**The Hashemite University**
**Faculty of Prince Al-Hussein Bin Abdallah II For Information Technology**
**Information Technology Department**

# Breast Cancer Detection System

**A project submitted**
**in partial fulfillment of the requirements for the**
**B.Sc. Degree in Data science and Artificial Intelligence**

**By**

Zaina Husam Hanna Al-Nimri (2136713)
Amro Ahmad Hussien Al-Mari (2137979)
Mohammad Loay Omar Silawy (2136861)

**Supervised by**

Dr. Esraa Ahmad Helael Al-Shdaifat

**Committee Member Names**

Dr. Zaher Salah

Dr. Eman Omar

**January 2025**

# CERTIFICATE

It is hereby certified that the project titled *Breast Cancer Detection System*, submitted by undersigned, in partial fulfillment of the award of the degree of "Bachelor in Data science and Artificial Intelligence" embodies original work done by them under my supervision.

All the analysis, design and system development have been accomplished by the undersigned. Moreover, this project has not been submitted to any other college or university.

| | |
|---|---|
| *Zaina* (2136713) | **Signature** |
| *Amro* (2137979) | **Signature** |
| *Mohammad* (2136861) | **Signature** |

# ABSTRACT

Breast cancer is the most common cancer type influencing ladies around the world **[1]**. Agreeing to the most recent insights, 2.3 million cases were analyzed in 2020. Among cancers influencing ladies, 25% of cases are breast cancer**[2]**. If cases are detected and treated early, two-fifths of these passing can potentially be survived, meaning that 274,000 ladies might be survived each year. Survival rates increment essentially when breast cancer is recognized at an early arrange (Stages 0-II) **[3]**. Mammography is the gold standard for identifying early signs of breast cancer, which can offer the detection of cancer at early stages. Consequently, building a breast cancer detection model that is able to predict abnormality in mammography can be useful for patients and medical field in general. The main objective of this project is to investigate different machine learning and deep learning strategies to create the desired model.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# ABBREVIATIONS

**MIAS** Mammographic Imaging Analysis Society

**NA** Not Available

**DDSM** Digital Database for Screening Mammography

**CBIS-DDSM** Curated Breast Imaging Subset of DDSM

**CADx** Computer-Aided Diagnosis

**CADe** Computer-Aided Detection

**ROI** Region-Of-Interest

**DICOM** Digital Imaging and Communications in Medicine

**BI-RADS** Breast Imaging Reporting and Data System

**CC** Bilateral craniocaudal (mammogram)

**MLO** Mediolateral oblique (mammogram)

**CALC** Calcification

**CIRC** Well-defined/circumscribed masses

**SPIC** Spiculated masses

**MISC** Other, ill-defined masses

**ARCH** Architectural distortion

**ASYM** Asymmetry

**NORM** Normal

**CLAHE** Contrast Limited Adaptive Histogram Equalization

**PII** Personally Identifiable Information

**GDPR** General Data Protection Regulation

**HIPAA** Health Insurance Portability and Accountability Act

**YOLO** You Only Look Once

**mAP** Mean Average Precision

# LIST OF FIGURES

IX

X

XI

# LIST OF TABLES

# CHAPTER 1: INTRODUCTION

## 1.1  Overview

Breast cancer is considered a big concern in the whole world, and detecting it early is really important for helping patients and increasing survival rates. In this project machine learning and deep learning approaches are utilized to create a model that can predict malignancy from mammographic images. This chapter explains what the project aims to do and why it's important for making healthcare better.

## 1.2 Project Motivation

Breast cancer is a pressing issue for women in Jordan, according to insights from the King Hussein Cancer Center (KHCC). KHCC data reveals that breast cancer ranks among the most common cancers affecting Jordanian women, with a significant number of new cases emerging annually[4].

Detecting breast cancer early is crucial for improving patient outcomes and lowering mortality rates. Screening tools like mammograms are vital for spotting breast cancer in its early stages, enhancing the chances of successful treatment and survival. However, ensuring the accuracy of mammogram interpretations poses challenges, as factors such as fatigue and human error can lead to misdiagnosis.

In Jordan, studies show that about 30% of breast cancers might not be spotted during the first mammogram screenings[5]. Thus, the importance of reducing mistakes in diagnosis is required.

1

Misinterpreting mammograms can result in unnecessary procedures like biopsies and treatments for benign conditions, as well as delays in identifying and treating actual cases of breast cancer.

To tackle these challenges, there's a proposal to use Computer-Assisted Detection (CAD) software to enhance the accuracy of mammography screening in Jordan. CAD software employs advanced algorithms to reduce incorrect interpretations and enhance the reliability of breast cancer detection.

The aim of implementing CAD software in Jordan is to enhance early breast cancer detection and lower the chances of both false positives and false negatives in mammogram readings. Additionally, there's a broader goal of integrating CAD technology with other deep learning algorithms, with the ultimate aim of creating an artificial intelligence system capable of detecting breast cancer more accurately than traditional methods.

## 1.3 Problem Statement

The current methods of breast cancer detection heavily rely on manual interpretation of medical images, which can be time-consuming, subjective, and leading to variability in diagnoses. Moreover, access to specialized healthcare facilities for diagnosis may be limited in certain regions. This project seeks to address these challenges by developing an automated system capable of accurately detecting breast cancer from medical images, thereby improving diagnostic effectiveness, efficiency and accessibility.

CAD systems utilizing deep learning techniques hold significant promise in enhancing the accuracy of mammogram screenings for detecting early signs of breast cancer. However, these advanced techniques require vast amounts of data to learn the

underlying patterns of cancer and adapt to new cases. Additionally, they demand powerful computing resources to perform the learning process, posing challenges in optimization. Despite these complexities, the potential benefits of utilizing deep learning to predict breast cancer emphasize the importance of exploring and overcoming these technological challenges.

**1.4 Project Aim and Objectives**

The primary goal of this project is to develop a robust predictive model for early detection of breast cancer from medical images. This project will utilize advanced machine learning algorithms and deep learning techniques to analyze medical images and identify patterns indicative of breast cancer. By training the model on a comprehensive dataset of annotated images, it aims to optimize accuracy and reliability. Through rigorous validation and testing, the project seeks to develop a scalable and clinically relevant solution for breast cancer detection. To achieve this aim, the project will focus on the following specific objectives:

- Developing an effective breast cancer detection model from medical images
- Utilize a comprehensive dataset for medical images and preprocess it in an effective and efficient way.
- Design a user-friendly interface that allowing users to access and utilize the detection system effectively

3

**1.5 Project Limitations**

In this project, it's essential to recognize and address the limitations and constraints inherent in the realm of data science and AI. While the goal is to develop a system capable of accurately predicting breast cancer from mammogram images uploaded by users, several challenges must be acknowledged. These challenges include:

- the availability and quality of training data (mammogram images).
- the complexity of interpreting medical images.
- the computational resources required for training and inference.
- ethical considerations surrounding patient privacy and data security must be carefully navigated.

