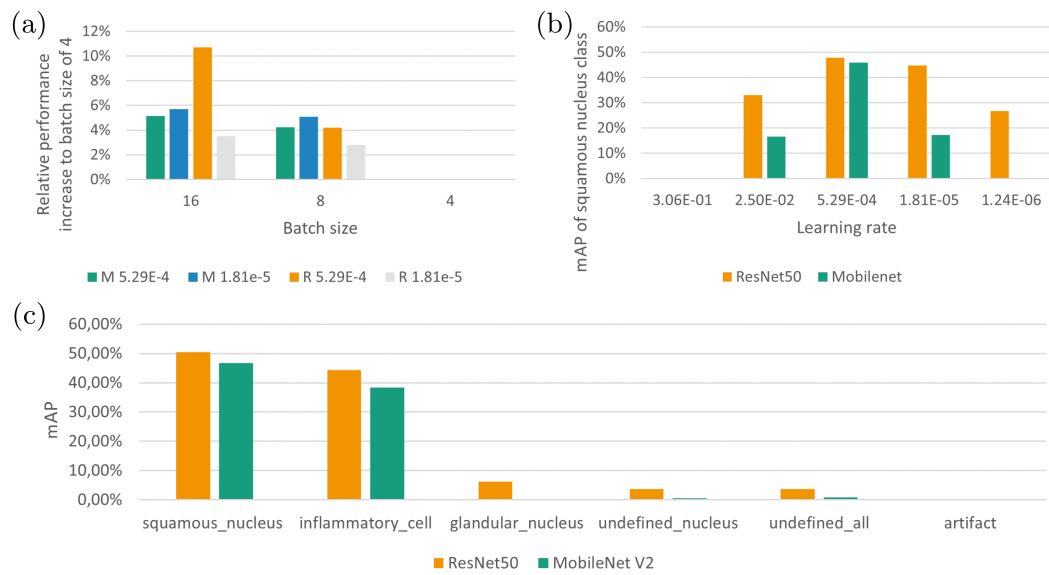


Figure 6.3: Results of hyperparameters tuning evaluated on the test set. Image (a) outlines relative performance increase of bigger batch size values when compared to batch size of 4. Image (b) depicts mAP of squamous nucleus class across different values of LR with batch size of 16. Image (c) shows mAP across different classes with LR = 5.29×10^{-4} and batch size = 16.

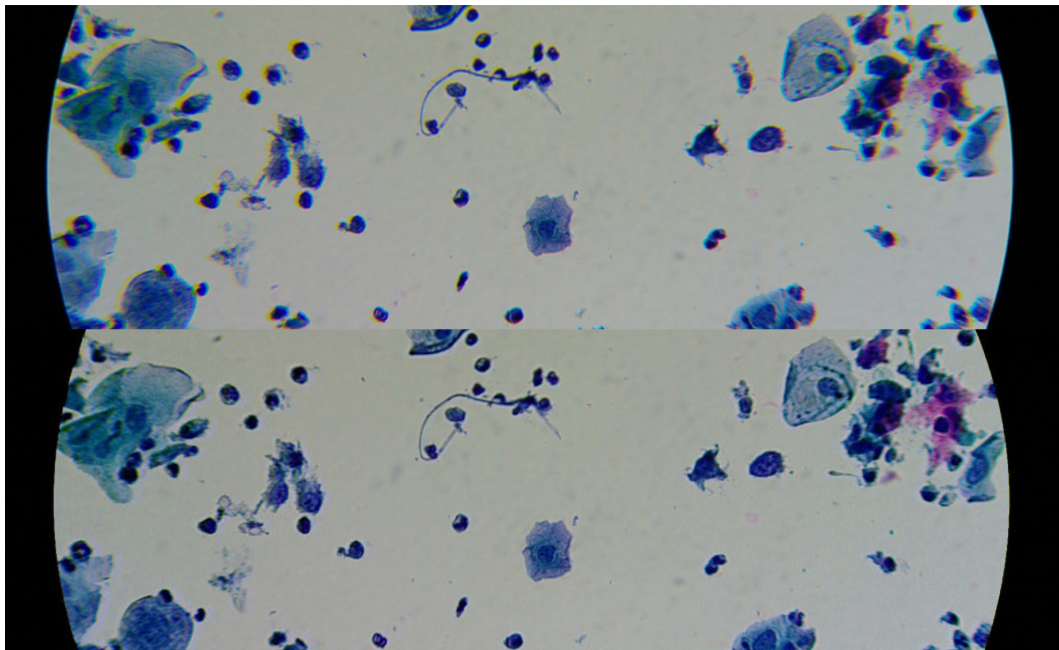
6.4 Chromatic aberration correction test

One of the internal Fraunhofer's projects addressed the chromatic aberration problem encountered on the images of the dataset. Visually, the developed algorithm to eliminate the chromatic aberrations revealed promising results, see Figure 6.4. We found that it would be interesting to compare how the model will behave when tested on the images without chromatic aberrations.

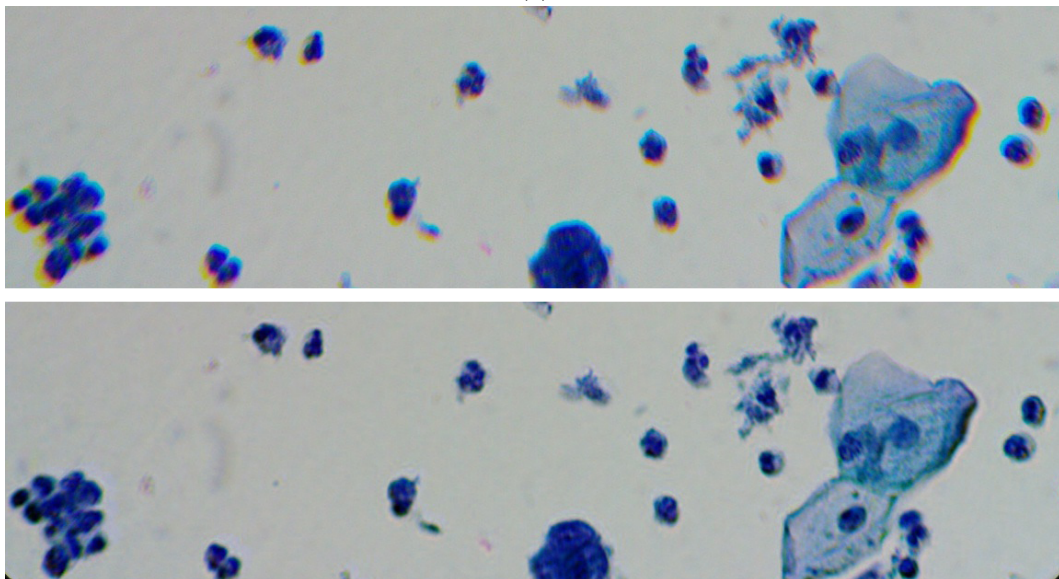
Unfortunately, the results of the test revealed slightly worse performance when compared to images with chromatic aberrations. This is not very surprising since the model used for this test was trained on images with chromatic aberration and probably accounts for the features of nuclei with blue and orange fringing.

Also, in a meeting with medical personnel, we had the opportunity to show the corrected images and get feedback. They informed that corrected images appeared to be more clean looking, and on some occasions, the classification of the nuclei was easier.

Since chromatic aberration is an unwanted transforming element to the image, which results in slightly distorted information, or in this particular case, nuclei, then the correction, which eliminates this transformation and brings the information to its initial state, should be accounted in the image post-processing after its acquisition. However, as further work, it would be interesting to train the model on the corrected images and assess the impact in the detection performance.



(a)



(b)

Figure 6.4: Example of chromatic aberration correction. Upper part of 6.4a and 6.4b is the original image and bottom part is the same image after correction. The images were cropped to its actual size to facilitate observation. As can be noticed, correction did an excellent job in eliminating blue and orange fringing around nuclei and other objects.

6.5 Data-centric experiments

6.5.1 Merging and eliminating classes

As seen from previous experiments, only two classes: squamous and inflammatory, are capable of getting any meaningful predictions. The other four classes, unfortunately, can not show this level of performance.

Probably the two classes that can exhibit a bigger level of concern are the undefined nucleus and undefined all. It was defined, at the annotation stage, that any object in these two classes can belong to squamous, inflammatory or glandular classes. So, if a model classifies one of the undefined objects as squamous, inflammatory, or glandular, it should be counted as true positive. However, due to the way the evaluation process is conducted, the case described above would be a misclassification.

To try to solve this issue, we contemplated several solutions, one of which consisted of duplicating undefined annotations. For example, instead of having just one undefined nucleus, create two annotations on the same spot with squamous and glandular classification. However, it would not solve the problem because, in the case of correct prediction of one of the annotations, another will always be classified as false positive. Therefore, two other approaches were subjects of a study:

- Eliminate all classes but squamous and inflammatory. This approach was motivated by the hypothesis that, because of the visual similarities between squamous, inflammatory and two undefined classes, the latter were introducing noise, and instead of helping distinguish different classes, only helped to confuse the model. For the other classes, glandular and artifact, two more tests were conducted, with and without them.
- Another approach consisted in, instead of eliminating similar classes, merging the similar ones, such as two undefined, or even merging all underrepresented, i.e., all but squamous nucleus and inflammatory cell classes, and calling them as "other" class.

The results we obtained (see Table D.1 in Appendix D) with these tests were around the previously set baseline or worse, and, considering that all classes aside from two main ones, account only for 5.9% of all annotations, which presence or lack also does not contribute directly to the main ones, and then, probably it was expected to have negligible impact. Interestingly enough, the test with only one class of the squamous nuclei performed noticeably worse than the test where we had two classes (squamous nucleus and inflammatory cells) or more. It is also worth mentioning that the test set was left untouched and the only metrics assessed here were metrics strictly related to the squamous and inflammatory classes. Having that said, we decided to keep all classes as they were originally.

6.5.2 Semi-supervised approach

Although there were more 99 samples available for Fraunhofer, we found out that only 30 of them were already digitized, that is, suitable for us to extend the annotations. These new samples had around 100 images each and to preserve homogeneous sample representations we decided to take only 19 randomly selected images from them and annotate with the last and best available model, which was still the model trained and selected on hyperparameter tuning stage. The precise number of 19 images per sample was used due to the fact that in the original dataset the average number of images that each annotated sample has is 19.

Having automatically annotated these new 30 samples, we were able to increase our train set by 570 images, totaling 1335 images, and, as can be seen in the Table 6.3, increase by 66% the number of annotations in the train set.

Table 6.3: Dataset distribution by nucleus type.

Class name	Train set	+30 Samples	Change	% Increase
Squamous nucleus	15543	29124	13581	87.38%
Inflammatory cell	14577	22292	7715	52.93%
Glandular nucleus	538	540	2	0.37%
Undefined nucleus	491	492	1	0.20%
Undefined all	831	832	1	0.12%
Artifact	123	123	0	0.00%
Total	32103	53403	21300	66.35%

One more time, the poor performance of underrepresented classes can be noticed here as there are almost no predictions of them.

Using the best hyperparameters obtained in the hyperparameter tuning step, we trained the model with this new train set and evaluated once more on the test set. Surprisingly, comparing to the best obtained result, so far, on the final system (can be seen in the next section) we noticed a bit better performance on inflammatory cell class, but at the cost of a little worse performance on squamous nucleus class.

For now, it is difficult to say the reason of why the model did not increased overall performance, but only traded, a bit, one class for another. But it could happen that the quality of original annotations was not satisfactory, and therefore, the quality of automatic annotations was even worse, which may justify the outcome of this test. The results of this test, obtained with the final system, can be observed in the Table 6.4 and Figure 6.5, and compared with the final image level results in the next section, in particular with the Table 6.5 and Figure 6.6.

Considering the importance of squamous nucleus class, further tests on this direction were discarded.

Table 6.4: Performance metrics of the model after semi-supervised training, taken on the final system.

Class	Squam.	Inflam.	Gland.	Und. n.	Und. all	Artifact
AP	0.7948	0.7060	0	0.0	0.0	0.0
Recall	0.6660	0.7096	0.0	0.0	0.0	0.0
Accuracy	0.7712	0.8198	0.9868	0.9842	0.985	0.997
Specificity	0.9214	0.8681	1.0	1.0	1.0	1.0
F1	0.7739	0.7060	0.0	0.0	0.0	0.0
Y. index	0.5873	0.5777	0.0	0.0	0.0	0.0
TP	4280	2365	0	0	0	0
FP	354	1002	0	0	0	0
FN	2147	968	1	173	164	29

Ground truth	Artifact	0	0	4	22	3	0	0
	Glandular n.	0	0	5	106	33	0	0
	Inflammatory c.	0	0	2365	888	80	0	0
	Nothing	0	0	481	0	180	0	0
	Squamous n.	0	0	455	1692	4280	0	0
	Und. n.	0	0	33	97	34	0	0
	Und. a.	0	0	24	125	24	0	0
		Artifact	Glandular n.	Inflammatory c.	Nothing	Squamous n.	Und. n.	Und. a.
		Detections						

Figure 6.5: Confusion Matrix of the model after semi-supervised training.

6.6 Final script evaluation

6.6.1 Final system parametrization

First tests and visual inspection of final model predictions revealed that some correct predictions were reported as false positives due to low intersection area and an IoU threshold set to 50%. Through testing (see Appendix E), we discovered that a threshold of 10% enables correct reporting of true positives without compromising performance and increasing misclassifications.

We also tested the overlap of sliced patches and a threshold for the NMS algorithm. It became apparent that an increasing overlap produced more predictions. By visually analyzing these results, we once again realized that a significant number of false positives detected by our model had very similar image characteristics to several annotated objects being detected as true positives. And as already said,

in a meeting with an experienced medical expert we discussed these results, and it was clarified that some of the predictions made by the system were indeed correct, meaning that the real performance is in fact higher than the reported. As per the NMS threshold, we discovered that a value of 97% of overlap or more was the best fit. It enabled the correct elimination of prediction duplicates while allowing the existence of prediction for two objects next to each other.

In Appendix E, especially in the Table E.5, which is ordered by decreasing Youden Index, we can see that the best set of overlap, NMS threshold and metric IoU for the ResNet50 is 29%, 97% and 10% correspondingly, when considering Youden Index and F1 Score metrics as key metrics.

Image level results

At the image level, the system reported an AP of 82.4%, Accuracy of 79.8%, and F1 score of 81.5% for the class of squamous nucleus. In terms of raw detection, the system made a total of 5216 predictions for 5483 existing annotations with a true positive rate of 74% and 473 false positives, which can be explored with more details in Table 6.5 and confusion matrix present in Fig. 6.6(c). Given the visual results shown in Fig. 6.6(b) and the previous conclusions concerning unannotated objects, these values are however expected to be higher, as some of them could be considered true positives if the dataset was correctly annotated.

It is worth mentioning that better results in terms of the recall could be achieved with a higher overlap value. However, aside from increasing processing requirements or time for assessment, the false positives rate will also increase. Thus, to correctly evaluate if overlap increasing brings performance benefits, we consider that the required dataset ground-truth quality improvement previously discussed should be handled first, to avoid tuning this parameter based on incorrectly unannotated nuclei.

Table 6.5: Performance metrics of the final system.

Class	Squam.	Inflam.	Gland.	Und. n.	Und. all	Artifact
AP	0.8236	0.7072	0.2257	0.0	0.0	0.0
Recall	0.7380	0.6274	0.0347	0.0	0.0	0.0
Accuracy	0.7984	0.8310	0.9864	0.9837	0.9847	0.9973
Specificity	0.8893	0.9232	0.9993	0.9999	1.0	1.0
F1	0.8147	0.6982	0.0641	0.0	0.0	0.0
Y. index	0.6273	0.5505	0.0341	-0.0001	0.0	0.0
TP	4743	2091	5	0	0	0
FP	473	566	7	1	0	0
FN	1684	1242	139	173	164	29

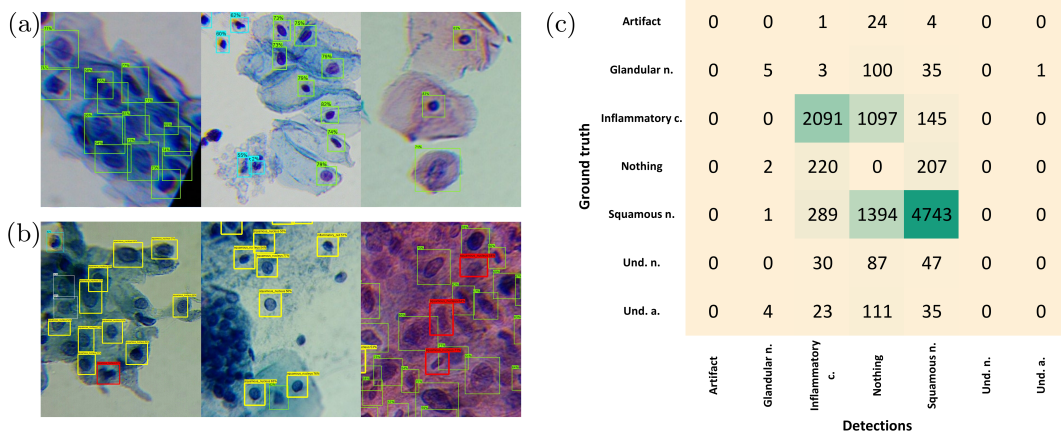
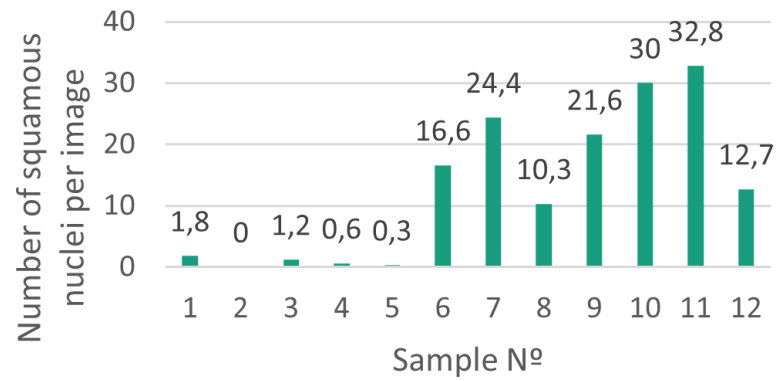


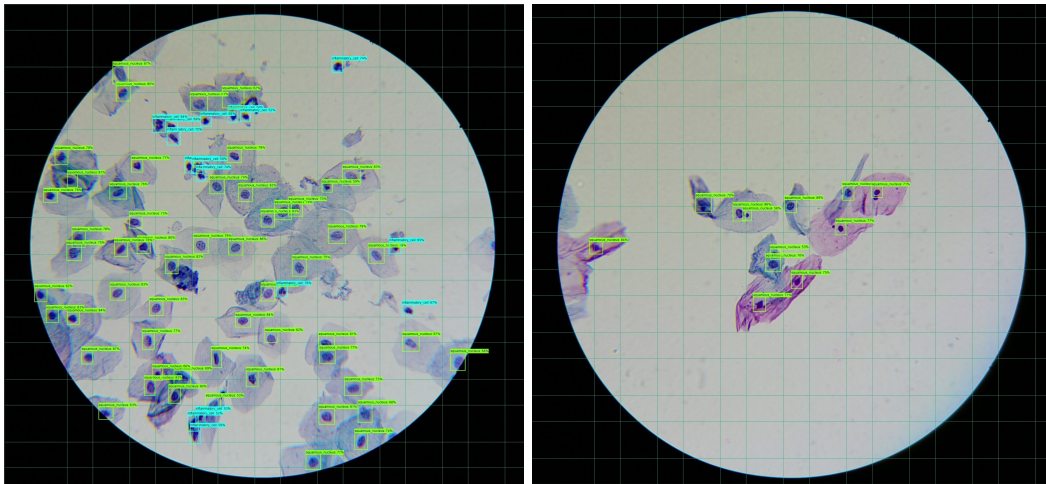
Figure 6.6: Detections and classifications made by the final system (a) and (b). Image (a) depicts mostly correct detection and classifications (green and blue bounding boxes means correct classifications of two different classes), while on the image (b) yellow and red bounding boxes correspond to misclassification and false positives detections respectively. Image (c) shows the confusion matrix of the test set. "Nothing" corresponds to false positives and false negatives.