Despite these constraints, the project aims to maximize its potential within defined boundaries, leveraging innovative approaches and methodologies to achieve reliable and meaningful predictions regarding breast cancer diagnosis.

**1.6 Project Expected Output**

In this section, we outline the expected outputs and deliverables of the project, highlighting its potential to yield valuable insights, innovative solutions, or impactful contributions within the field of data science and AI. These outcomes include:

1. development of a robust model capable of accurately predicting breast cancer from mammogram images. This includes:

   1.2 *A generalized breast cancer detection system that is able to detect abnormality in mammogram images.*

4

2.2 *A specialized breast cancer detection system that is able to classify a mammogram image into several specific categories such as calcification and tumor.*

2. Insights into the application of advanced data science techniques in healthcare, potentially paving the way for improved diagnostic accuracy and patient outcomes.

3. A comprehensive preprocessed dataset that will be available for machine learning researchers to conduct more researches in the field of cancer detection.

## 1.7 Project Schedule

This Milestone graph shows the key milestones of our project



Figure 1.1 Project Schedule

## 1.8 Report Organization

**Introduction** provides a comprehensive overview of the background of the subject matter, elucidating the problem at hand and delving into the motivation driving this project. It is followed by a delineation of the objectives that the project endeavors to accomplish.

**Literature Review** a crucial foundation for our research on breast cancer prediction, offering a comprehensive exploration and evaluation of existing literature pertinent to the subject or chosen topic area. Through summarizing prior research and highlighting the connections to our own project, it demonstrates a nuanced

understanding of the field while paving the way for new ideas and insights within breast cancer prediction. By integrating and synthesizing what is known about the subject, the literature review establishes a robust framework for further exploration and contributes to the advancement of knowledge within breast cancer prediction.

**Requirement Engineering and Analysis** This chapter focuses on preparing the groundwork for our project. We start by understanding and preparing the data, exploring various collection methods, preprocessing techniques, and integration strategies. Throughout the process, we emphasize ethical considerations and data privacy to ensure responsible project development. This chapter sets the stage for subsequent phases by establishing a solid foundation for our project's success.

**Model Development and Architectural Design** Examines high-level design considerations related to implementing deep learning and the associated software.

**Model/System Evaluation and Testing Plan** Evaluates the various outcomes to determine the efficiency of different techniques used in training the model, comparing it to other models.

**Model Deployment and Integration** Thoroughly details the steps taken in implementing and integrating deep learning models, and the deployment strategies for them.

**Conclusion and future work** summary of the achievements, results evaluation, lessons learned and potential future directions for the project.

# CHAPTER 2:LITERATURE REVIEW

## 2.1 Existing Systems

In the domain of computer-aided diagnosis (CAD) and detection (CADe) systems for breast cancer mammography, several solutions have been developed to assist radiologists in the interpretation of mammogram images. Notable examples are presented in this section.

Commencing with Karssemeijer and te Brake, who proposed a method based on analyzing the local texture patterns in mammographic images to identify these distortions. The algorithm models used in their system likely involved techniques from computer vision and pattern recognition, possibly including methods like texture analysis, edge detection, and possibly early forms of machine learning algorithms such as support vector machines or neural networks [6].

Another study is conducted by Yang et al., where they evaluated the sensitivity of a computer-aided detection (CAD) system applied to full-field digital mammograms for detecting cancers identified through screening mammography. The CAD system likely employed a combination of image processing techniques, such as noise reduction and contrast enhancement, alongside machine learning models for classification. Commonly utilized machine learning algorithms in CAD systems include Support Vector Machines (SVM), Artificial Neural Networks (ANN), Random Forests (RF), and Gradient Boosting Machines (GBM). These models would have been trained to identify suspicious regions indicative of potential cancerous lesions, utilizing features extracted through segmentation and pattern recognition. The primary objective of the

7

study was to assess the performance of the CAD system in enhancing cancer detection sensitivity within the context of screening mammography **[7]**.

Rangayyan et al investigated boundary modeling and shape analysis methods for the classification of mammographic masses. The research likely incorporated statistical classifiers such as Bayesian classifiers or k-Nearest Neighbors (k-NN) algorithms for classifying masses. Shape analysis techniques, including Fourier descriptors or geometric moments, may have been employed to quantify the shape characteristics of mammographic masses. Additionally, boundary modeling algorithms such as active contours (snakes) or level set methods were probably utilized to detect and represent the boundaries of masses. These models were likely integrated to develop a classification system capable of distinguishing between different types of mammographic masses based on their shape features **[8]**.

Mudigonda et al. investigated the utility of gradient and texture analysis for classifying mammographic masses. The research likely employed statistical classifiers such as Bayesian classifiers or k-Nearest Neighbors (k-NN) algorithms to categorize masses based on their gradient and texture characteristics. Additionally, techniques like Gabor filters for texture analysis and edge detection algorithms for gradient analysis may have been utilized to extract relevant features from mammographic images. These features were likely integrated into a classification framework to develop a system capable of accurately distinguishing between different types of mammographic masses based on their gradient and texture properties **[9]**.

Görgel et al. investigated the computer-aided classification of breast masses in mammogram images using the spherical wavelet transform and Support Vector Machines (SVM). The research likely utilized the spherical wavelet transform for

8

feature extraction from mammogram images, which captures both spatial and frequency information efficiently. These features were then fed into a SVM classifier, a supervised learning model known for its effectiveness in binary classification tasks. The SVM would have been trained on the extracted features to differentiate between benign and malignant breast masses. This approach aimed to develop a robust computer-aided classification system capable of accurately categorizing breast masses based on their characteristics in mammogram images **[10]**.

The following tables shows performance statistics of selected CADe CADx methods. Table 1 presents the performance statistics of selected CADe methods for the detection of abnormalities. While Table 2 shows the performance statistics of selected CADx methods for the classification of masses.

| Authors | Size of Data set (Cases) | Public or private data | Accuracy |
|---|---|---|---|
| **Karssemeijer and te Brake** | **50** | **Public MIAS** | **NA** |
| **Yang et al** | **203** | **Private** | **96.1%** |

Table1: Performance statistics of selected CADe methods for the detection of

abnormalities **[11]**.

| Authors | Size of Data set (Cases) | Public or private data | Classification Accuracy |
|---|---|---|---|
| Rangayyan et al. | 54 | Public MIAS and Private | 91% |
| Mudigonda et al. | 56 | Public MIAS | 82.1% |
| Görgel et al | 78, 65 | Private, Public MIAS | 91.4%, 90.1% |

Table 2: Performance statistics of selected CADx methods for the classification

of masses [11].

## 2.2 Limitations of Existing Systems

Despite the notable progress achieved by these systems, there are several limitations.

*With respect to Karssemeijer and Brake*, their work is limited by the computational resources available, potentially resulting in slower processing speeds and less sophisticated algorithms, which could impact the system's accuracy and efficiency [6].