Sample level results:

After the evaluation of the proposed approach at the image level, a study to assess sample adequacy at the sample level was conducted. In particular, a set of 12 samples (5 inadequate and 7 adequate) was used for that purpose, with each sample consisting of around 100 images of different microscopic fields. The final model was then used to count the number of squamous nuclei on each image, and the average number of squamous nuclei per sample was computed. The results are depicted in Fig. 6.7a, where it is visible that we were able to classify all samples correctly as adequate or inadequate, by using the threshold of 3.8 squamous nuclei per image stated by the The Bethesda System guidelines. Images present on Figures 6.7b and 6.7c are illustrative images of one of the outputs of the system. The grid present in the images exhibits patches boundaries.



(a)



(b)

(c)

Figure 6.7: Results for the proposed solution: (a) average number of detected squamous nuclei per sample (sample 1 to 5 are inadequate and 6 to 12 adequate); (b) and (c) illustrative images of automated nuclei detection at the image level.

Chapter 7

Conclusion and Future Work

The growing case rate of cervical cancer case rate and shortage of trained medical personnel and equipment in low-to-middle-income countries to perform screening and therefore prevent worsening of the disease stage have been one of the main factors that motivated this work. Moreover, the lack of research and already available cheap and mobile solutions for CADx systems for cervical cancer screening that complies with TBS specifications, was another crucial motivation for this dissertation.

Cytology screening comprises several steps, from obtaining the sample, passing through sample adequacy assessment and ending in the examination step where the search for malignant cells is performed. It is important that the sample is adequate before entering the last stage, otherwise, the chances of finding possible abnormal lesions are greatly compromised. While there are some work done on automated examinations of cervical cytology samples, almost none is accounting for sample adequacy.

Therefore, the dissertation described in the present document is centered around developing an automated adequacy assessment of cervical cytology samples. Considering the state of the art of object detection algorithms, another key aspect of the work is using deep learning models to perform the detection and classification of different types of nuclei and objects present in the cytology samples. The proposed work consists on using a recently created image dataset of cervical cytology samples, containing images acquired with μ SmartScope device, to train the model and performing several optimisation experiments in order to obtain the best performance.

The dataset incorporates 139 samples, 41 of which are annotated by specialists in terms of nuclei type and feature more than forty thousand annotations.

The developed system is composed of three main stages: pre-processing, object detection and classification, and post processing. The need of significant pre and post processing stages arose due to the decision to process images in patches of smaller sizes than the original image. Thus, the pre-processing stage is responsible for slicing the image in patches and post-processing reconstructing the image. The object detection and classification stage utilizes an CNN-based object detection DL model: SSD ResNet50, which was extensively trained and fine-tuned, over more than one hundred unique train and test runs. The hyperparameter tuning included batch size and learning rate optimizations, as well as preliminary search of best stop criteria, learning rate scheduler and the optimizer. After selecting best parameters using cross validation technique, the final system parametrization, mainly on the post-processing stage was performed, selecting the best overlapping percentage and the thresholds for the NMS algorithm to eliminate detection duplicates.

Lastly, the data centered approach was engaged and a few more experiments were performed with merged and eliminated underrepresented classes. A set of experiments were also performed using chromatic aberration corrected images. And finally, a semi-supervised training was employed, which used 30 additional sample without manual annotations. These images were then automatically annotated by the best model version found so far, and added to the original train set.

After training the model, finding best hyperparameters, experimenting with data-centric approaches, the best obtained object detection model showed, as expected, low performance on underrepresented classes, but also yielded acceptable results for two main classes: squamous nucleus and inflammatory cell. By using patches with overlap and NMS algorithm, we were able to increase performance, mainly by allowing the model to see some objects more than once.

In the last evaluation, performed on the sample level, we were able to classify all 12 samples correctly as adequate and inadequate.

Despite the apparent success of the proposed approach, several shortcomings of the current version of the dataset were identified and should be addressed in future works. Improving annotations quality, increasing dataset size, or balancing underrepresented classes will undoubtedly foster performance gains, both on automatic nuclei detection and adequacy assessment.

As future work it would be interesting to train the model on chromatic aberration corrected images and explore other object detection models. Furthermore, additional model-centric techniques such as fine-tuning online data augmentation or anchor box optimization could also be explored. However, the importance of quality data to this task cannot be overstated. We believe that innovative data-centric approaches are

still holding some potential for improvements and getting a better quality dataset could increase the performance.

Regarding the future deployment and usage of the proposed solution in a clinical context, further improvements and integration efforts are obviously required. The mobile integration with according optimizations are steps to evaluate the performance in real world environment. Nevertheless, it becomes clear that this work corresponds to a significant step forward in the development of innovative solutions that bring real efficiency gains to current cervical cancer screening procedures.

As an additional outcome of this dissertation, a paper has been published and an invited extended version to be submitted to a journal is under preparation. The first covers most of the work described in this dissertation while the second aims to provide extended results and to cover other aspects and experiments, including transfer learning to the cervical lesions detection and classification in microscopic images:

Automated Adequacy Assessment of Cervical Cytology Samples Using Deep Learning has been published in a conference: Mosiichuk, V., Viana, P., Oliveira, T., Rosado, L. (2022). Automated Adequacy Assessment of Cervical Cytology Samples Using Deep Learning. In: Pinho, A.J., Georgieva, P., Teixeira, L.F., Sánchez, J.A. (eds) Pattern Recognition and Image Analysis. IbPRIA 2022. Lecture Notes in Computer Science, vol 13256. Springer, Cham. https://doi.org/10.1007/978-3-031-04881-4_13.

Mosiichuk, V., et al, Cervical Lesions Detection and Classification in Microscopic Images Using Deep Learning and Transfer Learning. Pattern Analysis and Applications (PAA) journal.

References

- [1] R. Nayar and D. C. Wilbur, *The bethesda system for reporting cervical cytology: Definitions, criteria, and explanatory notes*. Springer International Publishing, 1 2015. [Cited on pages ix, 8, 9, 10, 11, and 12]
- [2] H. Kinsley and D. Kukiela, *Neural Networks from Scratch in Python: Building Neural Networks in Raw Python*. Harrison Kinsley, 2020. [Cited on pages ix, 16, 17, and 18]
- [3] A. H. Yazdani Abyaneh, A. Hossein Gharari, and V. Pourahmadi, “Deep neural networks meet csi-based authentication,” 11 2018. [Cited on pages ix and 20]
- [4] L. Rosado., J. Oliveira., M. João M. Vasconcelos., J. M. Correia da Costa., D. Elias., and J. S. Cardoso., “µsmartscope: 3d-printed smartphone microscope with motorized automated stage,” in *Proceedings of the 10th International Joint Conference on Biomedical Engineering Systems and Technologies - BIODEVICES, (BIOSTEC 2017)*, pp. 38–48, INSTICC, SciTePress, 2017. [Cited on pages ix and 26]
- [5] Adrian Rosebrock, Pyimagesearch, “Intersection over Union (IoU) for object detection .” Available at <https://pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/>. (Last accessed in 28/03/2021). [Cited on pages x and 51]
- [6] WHO, World Health Organization, “Cancer today.” Available at <https://gco.iarc.fr/today/fact-sheets-cancers>, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [7] WHO, World Health Organization, “Cervix uteri.” Available at <https://gco.iarc.fr/today/data/factsheets/cancers/23-Cervix-uteri-fact-sheet.pdf>, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [8] H. Sung, J. Ferlay, R. L. Siegel, M. Laversanne, I. Soerjomataram, A. Jemal, and F. Bray, “Global cancer statistics 2020: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries,” *CA: A Cancer Journal for Clinicians*, vol. 71, no. 3, pp. 209–249, 2021. [Cited on page 1]
- [9] WHO, World Health Organization, “Low HDI.” Available at <https://gco.iarc.fr/today/data/factsheets/populations/>

- 984-low-hdi-fact-sheets.pdf, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [10] WHO, World Health Organization, “Medium HDI.” Available at <https://gco.iarc.fr/today/data/factsheets/populations/983-medium-hdi-fact-sheets.pdf>, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [11] WHO, World Health Organization, “High HDI.” Available at <https://gco.iarc.fr/today/data/factsheets/populations/982-high-hdi-fact-sheets.pdf>, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [12] WHO, World Health Organization, “Very high HDI.” Available at <https://gco.iarc.fr/today/data/factsheets/populations/981-very-high-hdi-fact-sheets.pdf>, 2021. (Last accessed in 11/02/2021). [Cited on page 1]
- [13] C. Marth, F. Landoni, S. Mahner, M. McCormack, A. Gonzalez-Martin, and N. Colombo, “Cervical cancer: Esmo clinical practice guidelines for diagnosis, treatment and follow-up,” *Annals of Oncology*, vol. 28, pp. iv72–iv83, July 2017. [Cited on pages 1, 6, 7, and 8]
- [14] M. M. Rahaman, C. Li, X. Wu, Y. Yao, Z. Hu, T. Jiang, X. Li, and S. Qi, “A survey for cervical cytopathology image analysis using deep learning,” *IEEE Access*, vol. 8, pp. 61687–61710, 2020. [Cited on page 1]
- [15] L. J. Mango, “Computer-assisted cervical cancer screening using neural networks,” *Cancer Lett.*, vol. 77, pp. 155–162, Mar. 1994. [Cited on page 2]
- [16] P. Sanyal, S. Barui, P. Deb, and H. C. Sharma, “Performance of a convolutional neural network in screening liquid based cervical cytology smears,” *J. Cytol.*, vol. 36, pp. 146–151, July 2019. [Cited on page 2]
- [17] Q. Miao, J. Derbas, A. Eid, H. Subramanian, and V. Backman, “Automated cell selection using support vector machine for application to spectral nanocytology,” *BioMed Research International*, vol. 2016, pp. 1–10, 01 2016. [Cited on page 2]
- [18] T. Conceição, C. Braga, L. Rosado, and M. J. M. Vasconcelos, “A review of computational methods for cervical cells segmentation and abnormality classification,” *International Journal of Molecular Sciences*, vol. 20, p. 5114, Oct. 2019. [Cited on pages 2 and 7]

- [19] Centers for Disease Control and Prevention, “Basic Information About Cervical Cancer.” Available at https://www.cdc.gov/cancer/cervical/basic_info/index.htm, Jan 2021. (Last accessed in 11/02/2021). [Cited on pages 5 and 7]
- [20] WHO, World Health Organization, “Human papillomavirus (HPV) and cervical cancer.” Available at [https://www.who.int/news-room/fact-sheets/detail/human-papillomavirus-\(hpv\)-and-cervical-cancer](https://www.who.int/news-room/fact-sheets/detail/human-papillomavirus-(hpv)-and-cervical-cancer), Nov 2020. (Last accessed in 12/02/2021). [Cited on pages 6 and 7]
- [21] WHO, World Health Organization, ed., *Comprehensive cervical cancer control. Second edition: A guide to essential practice*. Genève, Switzerland: World Health Organization, 2014. [Cited on page 6]
- [22] Cancer Research UK, “Types and grades | Cervical cancer | Cancer Research UK.” Available at <https://www.cancerresearchuk.org/about-cancer/cervical-cancer/stages-types-grades/types-and-grades#>, Jan 2020. (Last accessed in 12/02/2021). [Cited on page 6]
- [23] National Cancer Institute, “HPV and Pap Testing.” Available at <https://www.cancer.gov/types/cervical/pap-hpv-testing-fact-sheet>, Dec 2019. (Last accessed in 15/02/2021). [Cited on page 7]
- [24] U. Banik, P. Bhattacharjee, S. U. Ahamad, and Z. Rahman, “Pattern of epithelial cell abnormality in pap smear: A clinicopathological and demographic correlation,” *CytoJournal*, vol. 8, p. 8, 4 2011. [Cited on page 7]
- [25] F. Haghghi, N. Ghanbarzadeh, M. Ataee, G. Sharifzadeh, J. Mojarrad, and F. Najafi-Semnani, “A comparison of liquid-based cytology with conventional papanicolaou smears in cervical dysplasia diagnosis,” *Advanced Biomedical Research*, vol. 5, p. 162, 2016. [Cited on page 8]
- [26] S. Qureshi, U. Singh, S. Foruin, N. Negi, G. Agarwal, and N. Singh, “Liquid-based cytology vs conventional cytology as a screening tool for cervical cancer in postmenopausal women,” *Journal of SAFOG*, vol. 9, pp. 159–163, 2017. [Cited on page 8]
- [27] K. K. Aboobacker and M. H. Shariff, “A comparative study of conventional pap smear with liquid based cytology for early diagnosis of cervical cancer,” *IP Archives of Cytology and Histopathology Research*, vol. 5, pp. 141–146, 6 2020. [Cited on page 8]
- [28] C. V. Biscotti, A. E. Dawson, B. Dziura, L. Galup, T. Darragh, A. Rahemtulla, and L. Wills-Frank, “Assisted primary screening using the automated thinprep imaging system,” *American Journal of Clinical Pathology*, vol. 123, pp. 281–287, 2 2005. [Cited on page 8]

- [29] H. Nishio, T. Iwata, H. Nomura, T. Morisada, N. Takeshima, H. Takano, H. Sasaki, E. Nakatani, S. Teramukai, and D. Aoki, “Liquid-based cytology versus conventional cytology for detection of uterine cervical lesions: A prospective observational study,” *Japanese Journal of Clinical Oncology*, vol. 48, pp. 522–528, 6 2018. [Cited on page 8]
- [30] S. Pankaj, S. Nazneen, S. Kumari, A. Kumari, A. Kumari, J. Kumari, V. Choudhary, and S. Kumar, “Comparison of conventional pap smear and liquid-based cytology: A study of cervical cancer screening at a tertiary care center in bihar,” *Indian Journal of Cancer*, vol. 55, p. 80, 1 2018. [Cited on page 8]
- [31] V. G. Padubidri and S. Daftary, eds., *Shaw’s Textbook of Gynecology*. India: Elsevier, 16th ed., 2014. [Cited on page 8]
- [32] A. Chrysostomou, D. Stylianou, A. Constantinidou, and L. Kostrikis, “Cervical cancer screening programs in europe: The transition towards hpv vaccination and population-based hpv testing,” *Viruses*, vol. 10, p. 729, 12 2018. [Cited on page 8]
- [33] R. Nayar and D. C. Wilbur, “The pap test and bethesda 2014,” *Acta Cytologica*, vol. 59, pp. 121–132, 6 2015. [Cited on page 8]
- [34] C. M. Bishop, *Pattern recognition and machine learning*. Springer New York, NY, 2006. [Cited on page 14]
- [35] H. Li and K. Kang, “Notes and assignments for cuhk deep learning course eleg5491: Introduction to deep learning.” <https://github.com/eleg5491/eleg5491.github.io>, 2017. [Cited on page 15]
- [36] J. Patterson and A. Gibson, *Deep Learning: A Practitioner’s Approach*. O’Reilly Media, 2017. [Cited on page 19]
- [37] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, pp. 436–44, 05 2015. [Cited on page 19]
- [38] F. Rosenblatt, *The Perceptron, a Perceiving and Recognizing Automaton Project Para*. Report: Cornell Aeronautical Laboratory, Cornell Aeronautical Laboratory, 1957. [Cited on page 19]
- [39] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, pp. 84 – 90, 2012. [Cited on page 19]

- [40] A. Ghosh, A. Sufian, F. Sultana, A. Chakrabarti, and D. De, *Fundamental Concepts of Convolutional Neural Network*, pp. 519–567. 01 2020. [Cited on page 22]
- [41] K. O’Shea and R. Nash, “An introduction to convolutional neural networks,” *ArXiv e-prints*, 11 2015. [Cited on page 22]
- [42] T. Falk, D. Mai, R. Bensch, d. Çiçek, A. Abdulkadir, Y. Marrakchi, A. Böhm, J. Deubner, Z. Jaeckel, K. Seiwald, O. Dovzhenko, O. Tietz, C. Dal Bosco, S. Walsh, D. Saltukoglu, T. Tay, M. Prinz, K. Palme, M. Simons, and O. Ronneberger, “U-net: deep learning for cell counting, detection, and morphometry,” *Nature Methods*, vol. 16, pp. 67–70, 01 2019. [Cited on page 22]
- [43] C. X. Hernández, M. M. Sultan, and V. S. Pande, “Using deep learning for segmentation and counting within microscopy data,” 2018. [Cited on page 22]
- [44] J. Ke, Z. Jiang, C. Liu, T. Bednarz, A. Sowmya, and X. Liang, “Selective detection and segmentation of cervical cells,” *ICBBT’19: Proceedings of the 2019 11th International Conference on Bioinformatics and Biomedical Technology*, pp. 55–61, 05 2019. [Cited on page 22]
- [45] I. Huh, “Blood cell detection using singleshot multibox detector,” 2018. [Cited on page 22]
- [46] W. Xie, J. Noble, and A. Zisserman, “Microscopy cell counting and detection with fully convolutional regression networks,” *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*, pp. 1–10, 05 2016. [Cited on page 22]
- [47] X. Zhu, X. Li, K. H. Ong, W. Zhang, W. Li, L. Li, D. Young, Y. Su, B. Shang, L. Peng, W. Xiong, Y. Liu, W. Liao, J. Xu, F. Wang, Q. Liao, S. Li, M. Liao, Y. Li, and Y. Ding, “Hybrid ai-assistive diagnostic model permits rapid tbs classification of cervical liquid-based thin-layer cell smears,” *Nature Communications*, vol. 12, p. 3541, 06 2021. [Cited on page 22]
- [48] K. de Haan, H. Koydemir, Y. Rivenson, D. Tseng, E. Dyne, L. Bakic, D. Karınca, K. Liang, M. Ilango, E. Gumustekin, and A. Ozcan, “Automated screening of sickle cells using a smartphone-based microscope and deep learning,” *npj Digital Medicine*, vol. 3, p. 76, 12 2020. [Cited on page 23]
- [49] Q. Wei, H. Qi, W. Luo, D. Tseng, S. Ki, Z. Wan, Z. Göröcs, L. Bentolila, T.-T. Wu, R. Sun, and A. Ozcan, “Fluorescent imaging of single nanoparticles and viruses on a smart phone,” *ACS nano*, vol. 7, p. 9147–9155, 09 2013. [Cited on page 23]

- [50] L. Rosado, J. M. Correia da Costa, D. Elias, and J. Cardoso, “Automated detection of malaria parasites on thick blood smears via mobile devices,” *Procedia Computer Science*, vol. 90, pp. 138–144, 12 2016. [Cited on page 23]
- [51] O. Holmström, N. Linder, B. Ngasala, A. Mårtensson, E. Linder, M. Lundin, H. Moilanen, A. Suutala, V. Diwan, and J. Lundin, “Point-of-care mobile digital microscopy and deep learning for the detection of soil-transmitted helminths and schistosoma haematobium,” *Global Health Action*, vol. 10, p. 1337325, 06 2017. [Cited on page 23]
- [52] Y. Rivenson, H. Ceylan Koydemir, H. Wang, Z. Wei, Z. Ren, H. Günaydın, Y. Zhang, Z. Göröcs, K. Liang, D. Tseng, and A. Ozcan, “Deep learning enhanced mobile-phone microscopy,” *ACS Photonics*, vol. 5, no. 6, pp. 2354–2364, 2018. [Cited on page 23]
- [53] Fraunhofer AICOS, “CLARE - Computer-Aided Cervical Cancer Screening.” Available at https://www.aicos.fraunhofer.pt/en/our_work/projects/clare.html. (Last accessed in 02/12/2021). [Cited on page 25]
- [54] Fraunhofer AICOS, “DEMalariaScope - Automatic detection of malaria in blood smears using smartphones.” Available at https://www.aicos.fraunhofer.pt/en/our_work/projects/malariascope.html. (Last accessed in 02/12/2021). [Cited on page 25]
- [55] P. Brandão, P. T. Silva, M. Parente, and L. Rosado, “ μ smartscope – towards a low-cost microscopic medical device for cervical cancer screening using additive manufacturing and optimization,” *Proceedings of the Institution of Mechanical Engineers, Part L: Journal of Materials: Design and Applications*, vol. 236, no. 2, pp. 267–279, 2022. [Cited on page 27]
- [56] C. Pereira, P. T. Silva, L. Rosado, L. Mota, and J. Martins, “The design thinking process in the development of an intelligent microscopic equipment,” in *Advances in Design and Digital Communication II* (N. Martins and D. Brandão, eds.), (Cham), pp. 170–182, Springer International Publishing, 2022. [Cited on page 27]
- [57] J. Jantzen, J. Norup, G. Dounias, and B. Bjerregaard, “Pap-smear benchmark data for pattern classification,” in *Proc. NiSIS 2005*, pp. 1–9, NiSIS, 2005. Nature inspired Smart Information Systems : EU co-ordination action, Nisis 2005 ; Conference date: 01-01-2005. [Cited on page 27]
- [58] M. E. Plissiti, P. Dimitrakopoulos, G. Sfikas, C. Nikou, O. Krikoni, and A. Charchanti, “Sipakmed: A new dataset for feature and image based classification of normal and pathological cervical cells in pap smear images,” in *2018*

- 25th IEEE International Conference on Image Processing (ICIP)*, pp. 3144–3148, 2018. [Cited on page 27]
- [59] H. A. Phoulady and P. R. Mouton, “A new cervical cytology dataset for nucleus detection and image classification (cervix93) and methods for cervical nucleus detection,” 2018. [Cited on page 27]
- [60] Z. Lu, G. Carneiro, A. Bradley, D. Ushizima, M. S. Nosrati, A. Bianchi, C. Carneiro, and G. Hamarneh, “Evaluation of three algorithms for the segmentation of overlapping cervical cells,” *IEEE Journal of Biomedical and Health Informatics*, vol. 21, pp. 441–450, 01 2016. [Cited on page 27]
- [61] Z. Lu, G. Carneiro, and A. P. Bradley, “An improved joint optimization of multiple level set functions for the segmentation of overlapping cervical cells,” *Trans. Img. Proc.*, vol. 24, p. 1261–1272, apr 2015. [Cited on page 27]
- [62] V. Lee, S.-B. Ng, and M. Salto-Tellez, “Chapter 34 - new techniques,” in *Diagnostic Cytopathology (Third Edition)* (W. Gray and G. Kocjan, eds.), pp. 891–902, Edinburgh: Churchill Livingstone, third edition ed., 2010. [Cited on page 36]
- [63] D. H. Marimont and B. A. Wandell, “Matching color images: the effects of axial chromatic aberration,” *J. Opt. Soc. Am. A*, vol. 11, pp. 3113–3122, Dec 1994. [Cited on page 37]
- [64] Eurocytology, “Criteria for adequacy of a cervical cytology sample.” Available at <https://www.eurocytology.eu/en/course/1142>. (Last accessed in 02/11/2021). [Cited on page 39]
- [65] Andrew Ng, Alex Ratner, Joaquin Vanschoren, Salwa Nur Muhammad, Carlos Alzate, Dillon Laird, “Data-centric AI: Real World Approaches.” Available at <https://youtu.be/Yqj7Kyjznh4>. (Last accessed in 03/01/2022). [Cited on page 41]
- [66] Data-Centric AI, “Data-Centric AI Competition.” Available at <https://https-deeplearning-ai.github.io/data-centric-comp/>. (Last accessed in 03/01/2021). [Cited on page 41]
- [67] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International Journal of Computer Vision*, vol. 88, pp. 303–338, Jun 2010. [Cited on page 49]
- [68] COCO Common Objects in Context, “Detection Evaluation.” Available at <https://cocodataset.org/#detection-eval>. (Last accessed in 05/01/2021). [Cited on page 49]

- [69] R. Padilla, W. Lobato Passos, T. Dias, S. Netto, and E. da Silva, “A comparative analysis of object detection metrics with a companion open-source toolkit,” *Electronics*, vol. 10, pp. 279–306, 01 2021. [Cited on page 53]

Appendix A

Detailed description of dataset split

TRAIN - TEST SPLIT															
Images Divison			Annotations Division												
			Glandular nucleus			Squamous nucleus			Inflammatory cell			Undefined nucleus		Artifact	
%Images	#Images	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations
73%	556	18562	79%	15543	71%	14577	81%	491	74%	831	84%	123	81%	123	81%
27%	209	6729	21%	6427	29%	3333	19%	173	26%	164	16%	29	19%	29	19%
100%	765	25291	100%	21970	100%	17910	100%	664	100%	995	100%	132	100%	132	100%
FOLD 1 SPLIT (3-fold Cross-validation)															
Images Divison			Annotations Division												
			Glandular nucleus			Squamous nucleus			Inflammatory cell			Undefined nucleus		Artifact	
%Images	#Images	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations
65%	363	12120	67%	10267	66%	9773	67%	282	57%	558	67%	81	66%	81	66%
35%	193	6442	33%	5276	34%	4804	33%	209	43%	273	33%	43	34%	43	34%
100%	556	18562	100%	15543	100%	14577	100%	491	100%	831	100%	123	100%	123	100%
FOLD 2 SPLIT (3-fold Cross-validation)															
Images Divison			Annotations Division												
			Glandular nucleus			Squamous nucleus			Inflammatory cell			Undefined nucleus		Artifact	
%Images	#Images	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations
67%	370	12356	56%	10062	65%	9486	65%	376	77%	519	62%	82	67%	82	67%
33%	186	6202	44%	5481	35%	5091	38%	115	23%	312	38%	41	33%	41	33%
100%	556	18562	100%	15543	100%	14577	100%	491	100%	831	100%	123	100%	123	100%
FOLD 3 SPLIT (3-fold Cross-validation)															
Images Divison			Annotations Division												
			Glandular nucleus			Squamous nucleus			Inflammatory cell			Undefined nucleus		Artifact	
%Images	#Images	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations	#Annotations	%Annotations
68%	379	12648	77%	10757	69%	9895	68%	324	66%	585	70%	83	67%	83	67%
32%	177	5914	23%	4786	31%	4682	32%	167	34%	246	30%	40	33%	40	33%
100%	556	18562	100%	15543	100%	14577	100%	491	100%	831	100%	123	100%	123	100%

Table A.1: Dataset division: Stratification of nuclei annotations

Appendix B

Tensorboard charts: mAP for each class during training

Blue and Orange lines correspond to the training of model ResNet50 and MobileNet respectively, with Batch Size of 16 and LR of 5.29×10^{-4} . Training and validation was performed on fold 1. Horizontal axis represent the number of steps of training and vertical corresponds to the mAP value. The slightly faded lines correspond to the real values, while the vibrant ones to the is the representation of the same data but after the smooth coefficient applied.

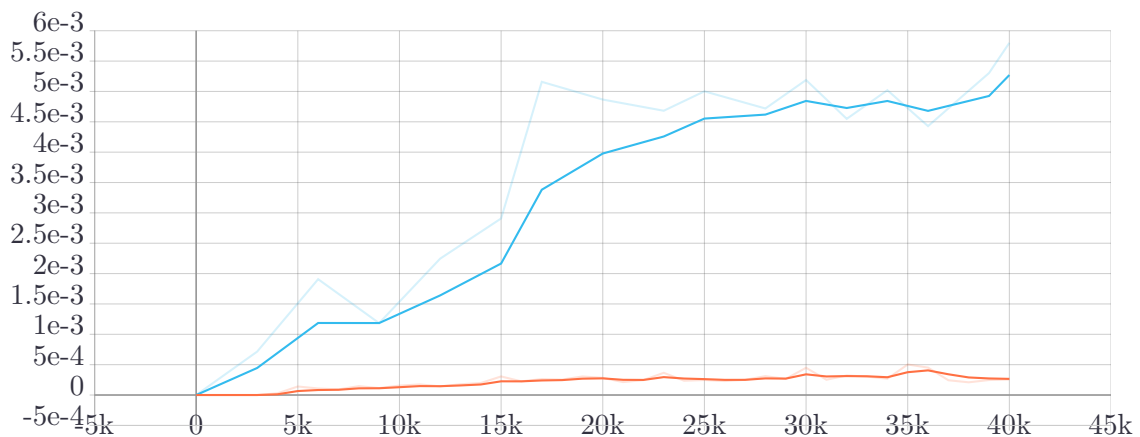


Figure B.1: Artifact

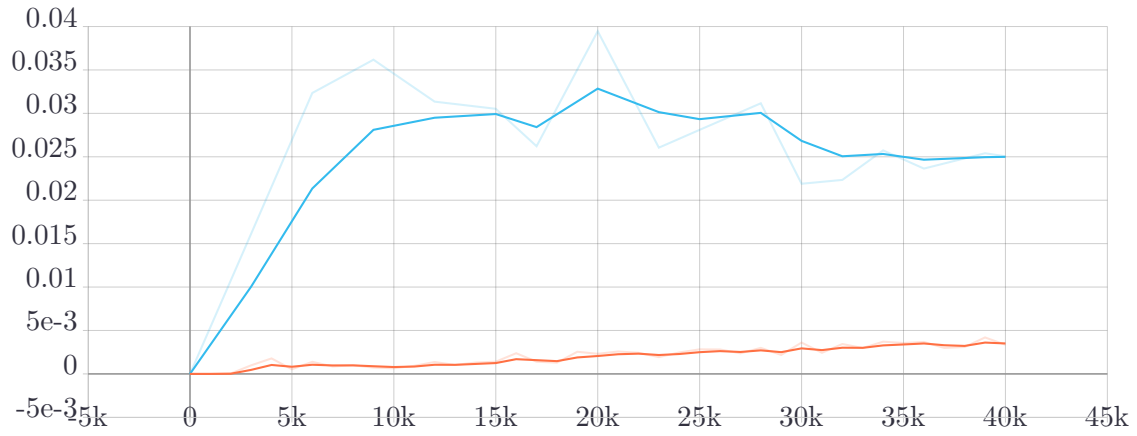


Figure B.2: Glandular nucleus

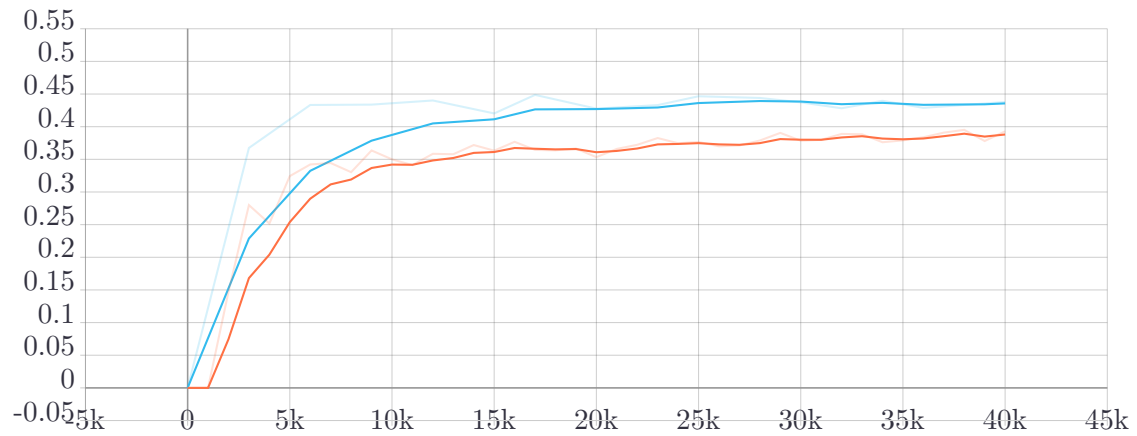


Figure B.3: Inflammatory cell

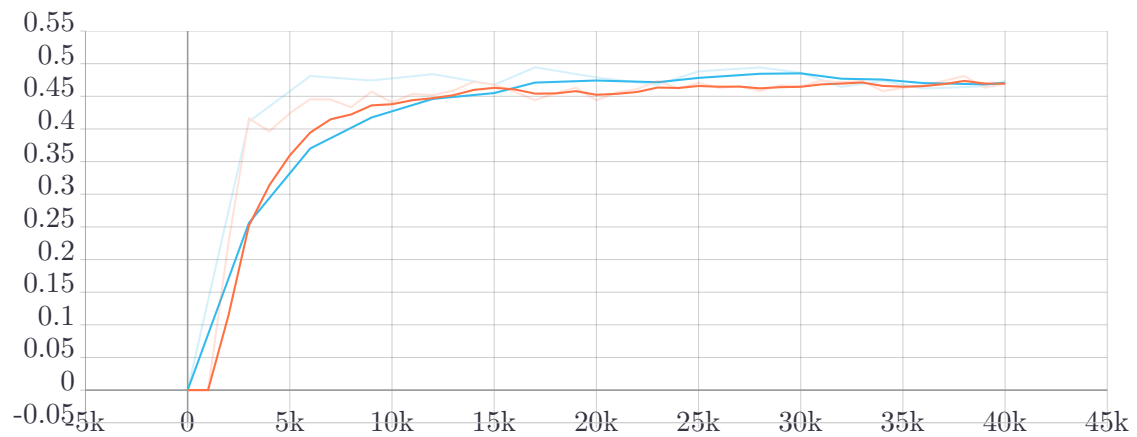


Figure B.4: Squamous nucleus

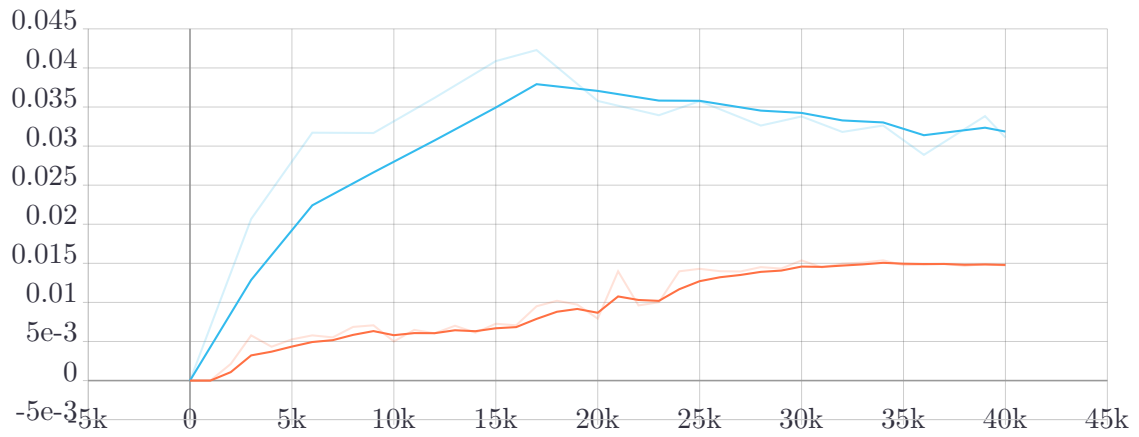


Figure B.5: Undefined all

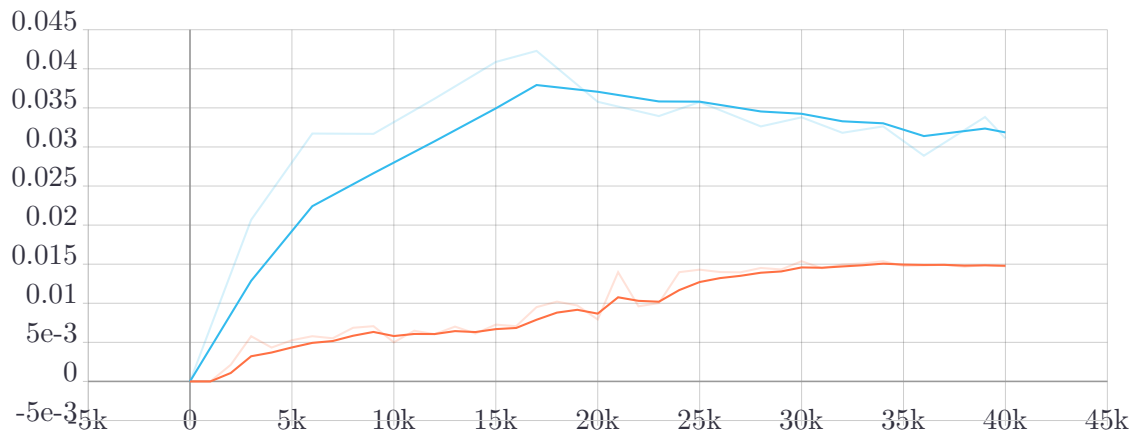


Figure B.6: Undefined nucleus

Appendix C

Hyperparameters tuning tests and results on test set

Meta-architecture	Backbone	Input Size / Image Size	#Training Steps	Batch Size	3 Folds validation results			3 fold average																							
					Learning Rate	Cross val	Best Checkpoint	MAP@50	squamous_nucleus	inflammatory_cell	glandular_nucleus	undefined_nucleus	artifact	squamous_nucleus	inflammatory_cell	glandular_nucleus	undefined_nucleus	artifact													
SSD	Resnet50	320 / 320	40000	16	0,305638	1	20000	0,430	0,356	0,002	0,001	0,005	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
						2	40000	0,453	0,375	0,010	0,013	0,007	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000			
						3	34000	0,447	0,419	0,013	0,013	0,013	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	
						1	80000	0,430	0,356	0,002	0,001	0,005	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	
						2	77000	0,410	0,336	0,001	0,011	0,004	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	
						3	68000	0,400	0,343	0,002	0,004	0,006	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	92000	0,394	0,315	0,001	0,003	0,003	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						2	14000	0,510	0,1829	0,445	0,065	0,047	0,030	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	17000	0,1756	0,494	0,449	0,026	0,037	0,042	0,005	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	68000	0,1465	0,474	0,384	0,001	0,006	0,014	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	66000	0,149	0,481	0,391	0,005	0,007	0,009	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	80000	0,1543	0,470	0,421	0,006	0,010	0,017	0,003	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						2	124000	0,1234	0,423	0,310	0,000	0,003	0,003	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	87000	0,1476	0,470	0,397	0,003	0,005	0,007	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	40000	0,1383	0,461	0,351	0,003	0,005	0,009	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						2	40000	0,1422	0,453	0,352	0,024	0,009	0,014	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	40000	0,1397	0,445	0,371	0,006	0,007	0,009	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	80000	0,136	0,457	0,342	0,003	0,006	0,009	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						2	48000	0,1342	0,453	0,333	0,004	0,007	0,007	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	55000	0,1364	0,440	0,358	0,004	0,009	0,007	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	64000	0,1272	0,445	0,308	0,001	0,004	0,006	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						2	142000	0,1236	0,442	0,287	0,002	0,005	0,006	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						3	73000	0,1261	0,426	0,321	0,001	0,003	0,005	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
						1	40000	0,07767	0,277	0,187	0,000	0,000	0,002	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000
2	40000	0,08203	0,295	0,193	0,000	0,002	0,002	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
3	40000	0,08282	0,278	0,215	0,000	0,001	0,002	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
1	80000	0,0782	0,298	0,170	0,000	0,000	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
2	80000	0,07878	0,292	0,178	0,000	0,001	0,002	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
3	80000	0,07525	0,263	0,185	0,001	0,000	0,003	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
1	160000	0,0625	0,245	0,129	0,000	0,000	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
2	160000	0,05975	0,246	0,111	0,000	0,000	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
3	160000	0,05559	0,206	0,126	0,000	0,000	0,001	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						
EfficientDet D1	Efficientnet	512x512	40000	1	0,000529	33000	0,02932	0,1039	0,07031	4,63E-03	8,22E-03	7,23E-04	6,11E-03	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000						