*Regarding the work conducted by Yang et al.* it was vulnerable to variations in the quality of mammographic images, which may affect the CAD system's performance and reliability, potentially leading to missed detections or false positives [7].

*Rangayyan et al*. study relied heavily on the accuracy of the segmentation process, which can be challenging and prone to errors, especially in cases where masses exhibit irregular shapes or indistinct boundaries, potentially leading to misclassification [8].

*With respect to Mudigonda et al.* work, the effectiveness of gradient and texture analysis methods may be limited in capturing subtle variations in mammographic

images, especially in cases where masses exhibit complex or heterogeneous textures, potentially leading to misclassification or reduced sensitivity **[9].**

***The work performed by Görgel et al.*** is susceptible to variations in image acquisition settings, including differences in imaging modalities, parameters, and techniques, which may introduce inconsistencies in feature extraction and classification, potentially affecting the CAD system's performance and generalizability **[10].**

## 2.3 Overall Solution Approach

To address the limitations of existing CADx and CADe systems for breast cancer detection, our project proposes a novel approach that leverages state-of-the-art deep learning techniques and large-scale mammography datasets. By training a convolutional neural network (CNN) on diverse and annotated mammogram images, our solution aims to improve both sensitivity and specificity in detecting breast cancer lesions.Additionally, our solution will prioritize user-centric design principles to ensure seamless integration into clinical workflows and enhance usability for healthcare professionals. Overall, our project serves as a starting point for the development of an advanced CADe system that addresses the current limitations and pushes the boundaries of breast cancer detection in mammography.

# CHAPTER 3: REQUIREMENT ENGINEERING AND ANALYSIS

## 3.1 Data Understanding and Preparation

### 3.1.1    Data Collection Methods and Sources

In this project, we utilized two datasets to ensure comprehensive analysis and robust results. The first dataset is the CBIS-DDSM (Curated Breast Imaging Subset of DDSM) is an updated and standardized version of the Digital Database for Screening Mammography (DDSM). The DDSM is a database of 2,620 scanned film mammography studies, containing normal, benign, and malignant cases with verified pathology information. The scale of the database, along with ground truth validation, makes the DDSM a useful tool for developing and testing decision support systems. The CBIS-DDSM collection includes a subset of the DDSM data, selected and curated by a trained mammographer. The images have been decompressed and converted to DICOM format. Updated ROI segmentation, bounding boxes, and pathologic diagnoses for training data are also included **[11]**.

Note here that the Digital Database for Screening Mammography (DDSM) is a resource for use by the mammographic image analysis research community. The primary purpose of the database is to facilitate research in the development of computer algorithms to aid in screening. Secondary purposes of the database may include the development of algorithms to aid in diagnosis and the creation of teaching or training aids. The database contains approximately 2,620 studies. Each study includes two images of each breast, along with some associated patient information (age at time of study, ACR breast density rating, subtlety rating for abnormalities, ACR keyword description of abnormalities) and image information (scanner, spatial resolution, etc.). Images containing suspicious areas have associated pixel-level "ground truth" information about the locations and types of suspicious regions. Also provided is software for accessing the mammogram and truth images and for calculating performance figures for automated image analysis algorithms **[11]**.

The dataset can be split as follows:

- By Case:
  - Mass: 1,318 images
  - Calcification: 1,545 images



Figure 3.1 Example of Mammogram image mass [11]



Figure 3.2 Example of Mammogram image calcification[11]

- By View:
  - MLO: 1,896 images
  - CC: 1,672 images



Figure 3.3 Example of Mammogram image MLO [11]



Figure 3.4 Example of Mammogram image CC [11]

- By Breast:
  - Left: 1,819 images
  - Right: 1,749 images



Figure 3.5 Example of

Mammogram image right [11]



Figure 3.6 Example of

Mammogram image left [11]

- By Pathology:
  - Malignant: 1,457 images
  - Benign: 1,429 images
  - Benign without callback: 682 images



Figure 3.7 Example of
Mammogram image
malignant [11]

Figure 3.8 Example of
Mammogram image
benign[11]

Figure 3.9 Example of
Mammogram image
benign without call
back[11]

*Note here that all figures were extracted from CBIS-DDSM dataset* **[11]***.*

- By BI-RADS Assessment:
    - Assessment 5: 573 images
    - Assessment 4: 1,633 images
    - Assessment 3: 477 images
    - Assessment 2: 644 images
    - Assessment 1: 3 images
    - Assessment 0: 238 images

- Additional Tags:
    - Calcification type
    - Calcification distrialcibution
    - Mass shape
    - Mass margins
    - Breast density
    - Subtlety
    - Patient ID

**Dataset Features and Descriptions**

The dataset is originally divided into two types:

I. Calcification (train and test), here are the descriptions for each feature:

**1. Patient ID**

- **Description:** A unique identifier assigned to each patient in the dataset, ensuring privacy and enabling specific data retrieval for individual patients.

**2. Breast Density**

- **Description:** Indicates the density of breast tissue, which is classified into four categories:
    - **1:** Least dense (mostly fatty tissue).
    - **2:** Scattered areas of fibro glandular density.

15

- o **3:** Heterogeneously dense breasts with many areas of fibro glandular density, which may obscure underlying abnormalities on a mammogram.

- o **4:** Most dense (mostly glandular and fibrous tissue).

## 3. Breast Side

- **Description:** Identifies whether the image is of the left or right breast.

  - o **Values:** Left, Right

## 4. Image View

- **Description:** Type of mammographic view, crucial for detecting abnormalities and assessing breast density accurately:

  - o **CC (Cranio-caudal view):** Image taken with the x-ray beam passing from the top of the breast (cranial) to the bottom (caudal), providing a straight-on view.

  - o **MLO (Medio-lateral oblique view):** Image taken with the x-ray beam passing from the side (medial) to the side of the breast (lateral) at an oblique angle, providing a side-angle view.

## 5. Abnormality ID

- **Description:** Identifies the Breast Imaging-Reporting and Data System (BI-RADS) categories, which standardize the interpretation of mammograms and other breast imaging studies:

  BI-RADS classifications range from 0 to 6:

  - o BI-RADS 0: Incomplete - Additional imaging evaluation needed.
  - o BI-RADS 1: Negative - No significant abnormality.
  - o BI-RADS 2: Benign - Non-cancerous findings.
  - o BI-RADS 3: Probably Benign - Findings have a very low likelihood of being cancer (less than 2%), short-term follow-up suggested.

- BI-RADS 4: Suspicious Abnormality - Findings not characteristic of breast cancer but with a reasonable probability of being malignant (subdivided into 4A, 4B, and 4C indicating increasing levels of suspicion).
- BI-RADS 5: Highly Suggestive of Malignancy - Findings have a high probability of being cancerous (greater than 95%).
- BI-RADS 6: Known Biopsy-Proven Malignancy - Findings where cancer has already been confirmed by biopsy.

- **Importance:** These categories help standardize reporting, making it easier for healthcare providers to interpret and compare results across different imaging studies and over time.

## 6. Abnormality Type

- **Description:** Type of abnormality detected.
  - **Values:** Calcification

## 7. Calcification Type

- **Description:** Represents different morphological types or characteristics of calcifications observed in breast imaging studies:
  - **Amorphous:** Calcifications appear shapeless or without a distinct form.
  - **Pleomorphic:** Calcifications have varying shapes and sizes within the same cluster.
  - **Round and Regular-Lucent Center-Dystrophic:** Calcifications have a round and regular shape with a lucent center, indicating degenerating or necrotic tissue.
  - **Punctate:** Calcifications appear as tiny dots or specks.
  - **Coarse:** Larger and have a coarse or rough texture.
  - **Vascular:** Calcifications have a vascular pattern or appearance.

17

- **Fine Linear Branching:** Calcifications appear as fine linear structures that branch out.

- **Large Rodlike:** Calcifications are large and rod-shaped.

- **Punctate-Lucent Center:** Calcifications have a punctate appearance with a central lucent area.

- **Vascular-Coarse-Lucent Center-Round and Regular-Punctate:** Calcifications exhibit a combination of vascular, coarse, lucent center, round and regular, and punctate features.

- **Round and Regular-Eggshell:** Calcifications have a round and regular shape with an eggshell appearance.

- **Punctate-Pleomorphic:** Calcifications exhibit a combination of punctate and pleomorphic features.

- **Pleomorphic-Fine Linear Branching:** Calcifications exhibit a combination of pleomorphic and fine linear branching features.

- **Dystrophic:** Calcifications occur in degenerating or necrotic tissue.

- **Lucent Center:** Calcifications have a central area of lucency.

- **Amorphous-Pleomorphic:** Calcifications exhibit a combination of amorphous and pleomorphic features.

- **Round and Regular:** Calcifications have a round and regular shape.

- **Vascular-Coarse-Lucent Centered:** Calcifications exhibit a combination of vascular, coarse, and centered lucent center features.

- **Coarse-Round and Regular:** Calcifications exhibit a combination of coarse and round and regular features.

- **Coarse-Pleomorphic:** Calcifications exhibit a combination of coarse and pleomorphic features.

- **Lucent Centered:** Calcifications have a centered area of lucency.

- **Vascular-Coarse:** Calcifications exhibit a combination of vascular and coarse features.

- **Round and Regular-Punctate:** Calcifications exhibit a combination of round and regular and punctate features.

- **Round and Regular-Lucent Center:** Calcifications exhibit a combination of round and regular and lucent center features.

- **Coarse-Round and Regular-Lucent Centered:** Calcifications exhibit a combination of coarse, round and regular, and centered lucent center features.

- **Skin:** Calcifications have an appearance similar to skin.

- **Lucent Center-Punctate:** Calcifications have a central area of lucency with punctate features.

- **Skin-Punctate:** Calcifications have an appearance similar to skin with punctate features.

- **Skin-Punctate-Round and Regular:** Calcifications have an appearance similar to skin with punctate and round and regular features.

- **Milk of Calcium:** Calcifications appear as clusters of calcified milk-like material.

- **Pleomorphic-Pleomorphic:** Calcifications exhibit a combination of pleomorphic features.

- **Skin-Coarse-Round and Regular:** Calcifications have an appearance similar to skin with coarse and round and regular features.

- **Amorphous-Round and Regular:** Calcifications exhibit a combination of amorphous and round and regular features.

- **Round and Regular-Pleomorphic:** Calcifications exhibit a combination of round and regular and pleomorphic features.

- **Round and Regular-Punctate-Amorphous:** Calcifications exhibit a combination of round and regular, punctate, and amorphous features.

- **Round and Regular-Amorphous:** Calcifications exhibit a combination of round and regular and amorphous features.

- o **Coarse-Round and Regular-Lucent Center:** Calcifications exhibit a combination of coarse, round and regular, and lucent center features.

- o **Large Rodlike-Round and Regular:** Calcifications exhibit a combination of large rodlike and round and regular features.

- o **Round and Regular-Lucent Center-Punctate:** Calcifications exhibit a combination of round and regular, lucent center, and punctate features.

- o **Coarse-Lucent Center:** Calcifications exhibit a combination of coarse and lucent center features.

- o **Punctate-Amorphous:** Calcifications exhibit a combination of punctate and amorphous features.

- o **Round and Regular-Lucent Centered:** Calcifications exhibit a combination of round and regular and centered lucent center features.

- o **Punctate-Round and Regular:** Calcifications exhibit a combination of punctate and round and regular features.

- o **Eggshell:** Calcifications have an eggshell-like appearance.

- o **Punctate-Fine Linear Branching:** Calcifications exhibit a combination of punctate and fine linear branching features.

- o **Importance:** These descriptions provide an understanding of the various morphological types of calcifications observed in breast imaging studies.

## 8. Calcification Distribution

- **Description:** Represents the spatial distribution patterns of calcifications observed in breast imaging studies:

  - o **Clustered:** Calcifications are grouped closely together in a localized area, often forming a cluster or cluster-like pattern.

  - o **Linear:** Calcifications are arranged in a linear or elongated pattern, typically following a straight or curved line.

20

- o **Regional:** Calcifications are distributed over a larger region of the breast rather than being confined to a specific area.

- o **Diffusely Scattered:** Calcifications are dispersed throughout the breast tissue in a diffuse or scattered manner, without a distinct pattern.

- o **Segmental:** Calcifications are distributed within a specific segment or quadrant of the breast, indicating a localized but larger area compared to clustered distribution.

- o **Clustered-Linear:** Calcifications exhibit a combination of clustered and linear distribution patterns.

- o **Clustered-Segmental:** Calcifications exhibit a combination of clustered and segmental distribution patterns.

- o **Linear-Segmental:** Calcifications exhibit a combination of linear and segmental distribution patterns.

- o **Regional-Regional:** Calcifications exhibit a combination of regional distribution patterns, suggesting involvement of multiple regions or quadrants of the breast.

- o **Importance:** These descriptions provide an understanding of the various distribution patterns of calcifications observed in breast imaging studies.

## 9. Assessment

- **Description:** In the context of breast imaging, the values in the "assessment" column typically correspond to the BI-RADS (Breast Imaging-Reporting and Data System) assessment categories, which help classify the level of suspicion for breast abnormalities:

  - o **Incomplete assessment:** Indicates that the assessment is incomplete, and further evaluation or additional imaging is needed to make a final assessment.

  - o **Negative:** No abnormality is detected, and the mammogram or imaging study is considered negative for any significant findings.

21

- **Benign:** The abnormality detected is benign, indicating that it is not cancerous or harmful.

- **Probably benign:** Suggests that the abnormality is likely benign, but further evaluation or short-term follow-up is recommended to confirm its benign nature.