Table C.2: Results for model ResNet50 and EfficientDet D1 using cross-validation folds.

Meta-architecture	Backbone	Input Size / Image Size	#Training Steps	#Training Epochs	Training Status	Experiment ID	Validation on TEST set and trained on TRAIN set				Best Checkpoint	mAP@50	squamous_nucleus	inflatmat_ory_cell	glandular_nucleus	undefined_nucleus_all	artifact	RECALL
							Batch Size	Learning Rate	Learning Rate	Learning Rate								
SSD	mobilenet_v2	320	40000	32,7	Done	mobilenet_1_16_02499999821_TEST_SET	16	0.024999998	39000	0.1351	0.4534	0.3405	4.90E-03	5.50E-03	5.90E-03	5.90E-03	5.50E-05	0.3745
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_2	16	0.024999998	45000	0.1562	0.4734	0.3694	0.0684	0.02092	5.00E-03	1.20E-04	0.4373	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_3	16	0.024999998	37000	0.1591	0.4736	0.3698	0.0815	0.02558	5.22E-03	8.14E-05	0.4158	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_4	16	0.024999998	26000	0.1431	0.4454	0.3492	0.04217	0.01655	5.40E-03	2.90E-05	0.3899	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_4	16	0.024999998	46000	0.1592	0.4542	0.3577	0.1172	0.02031	5.75E-03	2.10E-04	0.427	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_4	16	0.024999998	37000	0.1534	0.4549	0.3408	0.08664	0.01954	6.20E-03	1.18E-04	0.417	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_5	16	0.024999998	28000	0.1581	0.4706	0.3734	0.07673	0.02197	5.47E-03	3.98E-04	0.4323	
SSD	mobilenet_v2	320	61000	49,9	Done	mobilenet_1_16_02499999821_TEST_SET_6	16	0.024999998	26000	0.1477	0.4614	0.358	0.03529	0.02518	6.20E-03	5.97E-05	0.4312	
SSD	mobilenet_v2	320	100000	81,8	Done	mobilenet_5_perclass_TEST_SET	16	0.15	92000	0.1402	0.4544	0.355	0.0138	0.01247	5.26E-03	1.70E-04	0.4088	
SSD	resnet50	320	40000	32,7	Done	resnet50_1_16_000528663502_TEST_SET	16	0.005286635	35000	0.1529	0.4661	0.3574	6.24E-02	2.67E-02	4.60E-04	9.50E-05	0.4093	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_2	16	0.005286635	61000	0.1566	0.4648	0.3596	8.41E-02	2.57E-02	5.20E-03	1.90E-04	0.3914	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_2	16	0.005286635	12000	0.1498	0.4742	0.3488	3.95E-03	3.13E-02	5.30E-03	2.60E-05	0.4089	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_3	16	0.005286635	50000	0.134	0.4385	0.3384	7.37E-03	3.85E-03	0.01563	1.98E-05	0.3701	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_3	16	0.005286635	58000	0.1408	0.4476	0.3629	0.01214	0.01452	7.47E-03	2.62E-05	0.3917	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_4	16	0.005286635	15000	0.1531	0.4612	0.3707	0.04801	0.03231	6.14E-03	1.07E-04	0.4499	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_4	16	0.005286635	28000	0.1576	0.4789	0.3689	0.0524	0.04014	5.23E-03	1.04E-04	0.4205	
SSD	resnet50	320	61000	49,9	Done	resnet50_1_16_000528663502_TEST_SET_5	16	0.005286635	48000	0.1569	0.4728	0.3773	0.062	0.02254	6.31E-03	3.26E-04	0.4297	

Table C.3: Results for model ResNet50 and MobileNet for best parameters found during hyperparameter tuning. Results obtained using TRAIN and TEST sets.

Appendix D

Results for merging and eliminating classes approach

Validation on TEST set and trained on TRAIN set														
Meta-architecture	Backbone	Input Size / Image Size	#Training Steps	#Training Epochs	Training Status	Experiment ID	Batch Size	Learning Rate	Best Checkpoint	mAP@50	squamous	inflatmat	RECALL	
											_nucleus	ory_cell	100	
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v6_1	16	0.005286635	15000	0.296	0.474	0.3541	5.99E-02	0.5848
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v6_2	16	0.005286635	35000	0.3066	0.4783	0.3589	8.28E-02	0.5773
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v6_3	16	0.005286635	35000	0.3058	0.4799	0.3675	7.23E-02	0.5799
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v7_1	16	0.005286635	34000	0.4538	0.4538			0.6314
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v7_2	16	0.005286635	57000	0.4086	0.4086			0.6173
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v7_3	16	0.005286635	59000	0.447	0.447			0.6341
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v8_1	16	0.005286635	35000	0.4252	0.4773	0.3731		0.63
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v8_2	16	0.005286635	33000	0.4174	0.4649	0.3699		0.6263
SSD	resnet50	320	61000	49.9	Done	resnet50_1_16_000528663502_TEST_SET_v8_3	16	0.005286635	25000	0.4209	0.4738	0.368		0.6285

Table D.1: Results for model ResNet50 with some classes merged and eliminated.

Appendix E

Results of testing Overlap, NMS threshold and IOU threshold for the final system

This appendix consists of a filtered table divided in four parts (From Table E.1 to Table E.4) and an small part of the complete table, also divided across four tables (From Table E.5 to Table E.8) of results with different Overlap, NMS threshold and IOU threshold values.

E.1 Filtered table

Since we knew that some of our predictions, even though reported as FP, were TP, and it could be reflected on our evaluation results if all objects in the dataset were annotated, we decided to filter out more than 2800 tests to only those that report number of predictions of squamous nucleus greater than roughly 6200 predictions, which is around the number of annotations in this class (6427).

The table is ordered from (top) biggest Youden value to the lowest (bottom), and represents the best set of parameters at the top. Although the values that represent bigger overlap (>0.3) seem to give the best set of parameters, they were not considered to be in the final system due to the fact that they yield much more images than a much lower percentage of overlap, which reflects negatively on the processing

speed and a slightly better results do not compensate for the time necessary to complete the detection stage.

After visualizing the following best set of parameters (23% of overlap, 0.99 for NMS threshold and 0.15 for IOU between predictions and ground truth boxes) we acknowledged that NMS threshold of 0.99 was a bad fit and it resulted in many duplicate predictions to be unchanged. For this reason we took a step back and tried to get the best set of parameters from all tests, which can be seen in Section E.2.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z																
	TOP (Total Predictions)				IGT - Predict				TPR (True positives rate - recall)				TNR (True negative rate - specificity)				PPV (Precision - positive predictive value)				Accuracy				F1		FN		FP		AUPR		Youden									
1	model	overlap	nms_th	metric_th	inflamm	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ	inflat	squ		
2	resnet50_1	0.7	0.975	0.1	4813	6289	619	55	0.769	0.8173	0.812	0.813	0.67	0.81	0.798	0.815	0.716	0.814	969	1139	1588	1194	0.72	0.8137	0.5808	0.6307																
3	resnet50_1	0.8	0.965	0.1	4646	6413	452	179	0.734	0.8261	0.815	0.804	0.663	0.803	0.788	0.815	0.697	0.814	1114	1084	1566	1263	0.699	0.8146	0.5495	0.6296																
4	resnet50_1	0.8	0.96	0.1	4225	6279	31	45	0.674	0.8181	0.831	0.811	0.669	0.812	0.778	0.814	0.671	0.815	1369	1134	1400	1179	0.671	0.8152	0.5044	0.629																
5	resnet50_1	0.8	0.97	0.1	5053	6588	859	354	0.781	0.8314	0.798	0.791	0.649	0.787	0.792	0.811	0.709	0.808	917	1051	1776	1405	0.715	0.8091	0.5789	0.6227																
6	resnet50_1	0.8	0.965	0.1	4010	6315	184	81	0.676	0.8181	0.856	0.801	0.707	0.808	0.795	0.81	0.691	0.813	1358	1134	1174	1215	0.692	0.8128	0.5323	0.6199																
7	resnet50_1	0.8	0.97	0.1	4304	6482	110	248	0.716	0.824	0.845	0.788	0.697	0.793	0.802	0.806	0.707	0.808	1192	1097	1302	1345	0.707	0.8083	0.5604	0.6119																
8	resnet50_1	0.8	0.975	0.1	5372	6826	1178	592	0.82	0.8337	0.788	0.77	0.64	0.761	0.798	0.8	0.719	0.796	754	1037	1932	1629	0.73	0.7975	0.6087	0.604																
9	resnet50_1	0.7	0.98	0.1	4961	6578	767	344	0.784	0.8187	0.81	0.782	0.663	0.776	0.802	0.799	0.719	0.797	904	1130	1671	1474	0.724	0.7973	0.5943	0.6003																
10	resnet50_1	0.5	0.985	0.1	4764	6378	570	144	0.768	0.8011	0.823	0.792	0.677	0.783	0.805	0.796	0.72	0.792	971	1240	1541	1384	0.723	0.792	0.5912	0.593																
11	resnet50_1	0.8	0.975	0.1	4535	6737	341	503	0.746	0.8268	0.839	0.763	0.69	0.765	0.808	0.794	0.717	0.795	1065	1080	1406	1583	0.718	0.7959	0.5846	0.584																
12	resnet50_1	0.8	0.98	0.1	6260	6602	2066	368	0.873	0.7955	0.736	0.79	0.585	0.751	0.777	0.792	0.7	0.773	534	1275	2600	1643	0.729	0.7733	0.6086	0.585																
13	resnet50_1	0.7	0.98	0.1	4067	6354	127	120	0.702	0.799	0.866	0.784	0.724	0.784	0.811	0.791	0.713	0.791	1251	1253	1124	1373	0.713	0.7915	0.5678	0.5829																
14	resnet50_1	0.6	0.985	0.1	4856	6593	662	359	0.777	0.8081	0.821	0.773	0.671	0.764	0.807	0.79	0.72	0.786	935	1196	1597	1555	0.724	0.785	0.5975	0.5814																
15	mobilenet	0.23	0.99	0.15	3869	6322	325	88	0.687	0.7924	0.882	0.783	0.744	0.781	0.817	0.788	0.714	0.787	1314	1294	989	1382	0.716	0.7869	0.5689	0.5751																
16	mobilenet	0.23	0.99	0.1	3869	6322	325	88	0.69	0.7927	0.883	0.782	0.747	0.782	0.819	0.787	0.717	0.787	1302	1292	977	1380	0.719	0.7872	0.5728	0.5748																
17	mobilenet	0.23	0.99	0.2	3869	6322	325	88	0.685	0.7918	0.882	0.783	0.743	0.781	0.816	0.787	0.713	0.786	1321	1298	996	1386	0.714	0.7863	0.5667	0.5745																
18	mobilenet	0.23	0.99	0.25	3869	6322	325	88	0.683	0.7913	0.881	0.783	0.741	0.78	0.815	0.787	0.711	0.786	1328	1301	1003	1389	0.711	0.7858	0.5644	0.5741																
19	mobilenet	0.23	0.99	0.3	3869	6322	325	88	0.682	0.7905	0.88	0.782	0.739	0.78	0.815	0.786	0.709	0.785	1335	1306	1010	1394	0.71	0.785	0.5621	0.5723																
20	mobilenet	0.23	0.99	0.35	3869	6322	325	88	0.681	0.79	0.88	0.782	0.738	0.779	0.814	0.786	0.709	0.784	1337	1309	1012	1397	0.71	0.7845	0.5615	0.573																
21	mobilenet	0.23	0.99	0.4	3869	6322	325	88	0.68	0.7892	0.88	0.782	0.737	0.778	0.813	0.786	0.707	0.784	1344	1314	1019	1402	0.708	0.7837	0.5593	0.5712																
22	mobilenet	0.2	0.99	0.15	3688	6256	506	22	0.658	0.7868	0.888	0.784	0.748	0.784	0.811	0.786	0.7	0.785	1434	1329	928	1351	0.703	0.7854	0.5463	0.5712																
23	mobilenet	0.2	0.99	0.1	3688	6256	506	22	0.66	0.7873	0.889	0.784	0.751	0.785	0.812	0.786	0.703	0.786	1424	1326	918	1348	0.706	0.7859	0.5495	0.5711																
24	resnet50_1	0.2	0.99	0.1	4468	6265	274	31	0.728	0.7798	0.839	0.791	0.684	0.776	0.803	0.786	0.705	0.778	1139	1373	1413	1404	0.706	0.784	0.567	0.5706																
25	mobilenet	0.2	0.99	0.2	3688	6256	506	22	0.656	0.786	0.887	0.784	0.746	0.783	0.81	0.785	0.698	0.785	1444	1334	938	1356	0.701	0.7846	0.5431	0.5705																
26	mobilenet	0.25	0.99	0.1	3741	6360	453	126	0.67	0.7936	0.889	0.777	0.752	0.778	0.816	0.785	0.709	0.786	1382	1287	929	1413	0.711	0.7845	0.5595	0.5702																
27	mobilenet	0.25	0.99	0.15	3741	6360	453	126	0.667	0.7926	0.888	0.777	0.748	0.777	0.814	0.785	0.705	0.785	1395	1293	942	1419	0.708	0.7847	0.5552	0.5689																
28	mobilenet	0.2	0.99	0.25	3688	6256	506	22	0.654	0.7849	0.887	0.784	0.743	0.782	0.809	0.784	0.696	0.784	1452	1341	946	1363	0.699	0.7833	0.5405	0.5695																
29	mobilenet	0.24	0.99	0.15	3876	6403	318	169	0.691	0.7947	0.884	0.774	0.748	0.774	0.82	0.784	0.718	0.784	1296	1280	978	1449	0.719	0.7842	0.5714	0.5688																
30	mobilenet	0.24	0.99	0.2	3876	6403	318	169	0.688	0.7939	0.883	0.774	0.745	0.773	0.819	0.784	0.715	0.783	1308	1285	990	1454	0.716	0.7834	0.5713	0.5681																
31	mobilenet	0.25	0.99	0.2	3741	6360	453	126	0.666	0.7915	0.887	0.776	0.747	0.776	0.814	0.784	0.704	0.784	1401	1300	948	1426	0.706	0.7836	0.5534	0.568																
32	mobilenet	0.25	0.99	0.25	3741	6360	453	126	0.664	0.7911	0.887	0.777	0.744	0.775	0.813	0.784	0.702	0.783	1410	1302	957	1428	0.704	0.7833	0.5504	0.5679																
33	mobilenet	0.24	0.99	0.1	3876	6403	318	169	0.695	0.7947	0.886	0.773	0.752	0.774	0.823	0.784	0.723	0.784	1278	1280	960	1449																				

98 Appendix E. Results of testing Overlap, NMS threshold and IOU threshold for the final system

1	model	overlap	nms_th	metric_th	TOP (Total Predictions)		TPR (True positives rate - recall)		TNR (True negative rate - Specificity)		PPV (Precision - positive predictive value)				Accuracy	F1	FN	FP	AUPR	Youden						
					inflam	squ	infla	squ	infla	squ	infla	squ	infla	squ							infla	squ	infla	squ		
101	mobilenet	0.26	0.99	0.4	3896	6445	298	211	0.682	0.7907	0.879	0.768	0.734	0.765	0.814	0.779	0.707	0.778	1334	1305	1036	1516	0.708	0.7739	0.6712	0.5589
102	resnet50_1	0.8	0.98	0.1	5568	7343	1374	1109	0.833	0.8353	0.788	0.724	0.628	0.709	0.801	0.773	0.716	0.767	700	1027	2074	2136	0.73	0.7222	0.6207	0.5588
103	mobilenet	0.22	0.99	0.1	4558	6297	364	63	0.736	0.774	0.824	0.785	0.677	0.766	0.803	0.78	0.706	0.77	1106	1409	1470	1472	0.707	0.7701	0.5708	0.5588
104	resnet50_1	0.26	0.99	0.1	4595	6417	401	183	0.743	0.7825	0.824	0.776	0.678	0.76	0.805	0.779	0.707	0.771	1077	1356	1478	1539	0.711	0.7713	0.5774	0.5587
105	mobilenet	0.3	0.99	0.25	3844	6450	350	216	0.67	0.791	0.879	0.767	0.73	0.764	0.81	0.779	0.699	0.778	1386	1303	1036	1519	0.7	0.7777	0.5487	0.5587
106	mobilenet	0.22	0.99	0.35	3814	6349	380	115	0.672	0.7846	0.883	0.774	0.739	0.77	0.813	0.779	0.704	0.777	1376	1343	996	1458	0.705	0.7775	0.5545	0.5583
107	mobilenet	0.19	0.99	0.4	3803	6301	391	67	0.669	0.7807	0.882	0.778	0.738	0.772	0.812	0.779	0.702	0.777	1389	1367	998	1434	0.703	0.7766	0.5512	0.5583
108	mobilenet	0.22	0.99	0.4	3814	6349	380	115	0.67	0.7842	0.882	0.774	0.737	0.77	0.812	0.779	0.702	0.777	1383	1345	1003	1460	0.704	0.7771	0.5522	0.5581
109	mobilenet	0.3	0.99	0.3	3844	6450	350	216	0.669	0.7903	0.879	0.767	0.729	0.764	0.81	0.779	0.698	0.777	1390	1307	1040	1523	0.699	0.7721	0.5474	0.5576
110	mobilenet	0.3	0.99	0.35	3844	6450	350	216	0.667	0.7894	0.878	0.767	0.728	0.763	0.809	0.778	0.696	0.776	1395	1313	1045	1529	0.698	0.7762	0.5459	0.5563
111	mobilenet	0.3	0.99	0.4	3844	6450	350	216	0.666	0.7889	0.878	0.767	0.727	0.762	0.808	0.778	0.695	0.775	1401	1316	1051	1532	0.696	0.7757	0.5439	0.5557
112	resnet50_1	0.25	0.99	0.1	4598	6401	404	167	0.74	0.779	0.883	0.776	0.675	0.759	0.803	0.778	0.706	0.769	1092	1378	1496	1545	0.707	0.7688	0.5724	0.5552
113	mobilenet	0.29	0.99	0.1	3990	6525	204	291	0.699	0.7944	0.877	0.761	0.734	0.759	0.818	0.777	0.716	0.776	1264	1282	1060	1575	0.716	0.7766	0.5755	0.5549
114	mobilenet	0.29	0.99	0.15	3990	6525	204	291	0.696	0.7932	0.876	0.761	0.732	0.758	0.817	0.776	0.713	0.775	1275	1289	1071	1580	0.714	0.7755	0.5721	0.5539
115	mobilenet	0.7	0.985	0.1	5707	6540	1513	306	0.838	0.7767	0.773	0.777	0.616	0.74	0.792	0.777	0.71	0.758	679	1392	2192	1608	0.727	0.7585	0.6108	0.5533
116	mobilenet	0.29	0.99	0.25	3990	6525	204	291	0.693	0.7934	0.875	0.761	0.728	0.757	0.816	0.776	0.71	0.774	1288	1294	1084	1585	0.711	0.7748	0.5678	0.5533
117	mobilenet	0.29	0.99	0.2	3990	6525	204	291	0.694	0.7936	0.875	0.761	0.73	0.757	0.816	0.776	0.711	0.775	1283	1293	1097	1584	0.712	0.7749	0.5695	0.5532
118	mobilenet	0.29	0.99	0.3	3990	6525	204	291	0.691	0.7916	0.874	0.761	0.726	0.756	0.815	0.776	0.708	0.774	1296	1299	1092	1590	0.709	0.774	0.5652	0.5523
119	resnet50_1	0.07	0.995	0.1	4625	6250	431	16	0.723	0.7628	0.825	0.789	0.655	0.761	0.793	0.776	0.687	0.762	1163	1479	1594	1495	0.689	0.7618	0.5477	0.5513
120	mobilenet	0.29	0.99	0.35	3990	6525	204	291	0.69	0.7908	0.874	0.76	0.725	0.756	0.814	0.775	0.707	0.773	1300	1304	1096	1595	0.708	0.7732	0.564	0.5512
121	mobilenet	0.29	0.99	0.4	3990	6525	204	291	0.69	0.7899	0.874	0.76	0.725	0.755	0.814	0.774	0.707	0.772	1302	1310	1098	1601	0.707	0.7722	0.5635	0.5497
122	mobilenet	0.07	0.995	0.2	3856	6257	338	23	0.664	0.77	0.876	0.779	0.722	0.767	0.807	0.775	0.692	0.769	1410	1434	1072	1457	0.693	0.7686	0.5398	0.5494
123	mobilenet	0.07	0.995	0.15	3856	6257	338	23	0.666	0.7698	0.877	0.779	0.724	0.767	0.808	0.774	0.694	0.768	1402	1435	1064	1458	0.695	0.7684	0.5424	0.5485
124	mobilenet	0.07	0.995	0.1	3856	6257	338	23	0.671	0.7705	0.876	0.778	0.73	0.768	0.81	0.774	0.699	0.769	1381	1431	1043	1454	0.7	0.769	0.5491	0.5481
125	mobilenet	0.07	0.995	0.25	3856	6257	338	23	0.663	0.7684	0.876	0.778	0.721	0.766	0.806	0.774	0.69	0.767	1415	1444	1077	1467	0.692	0.767	0.5383	0.5468
126	mobilenet	0.07	0.995	0.3	3856	6257	338	23	0.661	0.7679	0.875	0.778	0.719	0.765	0.805	0.773	0.689	0.766	1422	1447	1084	1470	0.69	0.7665	0.5359	0.5462
127	mobilenet	0.07	0.995	0.35	3856	6257	338	23	0.659	0.7668	0.874	0.778	0.717	0.764	0.804	0.772	0.687	0.765	1429	1454	1091	1477	0.688	0.7654	0.5337	0.5446
128	resnet50_1	0.4	0.99	0.1	4845	6796	651	562	0.758	0.7952	0.822	0.749	0.656	0.729	0.802	0.77	0.704	0.761	1014	1277	1665	1839	0.704	0.7623	0.5802	0.5437
129	mobilenet	0.07	0.995	0.4	3856	6257	338	23	0.658	0.7661	0.874	0.778	0.716	0.763	0.804	0.772	0.686	0.765	1434	1458	1090	1481	0.687	0.7647	0.5321	0.5437
130	resnet50_1	0.28	0.99	0.1	3748	6296	446	62	0.65	0.7659	0.882	0.772	0.728	0.761	0.806	0.771	0.687	0.765	1466	1440	1020	1502	0.685	0.7652	0.5323	0.5413
131	resnet50_1	0.7	0.985	0.1	5170	7255	976	1021	0.791	0.8211	0.807	0.719	0.641	0.706	0.802	0.765	0.708	0.759	878	1115	1854	2136	0.716	0.7634	0.5981	0.5398
132	mobilenet	0.08	0.995	0.2	3898	6425	296	191	0.664	0.7754	0.873	0.764	0.714	0.752	0.806	0.77	0.688	0.764	1409	1400	1113	1596	0.695	0.7603	0.5374	0.5396
133	mobilenet	0.08	0.995	0.1	3898	6425	296	191	0.67	0.7765	0.875	0.763	0.721	0.753	0.809	0.77	0.695	0.765	1384	1393	1088	1584	0.695	0.765	0.5453	0.5395
134	mobilenet	0.08	0.995	0.15	3898	6425	296	191	0.668	0.7759	0.875	0.764	0.718	0.753	0.808	0.77	0.692	0.764	1394	1397	1098	1588	0.693	0.7644	0.5422	0.5395
135	mobilenet	0.08	0.995	0.3	3898	6425	296	191	0.66	0.7745	0.872	0.764	0.71	0.751	0.804	0.769	0.684	0.763	1425	1406	1129	1597	0.685	0.763	0.5322	0.5389
136	mobilenet	0.08	0.995	0.25	3898	6425	296	191	0.663	0.7746	0.872	0.764	0.713	0.752	0.805	0.769	0.687	0.763	1415	1405	1119	1596	0.688	0.763	0.5355	0.5387
137	resnet50_1	0.27	0.99	0.1	3714	6287	480	53	0.643	0.7671	0.882	0.772	0.726	0.761	0.804	0.769	0.682	0.764	1496	1452	1016	1505	0.685	0.7639	0.5255	0.5386
138	resnet50_1	0.29	0.99	0.1	3720	6328	474	94	0.65	0.7697	0.885	0.778	0.733	0.758	0.808	0.769	0.689	0.764	1468	1436	994	1530	0.681	0.7639	0.5349	0.5377
139	mobilenet	0.08	0.995	0.35	3898	6425	296	191	0.66	0.7737	0.872	0.764	0.71	0.751	0.803	0.769	0.684	0.762	1428	1411	1132	1602	0.685	0.7622	0.5313	0.5377
140	mobilenet	0.3	0.99	0.1	3769	6270	425	36	0.653	0.7653	0.871	0.772	0.727	0.761	0.806	0.769	0.688	0.763	1454	1463	1029	1499	0.689	0.763	0.5333	0.5372
141	mobilenet	0.08	0.995	0.4	3898	6425	296	191	0.659	0.7732	0.872	0.764	0.709	0.75	0.803	0.768	0.683	0.762	1430	1414	1134	1605	0.684	0.7617	0.5307	0.537
142	resnet50_1	0.08	0.995	0.1	4660	6388	466	154	0.718	0.7623	0.822	0.774	0.646	0.744	0.79	0.768	0.683	0.763	1184	1482	1650	1636	0.682	0.7531	0.5337	0.536
143	resnet50_1	0.26	0.99	0.1	3622	6285	572	51	0.641	0.7666	0.891	0.769	0.742	0.76	0.809	0.768	0.688	0.763	1505	1455	933	1506	0.692	0.7635	0.5332	0.5359
144	mobilenet	0.5	0.99	0.1	5569	6465	1375	231	0.81	0.7584	0.779	0.777	0.61	0.731	0.788	0.769	0.696	0.745	798	1506	2173	1737	0.71	0.7449	0.5886	0.5353
145	resnet50_1	0.8	0.98	0.1	4664	7363	470	1129	0.759	0.8295	0.842	0.702	0.682	0.702	0.816	0.76	0.719	0.761	1011	1063	1481	2192	0.721	0.7659	0.6011	0.531
146	mobilenet	0.09	0.995	0.1	3904	6517	290	283	0.67	0.7772	0.876	0.753	0.72	0.743	0.81	0.765	0.693	0.762	1382	1389	1092	1672	0.695	0.7603	0.5467	0.5306

1	model	overlap	nms_th	metric_th	TOP (Total Predictions)		TNR (True negatives rate - recall)		PPV (Precision - positive predictive value)		Accuracy	F1	FN	FP	AUPR	Youden										
					inflam	squs	infla	squl	infla	squs							infla	squl	infla	squl						
201	mobilenet	0.15	0.995	0.3	3858	6731	336	497	0.643	0.774	0.873	0.733	0.699	0.717	0.801	0.752	0.67	0.744	1497	1409	1161	1906	0.671	0.7454	0.5165	0.5067
202	mobilenet	0.26	0.995	0.1	5718	6629	1524	395	0.791	0.7441	0.769	0.762	0.58	0.7	0.776	0.754	0.669	0.721	876	1595	2400	1990	0.686	0.722	0.5604	0.5062
203	mobilenet	0.13	0.995	0.1	4036	6806	158	572	0.68	0.7798	0.871	0.726	0.707	0.714	0.811	0.751	0.693	0.746	1341	1373	1183	1945	0.694	0.747	0.5509	0.506
204	resnet50_1	0.15	0.995	0.1	4880	6909	686	675	0.728	0.7732	0.815	0.733	0.626	0.698	0.789	0.751	0.673	0.733	1140	1414	1826	2089	0.677	0.7354	0.543	0.506
205	mobilenet	0.11	0.995	0.3	3858	6736	336	502	0.649	0.7748	0.875	0.731	0.705	0.717	0.804	0.752	0.676	0.745	1473	1404	1137	1906	0.677	0.7459	0.5242	0.5058
206	mobilenet	0.27	0.995	0.1	5854	6640	1660	406	0.794	0.7427	0.76	0.763	0.569	0.697	0.77	0.754	0.663	0.719	863	1604	2523	2010	0.682	0.72	0.5545	0.5058
207	mobilenet	0.29	0.995	0.1	5991	6676	1797	442	0.801	0.7425	0.754	0.763	0.561	0.693	0.767	0.755	0.66	0.717	834	1605	2631	2047	0.681	0.718	0.5548	0.5057
208	mobilenet	0.15	0.995	0.35	3858	6731	336	497	0.642	0.7732	0.873	0.732	0.698	0.716	0.801	0.751	0.669	0.744	1502	1414	1166	1911	0.67	0.7446	0.5149	0.5056
209	mobilenet	0.12	0.995	0.25	3948	6764	246	530	0.661	0.7756	0.872	0.73	0.702	0.715	0.806	0.751	0.681	0.744	1421	1399	1175	1929	0.682	0.7452	0.5333	0.5056
210	mobilenet	0.18	0.995	0.15	4064	6963	130	729	0.677	0.7875	0.868	0.718	0.698	0.705	0.809	0.75	0.687	0.744	1356	1325	1226	2054	0.688	0.7462	0.5451	0.5051
211	mobilenet	0.11	0.995	0.35	3858	6736	336	502	0.647	0.7741	0.875	0.731	0.704	0.716	0.803	0.751	0.674	0.744	1479	1408	1143	1910	0.676	0.7453	0.5223	0.5051
212	mobilenet	0.18	0.995	0.2	4064	6963	130	729	0.674	0.7868	0.867	0.718	0.695	0.704	0.807	0.75	0.684	0.743	1369	1329	1239	2058	0.684	0.7456	0.541	0.505
213	mobilenet	0.12	0.995	0.35	3948	6764	246	530	0.658	0.7748	0.871	0.73	0.699	0.714	0.805	0.751	0.678	0.743	1433	1404	1187	1934	0.679	0.7444	0.5294	0.505
214	mobilenet	0.15	0.995	0.4	3858	6731	336	497	0.64	0.7725	0.873	0.732	0.696	0.715	0.8	0.751	0.667	0.743	1508	1418	1172	1915	0.668	0.744	0.5143	0.5049
215	mobilenet	0.13	0.995	0.15	4036	6806	158	572	0.678	0.7783	0.87	0.727	0.704	0.713	0.81	0.751	0.691	0.744	1352	1382	1194	1954	0.691	0.7456	0.5476	0.5048
216	mobilenet	0.18	0.995	0.1	4064	6963	130	729	0.68	0.7879	0.869	0.717	0.701	0.705	0.81	0.75	0.69	0.744	1344	1322	1214	2051	0.687	0.7445	0.5488	0.5048
217	mobilenet	0.12	0.995	0.3	3948	6764	246	530	0.66	0.7749	0.872	0.73	0.702	0.714	0.806	0.751	0.68	0.743	1424	1403	1178	1933	0.681	0.7446	0.5324	0.5048
218	mobilenet	0.22	0.995	0.1	5655	6575	1461	341	0.79	0.7408	0.773	0.764	0.586	0.702	0.778	0.754	0.673	0.721	881	1616	2342	1957	0.688	0.7216	0.5631	0.5046
219	mobilenet	0.13	0.995	0.2	4036	6806	158	572	0.675	0.7775	0.869	0.727	0.701	0.712	0.808	0.751	0.687	0.743	1365	1387	1207	1959	0.688	0.7448	0.5436	0.5046
220	resnet50_1	0.2	0.995	0.1	4970	7056	776	822	0.735	0.7799	0.811	0.725	0.62	0.689	0.789	0.749	0.673	0.732	1111	1372	1887	2194	0.678	0.7345	0.546	0.5045
221	resnet50_1	0.14	0.995	0.1	4890	6905	696	671	0.733	0.7717	0.816	0.733	0.629	0.697	0.791	0.75	0.677	0.732	1119	1423	1815	2094	0.681	0.7342	0.5494	0.5045
222	mobilenet	0.12	0.995	0.4	3948	6764	246	530	0.657	0.7741	0.871	0.73	0.698	0.713	0.804	0.751	0.676	0.743	1440	1408	1194	1938	0.677	0.7438	0.5272	0.5045
223	mobilenet	0.11	0.995	0.4	3858	6736	336	502	0.645	0.773	0.874	0.731	0.701	0.715	0.802	0.75	0.672	0.743	1488	1415	1152	1917	0.673	0.7442	0.5194	0.5037
224	mobilenet	0.13	0.995	0.25	4036	6806	158	572	0.673	0.7765	0.869	0.727	0.699	0.711	0.808	0.75	0.686	0.742	1371	1393	1213	1965	0.686	0.7439	0.5418	0.5036
225	mobilenet	0.18	0.995	0.25	4064	6963	130	729	0.672	0.785	0.867	0.718	0.694	0.703	0.807	0.749	0.683	0.742	1374	1340	1244	2069	0.683	0.744	0.5396	0.5026
226	mobilenet	0.13	0.995	0.3	4036	6806	158	572	0.672	0.7757	0.868	0.727	0.698	0.711	0.807	0.749	0.685	0.742	1375	1398	1217	1970	0.685	0.7431	0.5406	0.5025
227	mobilenet	0.8	0.985	0.1	6819	7802	2625	1568	0.878	0.7988	0.728	0.703	0.54	0.638	0.768	0.741	0.668	0.71	513	1254	3138	2822	0.709	0.7186	0.6061	0.5022
228	mobilenet	0.2	0.995	0.1	3928	6972	266	738	0.663	0.7887	0.875	0.713	0.708	0.705	0.809	0.748	0.685	0.745	1414	1317	1148	2055	0.685	0.747	0.5382	0.5022
229	mobilenet	0.18	0.995	0.3	4064	6963	130	729	0.672	0.7844	0.867	0.717	0.693	0.702	0.807	0.748	0.682	0.741	1377	1344	1247	2073	0.682	0.7433	0.5386	0.5018
230	mobilenet	0.2	0.995	0.15	3928	6972	266	738	0.66	0.7878	0.875	0.714	0.705	0.704	0.808	0.748	0.682	0.744	1424	1323	1158	2061	0.683	0.7461	0.5351	0.5015
231	mobilenet	0.13	0.995	0.35	4036	6806	158	572	0.67	0.7748	0.868	0.727	0.696	0.71	0.806	0.749	0.683	0.741	1383	1404	1225	1976	0.683	0.7422	0.5381	0.5014
232	mobilenet	0.14	0.995	0.1	4010	6840	184	606	0.675	0.7785	0.872	0.723	0.706	0.71	0.81	0.749	0.689	0.742	1364	1381	1180	1987	0.689	0.744	0.5466	0.5012
233	mobilenet	0.13	0.995	0.4	4036	6806	158	572	0.669	0.7745	0.867	0.727	0.695	0.709	0.806	0.749	0.682	0.74	1389	1406	1231	1978	0.682	0.7419	0.5361	0.5012
234	mobilenet	0.2	0.995	0.2	3928	6972	266	738	0.658	0.7868	0.874	0.714	0.702	0.704	0.807	0.748	0.679	0.743	1435	1329	1169	2067	0.68	0.7452	0.5318	0.5011
235	mobilenet	0.18	0.995	0.35	4064	6963	130	729	0.67	0.7838	0.866	0.717	0.691	0.702	0.806	0.748	0.681	0.74	1384	1348	1254	2077	0.681	0.7427	0.5364	0.501
236	mobilenet	0.25	0.995	0.1	5593	6538	1399	304	0.783	0.7368	0.767	0.764	0.587	0.703	0.778	0.752	0.671	0.719	909	1641	2308	1945	0.685	0.7199	0.5589	0.5009
237	resnet50_1	0.6	0.99	0.1	5339	7637	1145	1403	0.78	0.8104	0.801	0.69	0.613	0.662	0.795	0.742	0.686	0.728	922	1182	2065	2585	0.697	0.736	0.5811	0.5006
238	mobilenet	0.14	0.995	0.2	4010	6840	184	606	0.668	0.7765	0.871	0.724	0.699	0.708	0.807	0.748	0.683	0.741	1393	1393	1209	1999	0.683	0.7421	0.5375	0.5001
239	resnet50_1	0.17	0.995	0.1	4928	7030	734	796	0.74	0.7759	0.817	0.724	0.629	0.688	0.794	0.747	0.68	0.729	1092	1397	1826	2193	0.685	0.732	0.5569	0.5001
240	mobilenet	0.18	0.995	0.4	4064	6963	130	729	0.668	0.783	0.866	0.717	0.69	0.701	0.805	0.747	0.69	0.74	1391	1353	1261	2082	0.679	0.742	0.5342	0.5
241	resnet50_1	0.23	0.995	0.1	5088	7090	894	856	0.751	0.7774	0.808	0.723	0.619	0.683	0.791	0.746	0.678	0.727	1046	1388	1940	2244	0.685	0.7304	0.559	0.4999
242	mobilenet	0.14	0.995	0.15	4010	6840	184	606	0.672	0.777	0.871	0.723	0.702	0.708	0.809	0.748	0.687	0.741	1377	1390	1193	1996	0.687	0.7426	0.5426	0.4998
243	resnet50_1	0.18	0.995	0.1	4928	6967	734	733	0.742	0.7722	0.818	0.727	0.632	0.691	0.795	0.747	0.683	0.729	1080	1420	1814	2153	0.687	0.7316	0.56	0.4997
244	mobilenet	0.14	0.995	0.25	4010	6840	184	606	0.667	0.7759	0.869	0.724	0.697	0.707	0.806	0.748	0.682	0.74	1398	1392	1214	2003	0.682	0.7415	0.5359	0.4995
245	mobilenet	0.23	0.995	0.1	4140	7137	54	903	0.692	0.795	0.869	0.704	0.701	0.694	0.814	0.746	0.697	0.741	1240	1378	1236	2181	0.697	0.7447	0.5612	0.4995
246	mobilenet	0.2	0.995	0.25	3928	6972	266	738	0.655	0.7854	0.873	0.714	0.7	0.702	0.805	0.747	0.677	0.741	1496	1328	1180	2076	0.677	0.744	0.5283	0.4993