- **Suspicious:** Indicates that the abnormality is suspicious for malignancy, meaning that it may be cancerous, and further evaluation or biopsy is warranted.

- **Highly suggestive of malignancy:** The abnormality is highly suspicious for malignancy, and immediate action such as biopsy or treatment may be necessary.

- **Importance:** These categories help radiologists communicate findings and guide patient management decisions based on the level of suspicion for breast abnormalities.

## 10. Pathology

- **Description:** The "pathology" column represents the pathological diagnosis of breast abnormalities identified through imaging studies such as mammograms:

  - **Malignant:** Indicates that the breast abnormality is cancerous or indicative of breast cancer. Further evaluation, such as biopsy and treatment, is usually necessary.

  - **Benign:** The breast abnormality is non-cancerous and not harmful. It does not pose a threat to health, and no further treatment may be required, although periodic monitoring may be recommended.

  - **Benign without callback:** Suggests that the abnormality is likely benign based on imaging findings, but further evaluation or follow-up imaging is recommended to confirm its benign nature and ensure no changes over time. It's a cautious approach taken when the abnormality

doesn't definitively appear benign but also doesn't clearly indicate malignancy.

- **Importance:** These categories help in classifying breast abnormalities based on their pathological characteristics and guide further management and treatment decisions accordingly.

## 11. Subtlety

- **Description:** In breast imaging, the "subtlety" column typically represents the level of subtlety or difficulty in identifying abnormalities on the imaging studies, such as mammograms:

  - **Not applicable or undefined:** This value may indicate that the subtlety level is not applicable to the specific case or that it's undefined.

  - **Minimal:** Abnormalities are easily identifiable and have minimal subtlety.

  - **Low:** Abnormalities are somewhat subtle but still detectable with moderate effort.

  - **Moderate:** Abnormalities have a moderate level of subtlety, requiring careful examination to identify.

  - **High:** Abnormalities are quite subtle and challenging to detect, requiring significant expertise and attention to detail.

  - **Extreme:** Abnormalities are extremely subtle and difficult to detect, often requiring specialized techniques or advanced imaging modalities for identification.

- **Importance:** These values help radiologists and clinicians understand the difficulty level associated with identifying abnormalities on imaging studies and may influence decisions regarding further evaluation and management.

II. Mass  (train and test), here are the descriptions for each feature:

## 1. Patient ID

- **Description:** A unique identifier assigned to each patient in the dataset.

23

- **Importance:** Ensures the privacy and confidentiality of patient information while allowing for the organization and tracking of patient records.

## 2. Breast Density

- **Description:** Categorizes the density of the breast tissue seen on a mammogram, which can affect the accuracy of the mammogram:

  - **1:** Least dense (mostly fatty tissue).

  - **2:** Scattered areas of fibro glandular density.

  - **3:** Heterogeneously dense breasts, meaning there are many areas of fibro glandular density that may obscure underlying abnormalities.

  - **4:** Most dense (mostly glandular and fibrous tissue).

- **Importance:** Dense breast tissue can make it more difficult to detect abnormalities and is associated with an increased risk of breast cancer.

## 3. Left or Right Breast

- **Description:** Indicates which breast (left or right) the image pertains to.

- **Importance:** Provides clarity on the location of the abnormality for accurate diagnosis and treatment planning.

## 4. Image View

- **Description:** Specifies the type of mammographic view:

  - **CC (Cranio-caudal):** Image taken from the top to the bottom of the breast.

  - **MLO (Medio-lateral oblique):** Image taken from the side at an oblique angle.

- **Importance:** Different views help radiologists detect abnormalities and assess breast density more accurately by providing multiple perspectives of the breast tissue.

## 5. Abnormality ID

- **Description:** Denotes the BI-RADS (Breast Imaging-Reporting and Data System) categories used to standardize the interpretation of mammograms:

BI-RADS classifications range from 0 to 6:

- o BI-RADS 0: Incomplete - Additional imaging evaluation needed.
- o BI-RADS 1: Negative - No significant abnormality.
- o BI-RADS 2: Benign - Non-cancerous findings.
- o BI-RADS 3: Probably Benign - Findings have a very low likelihood of being cancer (less than 2%), short-term follow-up suggested.
- o BI-RADS 4: Suspicious Abnormality - Findings not characteristic of breast cancer but with a reasonable probability of being malignant (subdivided into 4A, 4B, and 4C indicating increasing levels of suspicion).
- o BI-RADS 5: Highly Suggestive of Malignancy - Findings have a high probability of being cancerous (greater than 95%).
- o BI-RADS 6: Known Biopsy-Proven Malignancy - Findings where cancer has already been confirmed by biopsy.

- These categories help standardize reporting, making it easier for healthcare providers to interpret and compare results across different imaging studies and over time.

## 6. Abnormality Type

- **Description:** Specifies the type of abnormality detected, in this case, a "Mass."

- **Importance:** Identifying the type of abnormality is crucial for diagnosis, prognosis, and treatment planning.

## 7. Mass Shape

- **Description:** Describes the shape of the detected mass:

  - o **IRREGULAR-ARCHITECTURAL_DISTORTION:**Irregular shape with architectural distortion.

  - o **ARCHITECTURAL_DISTORTION:** Disruption in normal tissue arrangement.

  - o **OVAL:** Elongated with a smooth outline.

25

- o **IRREGULAR:** Lacking a defined or symmetrical form.

- o **LYMPH_NODE:** Resembling a lymph node.

- o **LOBULATED-LYMPH_NODE:** Lobulated with lymph node characteristics.

- o **LOBULATED:** Multiple rounded contours.

- o **FOCAL_ASYMMETRIC_DENSITY:** Localized density difference.

- o **ROUND:** Circular outline.

- o **LOBULATED-ARCHITECTURAL_DISTORTION:** Lobulated with architectural distortion.

- o **ASYMMETRIC_BREAST_TISSUE:** Non-symmetric distribution without a distinct shape.

- o **LOBULATED-IRREGULAR:** Combination of lobulated and irregular shapes.

- o **OVAL-LYMPH_NODE:** Oval shape resembling a lymph node.

- o **LOBULATED-OVAL:** Combination of lobulated and oval shapes.

- o **ROUND-OVAL:** Combination of round and oval shapes.

- o **IRREGULAR-FOCAL_ASYMMETRIC_DENSITY:**Irregular shape with focal asymmetric density.

- o **ROUND-IRREGULAR-ARCHITECTURAL_DISTORTION:** Combination of round, irregular, and architectural distortion.

- o **ROUND-LOBULATED:** Combination of round and lobulated shapes.

- **Importance:** The shape of the mass can provide significant clues regarding its nature (benign or malignant) and guide further diagnostic procedures.

### 8. Mass Margins

- **Description:** Describes the borders of the mass:

  - o **SPICULATED:** Sharp, spike-like projections.