100 Appendix E. Results of testing Overlap, NMS threshold and IOU threshold for the final system

		TOP (Total Predictions)				TPR (True positives rate - recall)				TNR (True negative rate - specificity)				PPV (Precision - positive predictive value)				Accuracy				F1		FN		FP		AUPR		Youden				
1	model	overlap	nms_th	metric_th	inflam	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ	infla	squ
300	mobilenet	0.21	0.995	0.15	4087	7040	107	806	0.683	0.7836	0.87	0.707	0.701	0.694	0.812	0.742	0.692	0.736	1328	1349	1221	2155	0.692	0.7387	0.5533	0.4904								
301	mobilenet	0.21	0.995	0.3	4087	7040	107	806	0.676	0.7818	0.868	0.708	0.693	0.692	0.809	0.742	0.684	0.734	1360	1360	1253	2166	0.685	0.7371	0.5432	0.4898								
302	resnet50_1	0.26	0.995	0.1	5174	7326	980	1092	0.748	0.7833	0.804	0.706	0.607	0.667	0.788	0.739	0.67	0.72	1055	1351	2035	2443	0.678	0.7249	0.5521	0.4898								
303	mobilenet	0.26	0.995	0.2	4226	7270	32	1036	0.689	0.7937	0.862	0.696	0.684	0.681	0.81	0.74	0.686	0.733	1305	1286	1337	2322	0.686	0.7372	0.5507	0.4897								
304	mobilenet	0.25	0.995	0.35	3994	7179	200	945	0.664	0.7907	0.873	0.699	0.697	0.687	0.809	0.741	0.68	0.735	1411	1305	1211	2250	0.68	0.7386	0.5363	0.4897								
305	mobilenet	0.27	0.995	0.1	4330	7273	136	1039	0.701	0.7945	0.856	0.695	0.679	0.681	0.809	0.74	0.69	0.733	1254	1281	1390	2320	0.69	0.7378	0.557	0.4897								
306	mobilenet	0.26	0.995	0.25	4226	7270	32	1036	0.687	0.7932	0.861	0.696	0.682	0.68	0.809	0.74	0.685	0.732	1312	1289	1344	2325	0.685	0.7367	0.5485	0.4893								
307	resnet50_1	0.27	0.995	0.1	5217	7333	1023	1099	0.75	0.7822	0.801	0.707	0.603	0.665	0.787	0.739	0.669	0.719	1047	1358	2070	2457	0.677	0.7236	0.5518	0.4892								
308	mobilenet	0.25	0.995	0.4	3994	7179	200	945	0.663	0.7902	0.873	0.699	0.696	0.686	0.809	0.74	0.679	0.735	1413	1308	1213	2253	0.68	0.7382	0.5357	0.4891								
309	mobilenet	0.27	0.995	0.15	4330	7273	136	1039	0.695	0.7927	0.854	0.696	0.673	0.679	0.806	0.739	0.684	0.732	1278	1292	1414	2331	0.684	0.7361	0.5496	0.4888								
310	mobilenet	0.26	0.995	0.3	4226	7270	32	1036	0.686	0.7926	0.861	0.696	0.68	0.68	0.808	0.739	0.683	0.732	1319	1293	1351	2329	0.683	0.7361	0.5463	0.4887								
311	resnet50_1	0.13	0.995	0.1	3759	6741	435	507	0.625	0.7631	0.877	0.726	0.698	0.706	0.799	0.743	0.659	0.733	1572	1477	1137	1984	0.661	0.7344	0.5025	0.4887								
312	mobilenet	0.21	0.995	0.35	4087	7040	107	806	0.675	0.7809	0.867	0.708	0.693	0.691	0.808	0.741	0.684	0.733	1363	1366	1256	2172	0.684	0.7362	0.5424	0.4886								
313	mobilenet	0.27	0.995	0.2	4330	7273	136	1039	0.694	0.7916	0.854	0.696	0.672	0.679	0.806	0.739	0.683	0.731	1285	1299	1421	2338	0.683	0.7351	0.5477	0.4878								
314	mobilenet	0.26	0.995	0.35	4226	7270	32	1036	0.684	0.7918	0.86	0.696	0.679	0.679	0.807	0.739	0.682	0.731	1324	1298	1356	2334	0.684	0.7354	0.5447	0.4877								
315	mobilenet	0.22	0.995	0.1	4066	7160	128	926	0.683	0.7902	0.873	0.697	0.705	0.688	0.814	0.74	0.694	0.736	1329	1308	1201	2234	0.694	0.7391	0.5557	0.4876								
316	mobilenet	0.21	0.995	0.4	4087	7040	107	806	0.674	0.7807	0.873	0.707	0.691	0.691	0.808	0.741	0.683	0.733	1368	1371	1261	2177	0.683	0.7354	0.5408	0.4875								
317	resnet50_1	0.25	0.995	0.1	5118	7300	924	1066	0.746	0.781	0.807	0.706	0.611	0.667	0.789	0.738	0.672	0.72	1067	1365	1991	2431	0.678	0.724	0.5526	0.4872								
318	mobilenet	0.26	0.995	0.4	4226	7270	32	1036	0.684	0.7913	0.86	0.696	0.679	0.679	0.807	0.739	0.681	0.731	1325	1301	1357	2337	0.681	0.7349	0.5445	0.4871								
319	mobilenet	0.27	0.995	0.25	4330	7273	136	1039	0.692	0.7907	0.854	0.696	0.67	0.678	0.805	0.738	0.681	0.73	1292	1305	1428	2344	0.681	0.7342	0.5455	0.4867								
320	resnet50_1	0.11	0.995	0.1	3764	6684	430	450	0.62	0.7573	0.875	0.729	0.691	0.706	0.795	0.742	0.653	0.731	1504	1513	1164	1963	0.655	0.7318	0.4946	0.4865								
321	mobilenet	0.28	0.995	0.1	4317	7320	123	1086	0.698	0.7948	0.857	0.692	0.678	0.677	0.809	0.738	0.688	0.731	1267	1279	1390	2365	0.688	0.7359	0.5547	0.4865								
322	mobilenet	0.22	0.995	0.15	4066	7160	128	926	0.68	0.7886	0.872	0.698	0.701	0.687	0.813	0.739	0.691	0.734	1342	1318	1214	2244	0.684	0.7362	0.5517	0.4863								
323	mobilenet	0.27	0.995	0.3	4330	7273	136	1039	0.691	0.7902	0.853	0.696	0.67	0.677	0.805	0.738	0.68	0.729	1295	1308	1431	2347	0.68	0.7337	0.5446	0.4862								
324	mobilenet	0.28	0.995	0.2	4317	7320	123	1086	0.693	0.7934	0.856	0.693	0.673	0.676	0.807	0.738	0.683	0.73	1287	1288	1410	2374	0.683	0.7345	0.5487	0.4862								
325	mobilenet	0.28	0.995	0.15	4317	7320	123	1086	0.695	0.7937	0.856	0.692	0.675	0.676	0.808	0.737	0.685	0.73	1278	1286	1401	2372	0.685	0.7348	0.5514	0.4857								
326	mobilenet	0.28	0.995	0.25	4317	7320	123	1086	0.691	0.7927	0.855	0.693	0.672	0.675	0.806	0.737	0.681	0.729	1295	1292	1418	2378	0.682	0.7349	0.5483	0.4856								
327	mobilenet	0.27	0.995	0.4	4330	7273	136	1039	0.687	0.7891	0.852	0.696	0.666	0.676	0.803	0.738	0.676	0.728	1311	1315	1447	2354	0.677	0.7327	0.5396	0.4854								
328	resnet50_1	0.12	0.995	0.1	3736	6700	458	466	0.617	0.7579	0.876	0.727	0.692	0.705	0.796	0.742	0.652	0.731	1607	1509	1149	1975	0.655	0.7316	0.4931	0.4854								
329	mobilenet	0.22	0.995	0.2	4066	7160	128	926	0.679	0.7875	0.871	0.698	0.7	0.686	0.812	0.739	0.689	0.733	1348	1325	1220	2251	0.689	0.7365	0.55	0.4852								
330	mobilenet	0.27	0.995	0.35	4330	7273	136	1039	0.69	0.7892	0.853	0.696	0.668	0.676	0.804	0.738	0.679	0.729	1301	1314	1437	2353	0.679	0.7328	0.5428	0.4851								
331	mobilenet	0.22	0.995	0.25	4066	7160	128	926	0.676	0.7865	0.871	0.698	0.697	0.685	0.811	0.738	0.687	0.732	1358	1331	1230	2257	0.687	0.7356	0.5469	0.4844								
332	mobilenet	0.29	0.995	0.1	4301	7353	107	1119	0.704	0.796	0.861	0.698	0.686	0.675	0.814	0.737	0.695	0.73	1243	1272	1350	2391	0.695	0.734	0.5646	0.4844								
333	mobilenet	0.28	0.995	0.3	4317	7320	123	1086	0.69	0.7918	0.855	0.693	0.671	0.674	0.805	0.737	0.68	0.728	1299	1298	1422	2384	0.68	0.7331	0.5451	0.4843								
334	mobilenet	0.3	0.995	0.1	4158	7315	36	1081	0.683	0.7947	0.866	0.699	0.689	0.677	0.81	0.737	0.686	0.731	1330	1280	1294	2361	0.686	0.736	0.5487	0.4841								
335	mobilenet	0.28	0.995	0.35	4317	7320	123	1086	0.689	0.7911	0.854	0.693	0.669	0.674	0.805	0.736	0.679	0.728	1304	1302	1427	2388	0.679	0.7325	0.5436	0.4837								
336	mobilenet	0.3	0.995	0.15	4158	7315	36	1081	0.677	0.7934	0.864	0.693	0.69	0.683	0.796	0.737	0.68	0.73	1353	1288	1317	2369	0.68	0.7348	0.5415	0.4836								
337	mobilenet	0.22	0.995	0.3	4066	7160	128	926	0.675	0.7857	0.87	0.698	0.697	0.684	0.811	0.738	0.686	0.731	1362	1336	1234	2262	0.686	0.7349	0.5457	0.4834								
338	mobilenet	0.3	0.995	0.2	4158	7315	36	1081	0.674	0.7926	0.863	0.691	0.68	0.675	0.806	0.736	0.677	0.729	1368	1293	1332	2374	0.677	0.734	0.5368	0.4834								
339	mobilenet	0.28	0.995	0.4	4317	7320	123	1086	0.687	0.7903	0.854	0.693	0.668	0.673	0.804	0.736	0.677	0.727	1311	1307	1434	2393	0.678	0.7317	0.5414	0.4829								
340	resnet50_1	0.16	0.995	0.1	3854	6781	340	547	0.634	0.7605	0.872	0.722	0.689	0.699	0.799	0.74	0.66	0.729	1537	1493	1197													

E.2 Complete table

Since the table has more than 2800 rows, it was not possible to include this table entirely, and for this reason, it is divided in four tables with 100 rows each.

102 Appendix E. Results of testing Overlap, NMS threshold and IOU threshold for the final system

1	TOP (Total Predictions)				TPR (True positives rate - recall)				TNR (True negative rate - specificity)				PPV (Precision - positive predictive value)				Accuracy	F1	FN	FP	AUPR				Youden			
	model	overl	nms_th	metric	inflar	squar	inflar	squar	inflar	squar	inflar	squar	inflar	squar	inflar	squar					inflar	squar	inflar	squar	inflar	squar	inflar	squar
3	mobilenet	0.28	0.97	0.1	4351	5066	157	1168	0.729	0.74	0.834	0.911	0.7028	0.911	0.797	0.827	0.7157	0.817	1136	1620	1293	452	0.716	0.825	0.563	0.6612		
4	mobilenet	0.28	0.975	0.1	4624	5100	430	1134	0.766	0.742	0.821	0.9186	0.6944	0.906	0.802	0.827	0.7283	0.816	983	1611	1413	477	0.73	0.824	0.587	0.6602		
5	mobilenet	0.28	0.965	0.1	4012	5025	182	1209	0.68	0.736	0.848	0.924	0.7111	0.913	0.789	0.824	0.6953	0.815	1341	1644	1159	435	0.696	0.825	0.529	0.6587		
6	mobilenet	0.24	0.97	0.1	4327	5011	133	1223	0.724	0.734	0.834	0.924	0.7016	0.913	0.796	0.825	0.7126	0.814	1158	1659	1291	436	0.713	0.823	0.558	0.6581		
7	mobilenet	0.24	0.975	0.1	4586	5039	392	1195	0.759	0.735	0.822	0.923	0.6939	0.91	0.8	0.826	0.7248	0.813	1012	1651	1404	456	0.726	0.822	0.581	0.6574		
8	mobilenet	0.23	0.97	0.1	4289	5001	95	1233	0.719	0.733	0.836	0.9243	0.703	0.913	0.795	0.824	0.7108	0.813	1179	1667	1274	434	0.711	0.823	0.555	0.6569		
9	mobilenet	0.23	0.975	0.1	4568	5034	374	1200	0.757	0.734	0.823	0.923	0.6948	0.909	0.8	0.825	0.7245	0.812	1020	1657	1394	457	0.726	0.822	0.58	0.6562		
10	mobilenet	0.13	0.975	0.1	4390	4963	196	1271	0.738	0.729	0.834	0.9273	0.7052	0.916	0.8	0.824	0.7213	0.812	1098	1690	1294	419	0.722	0.822	0.572	0.6561		
11	mobilenet	0.14	0.975	0.1	4419	4967	225	1267	0.743	0.729	0.833	0.9268	0.7049	0.915	0.801	0.824	0.7233	0.812	1079	1689	1304	422	0.724	0.822	0.576	0.6559		
12	mobilenet	0.27	0.97	0.1	4551	5031	157	1203	0.727	0.734	0.833	0.921	0.7012	0.91	0.796	0.824	0.7141	0.813	1143	1656	1300	453	0.714	0.822	0.56	0.6554		
13	mobilenet	0.28	0.96	0.1	3693	4971	501	1263	0.632	0.731	0.861	0.9245	0.7173	0.917	0.779	0.822	0.6717	0.813	1545	1678	1044	415	0.674	0.824	0.493	0.6553		
14	mobilenet	0.13	0.97	0.1	4136	4927	58	1307	0.703	0.726	0.846	0.9272	0.7133	0.919	0.795	0.823	0.708	0.811	1245	1707	1187	400	0.708	0.822	0.549	0.6553		
15	resnet50_1	0.29	0.97	0.1	4029	5555	165	679	0.693	0.778	0.856	0.8772	0.7213	0.873	0.799	0.826	0.7068	0.823	1288	1385	1123	706	0.707	0.825	0.549	0.6551		
16	mobilenet	0.14	0.97	0.1	4183	4937	11	1297	0.711	0.727	0.844	0.9282	0.7134	0.918	0.797	0.823	0.7124	0.811	1210	1703	1199	406	0.712	0.822	0.556	0.6551		
17	mobilenet	0.26	0.975	0.1	4520	5089	326	1145	0.76	0.738	0.83	0.9165	0.7051	0.904	0.806	0.824	0.7315	0.813	1007	1631	1333	486	0.732	0.821	0.59	0.6548		
18	resnet50_1	0.29	0.975	0.1	4259	5612	65	622	0.726	0.781	0.846	0.8734	0.7152	0.868	0.805	0.826	0.7207	0.822	1148	1364	1213	746	0.721	0.824	0.573	0.6546		
19	mobilenet	0.29	0.97	0.1	4361	5043	167	1191	0.732	0.735	0.834	0.9197	0.7042	0.908	0.799	0.824	0.7179	0.812	1123	1653	1290	462	0.718	0.822	0.567	0.6545		
20	mobilenet	0.19	0.97	0.1	4289	4980	95	1254	0.723	0.73	0.838	0.9247	0.7074	0.913	0.798	0.823	0.7153	0.811	1160	1685	1255	431	0.715	0.822	0.562	0.6544		
21	resnet50_1	0.17	0.975	0.1	4082	5484	112	750	0.711	0.772	0.859	0.8824	0.731	0.877	0.807	0.825	0.7211	0.821	1210	1423	1098	673	0.721	0.825	0.57	0.6542		
22	mobilenet	0.23	0.965	0.1	3974	4966	220	1268	0.673	0.729	0.85	0.925	0.7106	0.915	0.787	0.822	0.6915	0.812	1370	1689	1150	421	0.692	0.822	0.523	0.6541		
23	mobilenet	0.27	0.975	0.1	4628	5065	434	1169	0.765	0.736	0.82	0.9184	0.693	0.905	0.801	0.824	0.727	0.812	987	1648	1421	479	0.729	0.821	0.585	0.654		
24	mobilenet	0.24	0.965	0.1	4012	4970	182	1264	0.68	0.729	0.848	0.924	0.7104	0.915	0.789	0.822	0.6946	0.811	1344	1688	1162	424	0.695	0.821	0.528	0.6539		
25	mobilenet	0.19	0.975	0.1	4532	5014	338	1220	0.759	0.732	0.829	0.9222	0.7026	0.91	0.804	0.824	0.7298	0.811	1010	1673	1348	453	0.731	0.821	0.588	0.6538		
26	mobilenet	0.26	0.97	0.1	4261	5065	67	1169	0.723	0.737	0.842	0.9172	0.712	0.907	0.8	0.823	0.7177	0.813	1160	1642	1227	473	0.718	0.822	0.565	0.6538		
27	mobilenet	0.28	0.98	0.1	4746	5157	552	1077	0.779	0.742	0.816	0.9114	0.6886	0.897	0.803	0.825	0.7311	0.813	926	1606	1478	529	0.734	0.822	0.595	0.6538		
28	mobilenet	0.22	0.97	0.1	4278	5012	84	1222	0.721	0.732	0.839	0.9217	0.7069	0.91	0.797	0.823	0.7139	0.811	1170	1671	1254	449	0.714	0.821	0.56	0.6536		
29	resnet50_1	0.17	0.97	0.1	3874	5456	320	778	0.678	0.77	0.866	0.883	0.7336	0.88	0.799	0.824	0.7045	0.821	1352	1435	1032	657	0.706	0.825	0.544	0.6536		
30	resnet50_1	0.29	0.965	0.1	3704	5490	490	744	0.642	0.773	0.868	0.8808	0.7271	0.877	0.789	0.824	0.6819	0.822	1501	1417	1011	673	0.685	0.825	0.511	0.6535		
31	mobilenet	0.29	0.975	0.1	4640	5077	446	1157	0.769	0.736	0.822	0.9172	0.6953	0.904	0.804	0.824	0.7304	0.811	968	1645	1414	488	0.732	0.822	0.591	0.6533		
32	mobilenet	0.14	0.965	0.1	3877	4896	317	1338	0.677	0.723	0.858	0.9301	0.7212	0.921	0.79	0.821	0.6929	0.81	1398	1726	1091	388	0.694	0.822	0.524	0.6532		
33	mobilenet	0.18	0.975	0.1	4484	4960	290	1274	0.751	0.727	0.83	0.9253	0.7023	0.914	0.802	0.823	0.7257	0.81	1045	1702	1335	428	0.727	0.82	0.581	0.6532		
34	mobilenet	0.22	0.975	0.1	4543	5056	349	1178	0.758	0.735	0.827	0.9185	0.6995	0.906	0.803	0.824	0.7275	0.811	1016	1655	1365	477	0.729	0.82	0.585	0.6531		
35	mobilenet	0.29	0.965	0.1	4027	4998	167	1236	0.684	0.731	0.849	0.9217	0.7124	0.912	0.791	0.822	0.698	0.812	1325	1676	1158	440	0.698	0.822	0.533	0.6529		
36	mobilenet	0.18	0.97	0.1	4217	4940	23	1294	0.714	0.726	0.842	0.927	0.7097	0.916	0.797	0.822	0.7117	0.81	1201	1709	1224	415	0.712	0.821	0.555	0.6529		
37	mobilenet	0.21	0.97	0.1	4227	5033	33	1201	0.719	0.733	0.844	0.9197	0.7135	0.908	0.8	0.822	0.7163	0.811	1178	1663	1211	462	0.716	0.821	0.563	0.6525		
38	mobilenet	0.21	0.975	0.1	4485	5064	291	1170	0.755	0.735	0.832	0.9174	0.7059	0.905	0.806	0.823	0.7296	0.811	1028	1652	1319	482	0.73	0.822	0.587	0.6524		
39	mobilenet	0.6	0.97	0.1	4694	5410	500	824	0.763	0.762	0.814	0.8903	0.6815	0.878	0.797	0.825	0.7198	0.816	995	1483	1495	659	0.722	0.822	0.577	0.6524		
40	mobilenet	0.8	0.965	0.1	5202	5786	1008	448	0.776	0.788	0.771	0.8646	0.6259	0.849	0.773	0.827	0.6931	0.817	938	1323	1496	875	0.701	0.818	0.547	0.6523		
41	mobilenet	0.19	0.965	0.1	3955	4943	239	1291	0.675	0.726	0.853	0.9255	0.7153	0.916	0.79	0.821	0.6943	0.81	1365	1706	1126	415	0.695	0.821	0.527	0.6523		
42	mobilenet	0.13	0.965	0.1	3815	4892	379	1342	0.658	0.722	0.86	0.9297	0.7232	0.921	0.788	0.82	0.689	0.81	1435	1730	1056	388	0.691	0.822	0.518	0.6522		
43	mobilenet	0.27	0.965	0.1	4018	4996	176	1238	0.681	0.731	0.848	0.9214	0.7111	0.912	0.789	0.821	0.6958	0.811	1337	1679	1161	441	0.696	0.821	0.53	0.6521		
44	mobilenet	0.17	0.975	0.1	4473	4988	279	1246	0.752	0.729																		

1	2	model	overfit	nms	metric	TOP (Total Predictions)		IGT - Predict		TPR (True positives rate - recall)		TNR (True negative rate - specificity)		PPV (Precision - positive predictive value)		Accuracy	F1		FN		FP		AUPR		Youden	
						inflat	squar	inflat	squar	inflat	squar	inflat	squar	inflat	squar		inflat	squar	inflat	squar	inflat	squar	inflat	squar	inflat	squar
1000	mobilenet_v1	0.18	0.96	0.3	2906	5442	1288	792	0.525	0.761	0.907	0.8734	0.7581	0.872	0.77	0.814	0.6206	0.813	1991	1488	703	696	0.642	0.817	0.432	0.6347
1001	mobilenet_v1	0.28	0.965	0.3	3286	5652	908	582	0.584	0.777	0.891	0.8575	0.7456	0.857	0.783	0.815	0.6551	0.815	1744	1389	836	807	0.665	0.817	0.476	0.6347
1002	mobilenet_v1	0.03	0.965	0.1	2947	5321	1247	913	0.554	0.752	0.916	0.8822	0.7886	0.882	0.785	0.812	0.6509	0.812	1870	1543	623	630	0.671	0.817	0.47	0.6347
1003	resnet50_v1	0.05	0.97	0.1	3024	5210	1170	1024	0.564	0.743	0.911	0.892	0.7817	0.889	0.785	0.812	0.655	0.809	1830	1604	600	580	0.673	0.816	0.475	0.6347
1004	resnet50_v1	0.09	0.96	0.1	2602	5137	1592	1097	0.496	0.737	0.929	0.8973	0.8002	0.895	0.771	0.811	0.6127	0.809	2112	1637	520	540	0.648	0.816	0.425	0.6347
1005	mobilenet_v1	0.22	0.97	0.35	3414	5673	780	561	0.61	0.779	0.889	0.856	0.7499	0.856	0.791	0.816	0.673	0.815	1634	1380	854	819	0.68	0.817	0.5	0.6346
1006	mobilenet_v1	0.05	0.965	0.1	3641	4749	553	1485	0.643	0.703	0.874	0.932	0.7407	0.922	0.791	0.809	0.6884	0.798	1497	1854	369	699	0.682	0.817	0.517	0.6346
1007	mobilenet_v1	0.25	0.965	0.4	3079	5649	1115	585	0.555	0.777	0.902	0.8573	0.7564	0.858	0.78	0.815	0.6405	0.816	1865	1388	750	803	0.656	0.818	0.457	0.6346
1008	mobilenet_v1	0.29	0.975	0.4	3803	5805	391	429	0.668	0.787	0.873	0.8473	0.7363	0.845	0.802	0.816	0.7003	0.815	1394	1326	1003	897	0.702	0.816	0.541	0.6346
1009	mobilenet_v1	0.12	0.98	0.35	3635	5626	559	608	0.653	0.774	0.885	0.8606	0.7532	0.858	0.804	0.816	0.6995	0.814	1456	1409	897	801	0.703	0.816	0.538	0.6346
1010	mobilenet_v1	0.19	0.97	0.35	3416	5627	778	607	0.613	0.775	0.891	0.8596	0.7532	0.859	0.793	0.815	0.6762	0.815	1621	1403	843	796	0.683	0.817	0.504	0.6346
1011	mobilenet_v1	0.19	0.96	0.1	2903	5454	1291	780	0.541	0.763	0.914	0.8712	0.7809	0.873	0.78	0.813	0.6389	0.814	1927	1475	636	695	0.661	0.818	0.455	0.6346
1012	mobilenet_v1	0.07	0.97	0.1	3327	5474	867	760	0.62	0.764	0.904	0.8707	0.7812	0.87	0.802	0.814	0.6911	0.813	1595	1472	728	712	0.7	0.817	0.523	0.6346
1013	mobilenet_v1	0.29	0.98	0.3	3894	5852	300	382	0.686	0.791	0.872	0.8437	0.7393	0.842	0.808	0.817	0.7119	0.816	1315	1304	1015	922	0.713	0.817	0.559	0.6345
1014	mobilenet_v1	0.05	0.975	0.1	4200	4837	6	1397	0.723	0.709	0.849	0.926	0.7219	0.913	0.804	0.812	0.7224	0.798	1162	1817	1168	1020	0.722	0.817	0.572	0.6345
1015	mobilenet_v1	0.22	0.975	0.4	3639	5727	555	507	0.648	0.782	0.882	0.8525	0.7466	0.851	0.8	0.816	0.6937	0.815	1477	1359	922	852	0.697	0.817	0.53	0.6345
1016	mobilenet_v1	0.3	0.975	0.4	3664	5791	530	443	0.647	0.787	0.879	0.8478	0.7402	0.847	0.798	0.816	0.6903	0.816	1482	1330	952	887	0.693	0.817	0.526	0.6345
1017	mobilenet_v1	0.19	0.98	0.1	3713	5725	481	509	0.674	0.783	0.886	0.8518	0.7616	0.852	0.812	0.816	0.7153	0.816	1366	1355	885	846	0.718	0.817	0.56	0.6345
1018	mobilenet_v1	0.1	0.97	0.2	3188	5559	1006	675	0.586	0.771	0.904	0.8638	0.7713	0.864	0.791	0.815	0.6662	0.815	1735	1430	729	755	0.679	0.817	0.49	0.6345
1019	mobilenet_v1	0.23	0.965	0.4	3186	5644	1008	590	0.575	0.777	0.899	0.8579	0.7574	0.858	0.785	0.815	0.6539	0.815	1781	1393	773	803	0.666	0.817	0.475	0.6344
1020	mobilenet_v1	0.04	0.975	0.15	3420	5504	774	730	0.628	0.766	0.897	0.8688	0.7699	0.867	0.801	0.814	0.6916	0.813	1561	1461	787	731	0.699	0.816	0.524	0.6344
1021	mobilenet_v1	0.04	0.965	0.2	2975	5373	1219	861	0.547	0.756	0.909	0.8786	0.7711	0.877	0.779	0.813	0.64	0.812	1900	1522	681	661	0.659	0.816	0.456	0.6344
1022	mobilenet_v1	0.13	0.96	0.15	2903	5454	1291	780	0.535	0.763	0.912	0.8715	0.7733	0.872	0.777	0.813	0.6327	0.814	1498	1478	658	698	0.654	0.817	0.447	0.6344
1023	mobilenet_v1	0.08	0.98	0.1	5326	5593	932	641	0.813	0.764	0.796	0.8701	0.6648	0.852	0.801	0.818	0.7315	0.806	786	1459	1718	828	0.739	0.808	0.608	0.6344
1024	mobilenet_v1	0.09	0.98	0.1	3616	5631	578	603	0.664	0.775	0.892	0.859	0.7696	0.858	0.811	0.815	0.7127	0.815	1411	1400	833	797	0.717	0.817	0.555	0.6344
1025	mobilenet_v1	0.09	0.97	0.35	3314	5531	880	703	0.602	0.768	0.896	0.8668	0.7613	0.865	0.792	0.815	0.6721	0.813	1671	1449	791	746	0.681	0.816	0.498	0.6344
1026	resnet50_v1	0.04	0.975	0.1	3126	5263	1068	971	0.58	0.746	0.907	0.8881	0.778	0.884	0.79	0.813	0.6645	0.809	1762	1582	694	611	0.679	0.815	0.487	0.6344
1027	mobilenet_v1	0.27	0.98	0.25	3927	5833	267	401	0.685	0.789	0.867	0.8453	0.7313	0.843	0.804	0.816	0.7073	0.815	1322	1315	1055	914	0.708	0.816	0.552	0.6344
1028	mobilenet_v1	0.15	0.965	0.35	2984	5446	1210	788	0.54	0.761	0.905	0.8732	0.7594	0.871	0.775	0.814	0.6314	0.813	1928	1489	718	701	0.65	0.816	0.445	0.6343
1029	mobilenet_v1	0.1	0.98	0.15	3479	5624	715	610	0.64	0.775	0.896	0.8592	0.7718	0.859	0.806	0.815	0.6999	0.815	1509	1402	794	792	0.706	0.817	0.537	0.6343
1030	mobilenet_v1	0.11	0.965	0.25	3001	5477	1193	757	0.548	0.764	0.907	0.8703	0.7664	0.87	0.779	0.814	0.6393	0.813	1894	1471	701	714	0.657	0.817	0.455	0.6343
1031	mobilenet_v1	0.09	0.96	0.15	2830	5375	1364	859	0.53	0.757	0.918	0.8776	0.7848	0.878	0.778	0.813	0.6324	0.813	1973	1517	609	658	0.657	0.817	0.447	0.6343
1032	mobilenet_v1	0.19	0.98	0.25	3713	5725	481	509	0.668	0.782	0.883	0.8525	0.7541	0.851	0.808	0.816	0.7082	0.815	1394	1360	913	851	0.711	0.817	0.551	0.6343
1033	mobilenet_v1	0.22	0.965	0.2	3144	5595	1050	639	0.574	0.774	0.903	0.8608	0.7653	0.862	0.786	0.815	0.6558	0.815	1788	1412	738	773	0.669	0.818	0.476	0.6343
1034	mobilenet_v1	0.12	0.98	0.3	3635	5626	559	608	0.655	0.774	0.886	0.8603	0.7557	0.858	0.805	0.815	0.7017	0.814	1447	1409	888	801	0.705	0.816	0.541	0.6343
1035	mobilenet_v1	0.27	0.965	0.3	3297	5651	897	583	0.584	0.777	0.89	0.8574	0.7425	0.857	0.782	0.815	0.6536	0.815	1746	1391	849	808	0.663	0.817	0.474	0.6343
1036	mobilenet_v1	0.19	0.97	0.4	3416	5627	778	607	0.612	0.775	0.89	0.8596	0.7515	0.858	0.792	0.815	0.6746	0.814	1627	1405	849	798	0.682	0.816	0.502	0.6342
1037	mobilenet_v1	0.13	0.96	0.1	2888	5387	1306	847	0.544	0.758	0.918	0.8761	0.7895	0.877	0.782	0.813	0.6439	0.813	1914	1508	608	667	0.682	0.816	0.461	0.6342
1038	mobilenet_v1	0.2	0.98	0.3	3608	5739	586	495	0.647	0.783	0.885	0.8511	0.7525	0.851	0.802	0.816	0.696	0.816	1479	1352	893	857	0.7	0.817	0.533	0.6342
1039	resnet50_v1	0.6	0.975	0.1	3860	5832	334	402	0.675	0.789	0.87	0.845	0.7329	0.844	0.802	0.816	0.7025	0.816	1365	1314	1031	912	0.704	0.816	0.544	0.6342
1040	mobilenet_v1	0.13	0.98	0.4	3720	5626	474	608	0.664	0.774	0.88	0.8607	0.7484	0.857	0.805	0.815	0.7036									