  - o **ILL_DEFINED:** Lack of clear demarcation.

- **CIRCUMSCRIBED:** Well-defined boundary.

- **ILL_DEFINED-SPICULATED:** Combination of unclear margins and sharp projections.

- **OBSCURED:** Not clearly visible or distinguishable.

- **OBSCURED-ILL_DEFINED:** Combination of unclear visibility and ill-defined margins.

- **MICROLOBULATED:** Tiny lobulations or irregularities.

- **MICROLOBULATED-ILL_DEFINED-SPICULATED:** Combination of tiny irregularities, unclear margins, and sharp projections.

- **MICROLOBULATED-SPICULATED:** Combination of tiny irregularities and sharp projections.

- **CIRCUMSCRIBED-ILL_DEFINED:** Combination of well-defined and unclear margins.

- **MICROLOBULATED-ILL_DEFINED:** Tiny irregularities and unclear margins.

- **CIRCUMSCRIBED-OBSCURED:** Combination of well-defined boundaries and unclear visibility.

- **OBSCURED-SPICULATED:** Unclear boundaries with sharp spikes.

- **OBSCURED-ILL_DEFINED-SPICULATED:** Unclear boundaries with sharp projections.

- **CIRCUMSCRIBED-MICROLOBULATED:** Well-defined boundaries with tiny irregularities.

- **CIRCUMSCRIBED-OBSCURED-ILL_DEFINED:** Mix of well-defined boundaries, unclear visibility, and lack of demarcation.

- **CIRCUMSCRIBED-SPICULATED:** Well-defined boundaries and sharp projections.

- **Importance:** The margins of a mass are critical in determining the likelihood of malignancy. Spiculated margins are more commonly associated with malignant masses.

### 9. Assessment

- **Description:** Corresponds to the BI-RADS assessment categories, which classify the level of suspicion for breast abnormalities:

  - **0:** Incomplete assessment; further evaluation needed.

  - **1:** Negative; no significant findings.

  - **2:** Benign; not cancerous or harmful.

  - **3:** Probably benign; further evaluation recommended.

  - **4:** Suspicious; further evaluation or biopsy warranted.

  - **5:** Highly suggestive of malignancy; immediate action required.

- **Importance:** These categories help radiologists communicate findings and guide patient management decisions based on the level of suspicion for breast abnormalities.

### 10. Pathology

- **Description:** Represents the pathological diagnosis of breast abnormalities:

  - **1:** Malignant; cancerous.

  - **2:** Benign; non-cancerous.

  - **3:** Benign without callback; likely benign but further evaluation recommended.

- **Importance:** Classifies breast abnormalities based on their pathological characteristics, guiding further management and treatment decisions.

### 11. Subtlety

- **Description:** Represents the level of subtlety or difficulty in identifying abnormalities on imaging studies:

- o **0:** Not applicable or undefined.

- o **1:** Minimal; easily identifiable.

- o **2:** Low; somewhat subtle but detectable.

- o **3:** Moderate; requires careful examination.

- o **4:** High; challenging to detect.

- o **5:** Extreme; extremely difficult to detect.

- **Importance:** Helps radiologists and clinicians understand the difficulty level associated with identifying abnormalities on imaging studies, influencing further evaluation and management decisions.

The second dataset used in this project is the MIAS (Mammographic Image Analysis Society) database, which consists of 322 digitized mammogram images. These images include both normal and abnormal cases, and the dataset is widely used for research in the development of automated breast cancer detection algorithms. The MIAS dataset is designed to aid in evaluating the performance of image analysis techniques by providing ground truth annotations, such as lesion type and its location within the breast tissue. The images in the MIAS database come with a variety of characteristics, including different levels of image quality and resolution. This diversity makes the dataset a valuable resource for developing robust and generalized algorithms for mammography analysis and breast cancer detection**[12]**.

By popular request, the original MIAS Database (digitised at 50 micron pixel edge) has been reduced to 200 micron pixel edge and clipped/padded so that every image is $1024 \times 1024$ pixels **[12]**.

The dataset can be split as follows:

- **Fisrt column**: MIAS database reference number.

- **Second column:** Character of background tissue:
  - o   F  Fatty
  - o   G  Fatty-glandular
  - o   D  Dense-glandular



Figure 3.10 Example of Mammogram image Fatty [13].

Figure 3.11 Example of Mammogram image Fatty-glandular[13].

Figure 3.12 Example of Mammogram image Dense-glandular[13].

- **Third column**: Class of abnormality present:

  - o   **CALC** Calcification
  - o   **CIRC** Well-defined/circumscribed masses
  - o   **SPIC** Spiculated masses
  - o   **MISC** Other, ill-defined masses
  - o   **ARCH** Architectural distortion
  - o   **ASYM** Asymmetry
  - o   **NORM** Normal

Figure 3.13 Example of Mammogram image CALC [13].

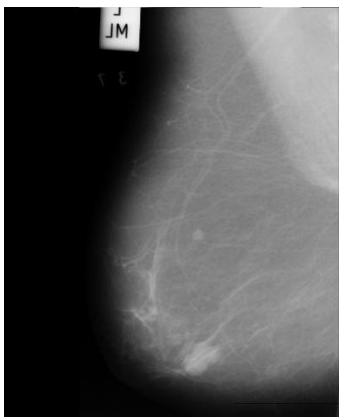

Figure 3.14 Example of Mammogram image CIRC [13].



Figure 3.15 Example of Mammogram image SPIC [13].



Figure 3.16 Example of Mammogram image MISC [13].



Figure 3.17 Example of Mammogram image ARCH [13].



Figure 3.18 Example of Mammogram image ASYM [13].



Figure 3.19 Example of Mammogram image NORM [13].

31

- **Fourth column**: Severity of abnormality:

    o **B** Benign

    o **M** Malignant



Figure 3.20 Example of
Mammogram image Benign **[13]**.



Figure 3.21 Example of Mammogram
Image Malignant **[13]**.

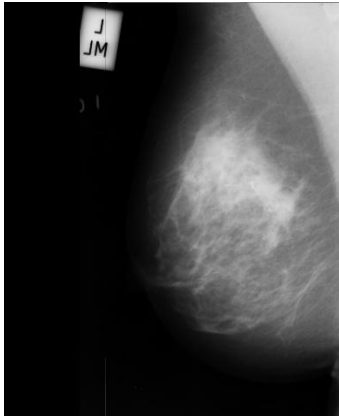- **Fifth & Sixth columns**: x,y image-coordinates of centre of abnormality.

- **Seventh column:** Approximate radius (in pixels) of a circle enclosing the abnormality.

Additional info:

- The list is arranged in pairs of films, where each pair represents the left (even filename numbers) and right mammograms (odd filename numbers) of a single patient **[12]**.

- The size of all the images is 1024 pixels x 1024 pixels. The images have been centered in the matrix**[12]** .

- Coordinate system origin is the bottom-left corner**[12]**.