104 Appendix E. Results of testing Overlap, NMS threshold and IOU threshold for the final system

1	2	model	overl	nms	metric	TOP (Total Predictions)		IGT - Predict		TPR (True positives rate - recall)		TNR (True negative rate - specificity)		PPV (Precision - positive predictive value)		Accuracy	F1	FN	FP	AUPR	Youden					
						inflar	squar	inflar	squar	inflar	squar	inflar	squar	inflar	squar											
2000	mobilenet_v1	0.23	0.95	0.3	2576	5274	1618	960	0.471	0.737	0.919	0.8746	0.7663	0.871	0.758	0.801	0.5852	0.798	2220	1641	602	681	0.618	0.804	0.39	0.6114
2001	mobilenet_v1	0.15	0.95	0.4	2421	5129	1773	1105	0.444	0.725	0.925	0.8863	0.7995	0.881	0.751	0.8	0.5633	0.796	2331	1714	558	609	0.607	0.803	0.369	0.6113
2002	resnet50_v1	0.19	0.985	0.1	3439	5594	755	640	0.631	0.76	0.898	0.8513	0.7694	0.847	0.805	0.804	0.6933	0.801	1548	1497	793	857	0.7	0.803	0.529	0.6112
2003	resnet50_v1	0.03	0.98	0.1	3158	5451	1036	783	0.586	0.749	0.809	0.8619	0.7783	0.857	0.795	0.803	0.6687	0.799	1736	1563	700	780	0.682	0.803	0.495	0.6111
2004	mobilenet_v1	0.19	0.95	0.25	2520	5201	1674	1033	0.46	0.731	0.921	0.8801	0.7663	0.876	0.754	0.8	0.5752	0.797	2263	1677	589	644	0.613	0.804	0.381	0.6111
2005	resnet50_v1	0.13	0.985	0.1	3335	5598	859	636	0.62	0.761	0.905	0.8503	0.7793	0.847	0.805	0.804	0.6904	0.802	1595	1491	736	855	0.7	0.804	0.525	0.6111
2006	mobilenet_v1	0.13	0.95	0.2	2529	5135	1665	1099	0.47	0.726	0.924	0.8852	0.779	0.881	0.759	0.799	0.586	0.796	2224	1709	550	610	0.624	0.804	0.394	0.6111
2007	mobilenet_v1	0.08	0.95	0.25	2483	5118	1711	1116	0.463	0.724	0.926	0.8868	0.7817	0.882	0.758	0.799	0.5814	0.795	2253	1719	542	603	0.622	0.803	0.389	0.611
2008	mobilenet_v1	0.04	0.955	0.4	2559	5170	1635	1064	0.467	0.728	0.919	0.8829	0.7659	0.878	0.756	0.8	0.5805	0.796	2234	1695	599	631	0.617	0.803	0.387	0.611
2009	mobilenet_v1	0.09	0.985	0.1	3642	5803	552	431	0.667	0.776	0.893	0.8347	0.7677	0.834	0.814	0.804	0.7136	0.804	1398	1395	846	964	0.717	0.805	0.559	0.6109
2010	mobilenet_v1	0.05	0.955	0.2	2519	5192	1675	1042	0.471	0.731	0.926	0.8802	0.7844	0.877	0.761	0.8	0.5887	0.797	2218	1679	543	637	0.628	0.804	0.397	0.6109
2011	mobilenet_v1	0.21	0.95	0.2	2515	5208	1679	1026	0.467	0.732	0.925	0.8791	0.7789	0.876	0.759	0.8	0.584	0.797	2235	1672	556	646	0.623	0.804	0.392	0.6109
2012	resnet50_v1	0.25	0.985	0.1	4364	5877	170	357	0.736	0.778	0.844	0.833	0.769	0.825	0.808	0.805	0.721	0.8	1109	1387	1279	1030	721	0.81	0.58	0.6108
2013	mobilenet_v1	0.02	0.955	0.3	2528	5119	1666	1115	0.469	0.724	0.924	0.8868	0.7789	0.882	0.76	0.799	0.5858	0.795	2225	1721	559	606	0.624	0.803	0.394	0.6108
2014	resnet50_v1	0.14	0.985	0.1	3373	5560	821	674	0.618	0.757	0.899	0.8536	0.7679	0.849	0.801	0.803	0.6846	0.8	1604	1514	783	840	0.693	0.803	0.517	0.6107
2015	mobilenet_v1	0.07	0.95	0.25	2472	5090	1722	1144	0.461	0.722	0.927	0.8891	0.782	0.884	0.758	0.799	0.58	0.795	2261	1735	539	591	0.621	0.803	0.388	0.6107
2016	mobilenet_v1	0.2	0.95	0.3	2456	5239	1738	995	0.452	0.734	0.924	0.8767	0.7712	0.873	0.753	0.8	0.5696	0.798	2300	1658	562	611	0.621	0.804	0.376	0.6107
2017	mobilenet_v1	0.1	0.95	0.2	2385	5158	1809	1076	0.447	0.728	0.93	0.8829	0.7853	0.88	0.754	0.799	0.5694	0.797	2321	1697	512	621	0.616	0.804	0.377	0.6107
2018	mobilenet_v1	0.3	0.95	0.3	2350	5300	1664	934	0.459	0.739	0.919	0.8719	0.7617	0.869	0.754	0.801	0.5732	0.799	2267	1629	603	695	0.611	0.804	0.379	0.6106
2019	mobilenet_v1	0.26	0.95	0.35	2566	5300	1628	934	0.469	0.739	0.92	0.872	0.7673	0.869	0.758	0.801	0.5825	0.798	2225	1630	597	696	0.618	0.804	0.39	0.6106
2020	mobilenet_v1	0.09	0.985	0.15	3642	5803	552	431	0.664	0.776	0.891	0.8349	0.7641	0.833	0.812	0.804	0.7103	0.803	1411	1399	859	968	0.714	0.804	0.555	0.6105
2021	mobilenet_v1	0.23	0.95	0.35	2576	5274	1618	960	0.47	0.736	0.919	0.8742	0.7651	0.87	0.758	0.801	0.5823	0.798	2223	1644	605	684	0.618	0.803	0.389	0.6105
2022	mobilenet_v1	0.09	0.95	0.35	2473	5139	1721	1095	0.456	0.726	0.924	0.8849	0.774	0.88	0.755	0.799	0.5742	0.795	2280	1711	559	616	0.615	0.803	0.381	0.6105
2023	mobilenet_v1	0.22	0.95	0.35	2530	5220	1664	1014	0.464	0.732	0.922	0.8784	0.77	0.874	0.757	0.8	0.5794	0.797	2246	1670	582	656	0.617	0.803	0.386	0.6105
2024	resnet50_v1	0.04	0.98	0.1	3189	5432	1005	802	0.589	0.747	0.806	0.8631	0.7745	0.858	0.794	0.802	0.6691	0.799	1724	1575	719	773	0.682	0.803	0.495	0.6104
2025	mobilenet_v1	0.3	0.985	0.25	3776	6005	418	229	0.666	0.789	0.879	0.8209	0.7397	0.819	0.806	0.805	0.7009	0.804	1401	1313	983	1084	0.703	0.804	0.544	0.6103
2026	mobilenet_v1	0.21	0.985	0.35	3754	5871	440	363	0.672	0.78	0.883	0.8303	0.7512	0.828	0.811	0.805	0.7096	0.803	1374	1372	934	1009	0.712	0.804	0.555	0.6103
2027	mobilenet_v1	0.14	0.95	0.35	2493	5128	1701	1106	0.457	0.724	0.922	0.8858	0.7682	0.881	0.754	0.799	0.5728	0.795	2279	1718	578	612	0.612	0.803	0.378	0.6102
2028	mobilenet_v1	0.29	0.985	0.35	3931	6026	263	208	0.689	0.79	0.872	0.82	0.7347	0.817	0.81	0.805	0.7109	0.804	1306	1308	1043	1100	0.712	0.804	0.561	0.6102
2029	mobilenet_v1	0.11	0.95	0.3	2425	5146	1769	1088	0.449	0.726	0.926	0.8839	0.7765	0.88	0.753	0.799	0.569	0.796	2311	1707	542	613	0.613	0.803	0.375	0.6101
2030	mobilenet_v1	0.05	0.965	0.4	2917	5390	1277	844	0.532	0.745	0.909	0.8651	0.7648	0.862	0.775	0.801	0.6275	0.799	1963	1590	686	746	0.648	0.803	0.441	0.6101
2031	mobilenet_v1	0.08	0.95	0.3	2483	5118	1711	1116	0.463	0.724	0.926	0.8863	0.7813	0.882	0.758	0.799	0.5811	0.795	2254	1722	543	606	0.622	0.803	0.389	0.6101
2032	mobilenet_v1	0.05	0.98	0.1	3442	5619	752	615	0.636	0.763	0.899	0.8473	0.7748	0.846	0.807	0.803	0.6985	0.802	1527	1480	775	865	0.705	0.804	0.535	0.6099
2033	mobilenet_v1	0.09	0.985	0.2	3642	5803	552	431	0.66	0.775	0.89	0.8349	0.7606	0.833	0.811	0.804	0.707	0.803	1424	1403	872	972	0.711	0.804	0.55	0.6099
2034	mobilenet_v1	0.13	0.95	0.25	2529	5135	1665	1099	0.467	0.725	0.923	0.8847	0.775	0.88	0.758	0.799	0.5831	0.795	2234	1714	569	615	0.621	0.803	0.39	0.6098
2035	mobilenet_v1	0.04	0.95	0.25	2404	5046	1790	1188	0.444	0.718	0.926	0.8921	0.7741	0.887	0.751	0.798	0.5641	0.793	2333	1760	544	572	0.609	0.802	0.37	0.6098
2036	mobilenet_v1	0.02	0.98	0.35	3453	5540	741	694	0.629	0.755	0.895	0.855	0.7637	0.849	0.802	0.803	0.6897	0.799	1557	1529	816	835	0.696	0.802	0.524	0.6097
2037	mobilenet_v1	0.27	0.95	0.35	2635	5266	1559	968	0.47	0.735	0.912	0.8749	0.748	0.87	0.753	0.8	0.5772	0.797	2223	1653	664	685	0.609	0.802	0.382	0.6097
2038	mobilenet_v1	0.1	0.95	0.25	2385	5158	1809	1076	0.444	0.727	0.928	0.8825	0.7807	0.879	0.753	0.799	0.566	0.796	2332	1701	523	625	0.612	0.803	0.373	0.6097
2039	mobilenet_v1	0.21	0.985	0.4	3754	5871	440	363	0.671	0.779	0.882	0.8302	0.7493	0.828	0.81	0.804	0.7079	0.803	1381	1375	941	1012	0.71	0.804	0.553	0.6096
2040	mobilenet_v1	0.05	0.955	0.25	2519	5192	1675	1042	0.469	0.73	0.925	0.8797	0.7809	0.876	0.76	0.799	0.586	0.796	2227	1684	552	642	0.625	0.803	0.394	0.6096
2041	mobilenet_v1	0.12	0.95																							

1	TOP (Total Predictions)				IGT - Predict		TPR (True positives rate - recall)		TNR (True negative rate - specificity)		PPV (Precision - positive predictive value)		Accuracy		F1		FN		FP		AUPR		Youden			
	model	over	nms	metric	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar	inflam	squar		
2745	mobilenet_v1	0.24	0.995	0.4	4181	7212	13	978	0.684	0.791	0.864	0.7	0.682	0.684	0.809	0.741	0.6551	0.734	1325	1302	1312	2280	0.685	0.738	0.548	0.4978
2746	mobilenet_v1	0.21	0.995	0.1	4087	7040	107	806	0.686	0.785	0.871	0.7062	0.7042	0.695	0.814	0.742	0.6951	0.737	1316	1341	1209	2147	0.695	0.74	0.557	0.4911
2750	resnet50_v1	0.1	0.995	0.1	3753	6562	441	328	0.616	0.753	0.872	0.7381	0.6882	0.715	0.792	0.745	0.6501	0.734	1611	1540	1170	1868	0.652	0.734	0.488	0.4911
2751	mobilenet_v1	0.26	0.995	0.15	4226	7179	31	1036	0.691	0.795	0.862	0.696	0.6853	0.682	0.81	0.74	0.6879	0.734	1298	1279	1300	2315	0.688	0.738	0.553	0.4908
2752	mobilenet_v1	0.25	0.995	0.3	3994	7279	200	945	0.664	0.791	0.873	0.6993	0.6975	0.687	0.809	0.741	0.6805	0.736	1408	1300	1230	2245	0.681	0.739	0.537	0.4908
2753	resnet50_v1	0.22	0.995	0.1	5061	7164	867	930	0.742	0.776	0.809	0.715	0.6147	0.675	0.789	0.741	0.6723	0.722	1083	1398	1590	2328	0.678	0.725	0.551	0.4907
2754	mobilenet_v1	0.21	0.995	0.2	4087	7040	107	806	0.688	0.783	0.869	0.7075	0.6978	0.693	0.811	0.742	0.6888	0.736	1342	1352	1235	2158	0.689	0.738	0.549	0.4906
2755	mobilenet_v1	0.21	0.995	0.25	4087	7040	107	806	0.677	0.782	0.868	0.708	0.6946	0.693	0.809	0.742	0.6857	0.735	1355	1356	1248	2162	0.686	0.738	0.545	0.4905
2756	mobilenet_v1	0.21	0.995	0.15	4087	7040	107	806	0.683	0.784	0.87	0.7068	0.7012	0.694	0.812	0.742	0.6922	0.736	1328	1349	1221	2155	0.692	0.739	0.553	0.4904
2757	mobilenet_v1	0.21	0.995	0.3	4087	7040	107	806	0.676	0.782	0.868	0.708	0.6934	0.692	0.809	0.742	0.6845	0.734	1360	1360	1253	2166	0.685	0.737	0.543	0.4898
2758	resnet50_v1	0.26	0.995	0.1	5174	7326	980	1092	0.748	0.783	0.804	0.7065	0.6067	0.667	0.788	0.739	0.6702	0.72	1055	1351	2035	2443	0.678	0.725	0.552	0.4898
2759	mobilenet_v1	0.26	0.995	0.2	4226	7279	32	1036	0.689	0.794	0.862	0.696	0.6836	0.681	0.81	0.74	0.6862	0.733	1305	1286	1337	2322	0.686	0.737	0.551	0.4897
2760	mobilenet_v1	0.25	0.995	0.35	3994	7179	200	945	0.664	0.791	0.873	0.699	0.6968	0.687	0.809	0.741	0.6798	0.735	1411	1305	1211	2250	0.68	0.739	0.536	0.4897
2761	mobilenet_v1	0.27	0.995	0.1	4330	7273	136	1039	0.701	0.795	0.856	0.6951	0.679	0.681	0.809	0.74	0.6898	0.733	1254	1281	1390	2320	0.69	0.738	0.557	0.4897
2762	mobilenet_v1	0.26	0.995	0.25	4226	7279	32	1036	0.687	0.793	0.861	0.6961	0.682	0.68	0.809	0.74	0.6846	0.732	1312	1289	1344	2325	0.687	0.737	0.548	0.4893
2763	resnet50_v1	0.27	0.995	0.1	5217	7333	1023	1099	0.75	0.782	0.801	0.707	0.6032	0.665	0.787	0.739	0.6688	0.719	1047	1358	2070	2457	0.677	0.724	0.552	0.4892
2764	mobilenet_v1	0.25	0.995	0.4	3994	7179	200	945	0.663	0.79	0.873	0.699	0.6963	0.686	0.809	0.74	0.6793	0.735	1413	1308	1213	2253	0.68	0.738	0.536	0.4891
2765	mobilenet_v1	0.27	0.995	0.15	4330	7273	136	1039	0.695	0.793	0.854	0.696	0.6734	0.679	0.806	0.739	0.6842	0.732	1278	1292	1414	2331	0.684	0.736	0.55	0.4888
2766	mobilenet_v1	0.26	0.995	0.3	4226	7279	32	1036	0.686	0.793	0.861	0.6962	0.6803	0.68	0.808	0.739	0.6829	0.732	1319	1293	1351	2329	0.683	0.736	0.546	0.4887
2767	resnet50_v1	0.13	0.995	0.1	3759	6741	435	507	0.625	0.763	0.877	0.7256	0.6975	0.706	0.799	0.743	0.6594	0.733	1572	1477	1137	1984	0.661	0.734	0.503	0.4887
2768	mobilenet_v1	0.21	0.995	0.35	4087	7040	107	806	0.675	0.781	0.867	0.7077	0.6927	0.691	0.808	0.741	0.6837	0.733	1363	1366	1256	2172	0.684	0.736	0.542	0.4886
2769	mobilenet_v1	0.27	0.995	0.2	4330	7273	136	1039	0.694	0.792	0.854	0.6962	0.6718	0.679	0.806	0.739	0.6825	0.731	1285	1299	1421	2338	0.683	0.735	0.548	0.4886
2770	mobilenet_v1	0.26	0.995	0.35	4226	7279	32	1036	0.684	0.792	0.86	0.6959	0.6791	0.679	0.807	0.739	0.6817	0.731	1324	1298	1356	2334	0.682	0.735	0.545	0.4877
2771	mobilenet_v1	0.22	0.995	0.1	4066	7160	128	926	0.683	0.79	0.873	0.6974	0.7046	0.688	0.814	0.74	0.6937	0.736	1299	1308	1201	2234	0.694	0.739	0.556	0.4876
2772	mobilenet_v1	0.21	0.995	0.4	4087	7040	107	806	0.674	0.78	0.867	0.7074	0.6915	0.691	0.808	0.741	0.6825	0.733	1368	1371	1261	2177	0.683	0.735	0.541	0.4875
2773	resnet50_v1	0.25	0.995	0.1	5118	7300	924	1066	0.746	0.781	0.807	0.7062	0.611	0.667	0.789	0.738	0.6716	0.72	1067	1365	1991	2431	0.673	0.724	0.553	0.4872
2774	mobilenet_v1	0.26	0.995	0.4	4226	7279	32	1036	0.684	0.791	0.86	0.6958	0.6789	0.679	0.807	0.739	0.6815	0.731	1325	1301	1357	2337	0.681	0.735	0.544	0.4871
2775	mobilenet_v1	0.27	0.995	0.25	4330	7273	136	1039	0.692	0.791	0.854	0.696	0.6702	0.678	0.805	0.738	0.6809	0.73	1292	1305	1428	2344	0.681	0.734	0.545	0.4867
2776	resnet50_v1	0.11	0.995	0.1	3764	6684	430	450	0.62	0.757	0.875	0.7292	0.6908	0.706	0.795	0.742	0.6534	0.731	1594	1513	1164	1965	0.655	0.732	0.495	0.4865
2777	mobilenet_v1	0.28	0.995	0.1	4317	7320	123	1086	0.698	0.795	0.857	0.6916	0.678	0.677	0.809	0.738	0.6878	0.731	1267	1279	1390	2365	0.688	0.736	0.555	0.4865
2778	mobilenet_v1	0.22	0.995	0.15	4066	7160	128	926	0.68	0.789	0.872	0.6977	0.7014	0.687	0.813	0.739	0.6906	0.734	1342	1318	1214	2244	0.691	0.738	0.552	0.4863
2779	mobilenet_v1	0.27	0.995	0.3	4330	7273	136	1039	0.691	0.79	0.853	0.696	0.6695	0.677	0.805	0.738	0.6802	0.729	1295	1308	1431	2347	0.68	0.734	0.545	0.4862
2780	mobilenet_v1	0.28	0.995	0.2	4317	7320	123	1086	0.693	0.793	0.856	0.6926	0.6734	0.676	0.807	0.738	0.6831	0.73	1287	1288	1401	2374	0.683	0.735	0.549	0.486
2781	mobilenet_v1	0.28	0.995	0.15	4317	7320	123	1086	0.695	0.794	0.856	0.6919	0.6755	0.676	0.808	0.737	0.6852	0.73	1278	1286	1410	2372	0.685	0.735	0.551	0.4857
2782	mobilenet_v1	0.28	0.995	0.25	4317	7320	123	1086	0.691	0.793	0.855	0.6929	0.6715	0.675	0.806	0.737	0.6812	0.729	1295	1292	1418	2378	0.681	0.734	0.546	0.4856
2783	mobilenet_v1	0.27	0.995	0.4	4330	7273	136	1039	0.687	0.789	0.852	0.6963	0.6658	0.676	0.803	0.738	0.6764	0.728	1311	1315	1447	2354	0.677	0.733	0.54	0.4854
2784	resnet50_v1	0.12	0.995	0.1	3736	6700	458	466	0.617	0.758	0.876	0.7274	0.6925	0.705	0.796	0.742	0.6525	0.731	1607	1509	1149	1975	0.655	0.732	0.493	0.4854
2785	mobilenet_v1	0.22	0.995	0.2	4066	7160	128	926	0.679	0.787	0.871	0.6978	0.7	0.686	0.812	0.739	0.6891	0.733	1348	1325	1220	2251	0.689	0.737	0.55	0.4852
2786	mobilenet_v1	0.27	0.995	0.35	4330	7273	136	1039	0.69	0.789	0.853	0.6959	0.6681	0.676	0.804	0.738	0.6788	0.729	1301	1314	1437	2353	0.679	0.733	0.543	0.4851
2787	mobilenet_v1	0.22	0.995	0.25	4066	7160	128	926	0.676	0.786	0.871	0.6979	0.6975	0.685	0.811	0.738	0.6867	0.732	1358	1331	1230	2257	0.687	0.736	0.547	0.4844
2788	mobilenet_v1	0.29	0.995	0.1	4301	7353	107	1119	0.704	0.796	0.861	0.6884	0.6861	0.675	0.814</											