- In some cases calcifications are widely distributed throughout the image rather than concentrated at a single site. In these cases centre locations and radii are inappropriate and have been omitted**[12]**.

### *3.1.2   Data Preprocessing Techniques and Cleaning*

*In this section we present the conducted preprocessing. The performed preprocessing on the first dataset(CBIS-DDSM), includes: (i) fixing images path, (ii) handling missing values and (iii) image resizing.*

1. **Fixing Images Path**

The dataset was extracted from DICOM files. The mass and calcification files have paths where images end with. dcm extention, and the file containing images is in JPEG format. Therefore, we prepared the paths in the files.

There are individual files for mass and calcification training and test sets including:

- **mass_case_description_train_set.csv**

- **mass_case_description_test_set.csv**

- **calc_case_description_train_set.csv**

- **calc_case_description_test_set.csv**

2. **Handling Missing Values**

We performed a thorough check to identify any null values within the dataset. Upon inspection, it was determined that no null values were present in the dataset.

3. **Image Resizing**

Prior to analysis, all images in the dataset were resized to a uniform dimension of 640 x 640 x 3 to ensure consistency. The images were then saved in tensor format to facilitate efficient processing and analysis.

*The performed preprocessing on the second dataset (MIAS), includes: (i) Data Augmentation, (ii) Gaussian Blurring ,(iii) Median Filtering,(iv) Histogram Equalization,(v) CLAHE Equalization (Contrast Limited Adaptive Histogram Equalization) and (vi) Image Inversion.*

## 1. Data Augmentation

Data augmentation is a technique used to increase the diversity of a dataset by applying various transformations, enhancing model generalization and performance. In our project, we utilized several data augmentation techniques, including **(a)Random Resized Crop,(b) Rotate,(c)Random Brightness Contrast ,(d) Horizontal Flip**, and **(e)Blur**. These methods helped improve the robustness of our model by introducing variations in the training data.

(a) **Random Resized Crop**

**Random Resized Crop** involves randomly selecting a region of the image based on a specified scale and resizing it to a fixed size. The cropping scale determines how much of the original image is retained, adding diversity to the dataset.

- **Purpose:**
  - Simulates varying object sizes and positions within the image.
  - Helps the model become robust to objects appearing at different scales and locations.

- **Benefits:**
  - Encourages the model to focus on different parts of the image.
  - Improves generalization for tasks like object detection and classification.

Figure 3.22(A) Example of original Mammogram image CALC [13]



Figure 3.22(B)Example of Mammogram image CALC after applying Random Resized Crop.



Figure 3.23(A) Example of original Mammogram image CIRC [13]



Figure 3.23(B) Example of Mammogram image CIRC after applying Random Resized Crop.



Figure 3.24(A)Example of original Mammogram image NORM [13]



Figure 3.24(B)Example of Mammogram image NORM after applying Random Resized Crop.

(b) **Rotate**

The Rotate technique introduces a random rotation to the image, constrained within a specified angular range (e.g., ±15 degrees). The background is usually filled with a constant value or interpolated pixels.

- **Purpose:**
    o Addresses scenarios where objects may appear rotated in the real world.
    o Enhances model invariance to orientation changes.
- **Benefits:**
    o Reduces overfitting by augmenting the dataset with diverse object orientations.
    o Especially useful for tasks involving multi-view object recognition.



Figure 3.25(A)Example of

original Mammogram image

CALC **[13]**



Figure 3.25(B)Example of

Mammogram image CALC

after applying rotate.

Figure 3.26(A)Example of
original Mammogram image
CIRC [13]



Figure 3.26(B)Example of
Mammogram image CIRC
after applying rotate.



Figure 3.27(A)Example of
original Mammogram image
NORM [13]



Figure 3.27(B)Example of
Mammogram image NORM
after applying rotate.

## (c) Random Brightness Contrast

This technique randomly adjusts the brightness and contrast levels of an image to simulate real-world variations in lighting conditions.

- **Purpose:**
  - Mimics different lighting environments, such as shadows, dim light, or overexposure.
  - Prepares models to perform well under varying brightness and contrast levels.

- **Benefits:**
  - Makes models robust to lighting changes.
  - Improves performance for outdoor tasks or applications with inconsistent illumination, such as autonomous driving or drone imaging.
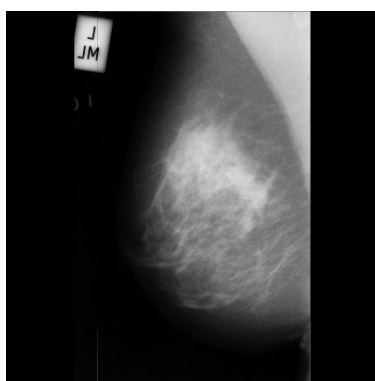


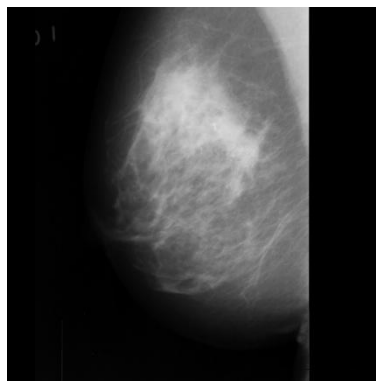Figure 3.28(A)Example of original Mammogram image CALC **[13]**



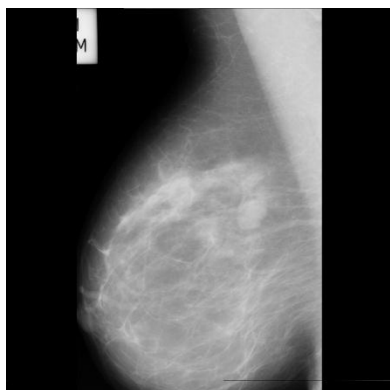Figure 3.28(B)Example of Mammogram image CALC after applying Random Brightness Contrast.



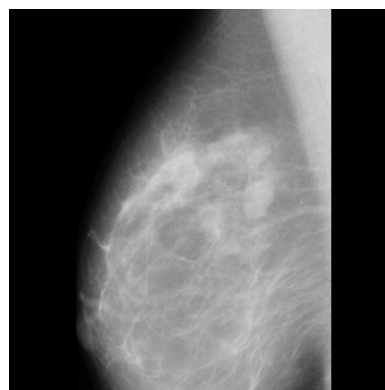Figure 3.29(A)Example of original Mammogram image CIRC **[13]**



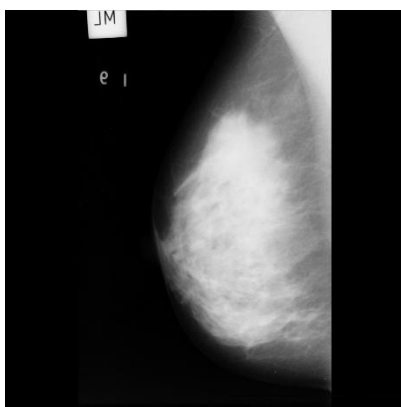Figure 3.29(B)Example of Mammogram image CIRC after applying Random Brightness Contrast.

Figure 3.30 (A)Example of
original Mammogram image
NORM **[13].**

Figure 3.30(B)Example of
Mammogram image NORM
after applying Random
Brightness Contrast.

(d) **Horizontal Flip**

Horizontal Flip mirrors the image along the vertical axis, effectively flipping it horizontally. This transformation is applied randomly based on a specified probability.

- **Purpose:**
  - o Introduces symmetry to the dataset.
  - o Useful for tasks where objects look similar when flipped, such as human faces, animals, or symmetrical objects.
- **Benefits:**
  - o Expands the dataset size by creating flipped versions of images.

  - o Reduces the risk of bias in models that might overfit to the original orientation of objects.

Figure 3.31(A)Example of
original Mammogram image
CALC **[13]**.

.



Figure 3.31(B)Example of
Mammogram image CALC
after applying Horizontal Flip.



Figure 3.32(A)Example of
original Mammogram image
CIRC **[13].**



Figure 3.32(B)Example of
Mammogram image CIRC
after applying Horizontal Flip.



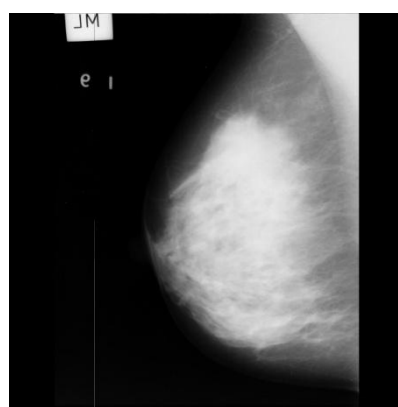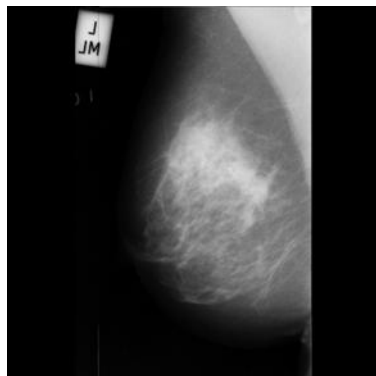Figure 3.33(A)Example of
original Mammogram image
NORM **[13].**



Figure 3.33(B)Example of
Mammogram image NORM
after applying Horizontal Flip.

**(e) Blur**

The Blur technique applies a soft blurring effect to the image using a randomly determined kernel size. This technique is often used to simulate slight defocusing or motion blur.

- **Purpose:**
  - Reduces the clarity of fine details, mimicking real-world imperfections like camera shake or out-of-focus areas.
  - Forces the model to rely on broader patterns rather than fine textures.
- **Benefits:**
  - Helps models generalize better to low-quality or blurry images.
  - Enhances robustness for applications in photography, surveillance, and low-resolution image analysis.



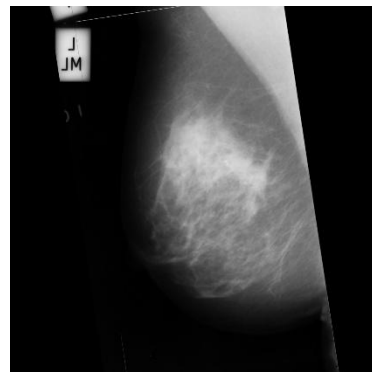Figure 3.34(A)Example of original Mammogram image CALC **[13]**



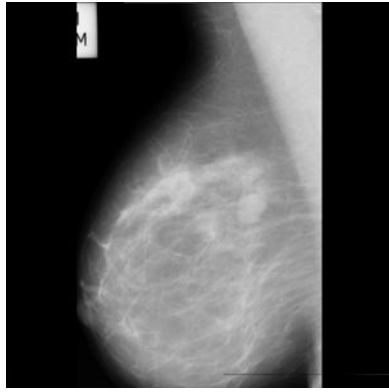Figure 3.34(B)Example of Mammogram image CALC after applying blur.



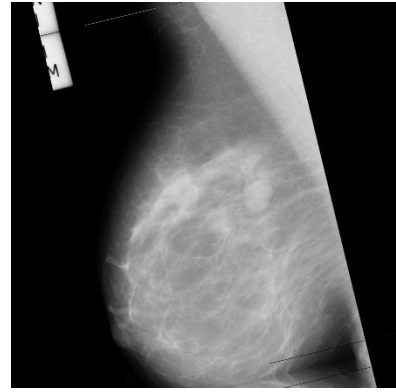Figure 3.35(A)Example of original Mammogram image CIRC **[13]**



Figure 3.35(B)Example of Mammogram image CIRC after applying blur.
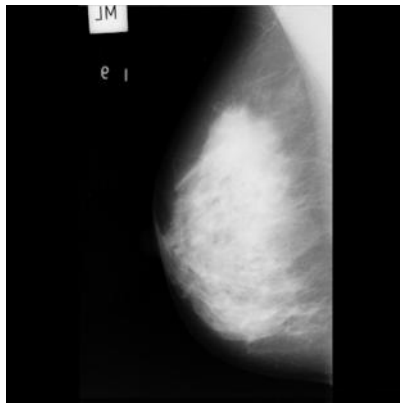
41

Figure 3.36(A)Example of
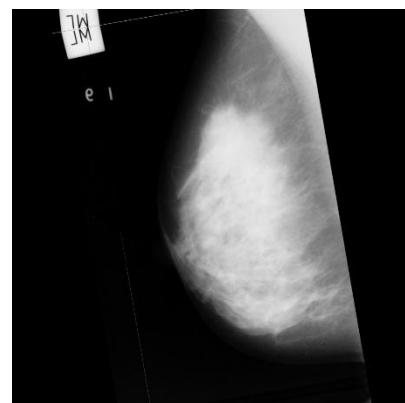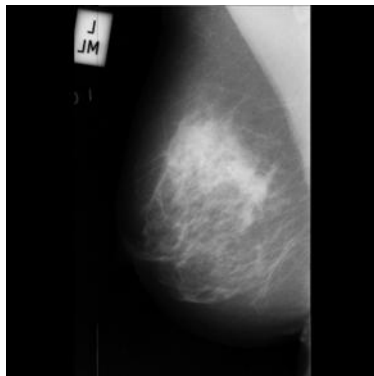original Mammogram image
NORM [13]



Figure 3.36(B)Example of
Mammogram image NORM
after applying blur.

2. **Gaussian Blurring**

   Gaussian Blurring is a widely used technique in image processing to reduce noise and detail in an image. It works by convolving the image with a Gaussian function, which results in a smooth blurring effect. The Gaussian filter assigns higher weights to the pixels closer to the center of the kernel, making it effective for reducing high-frequency noise.

   - **Applications:**
     - o Smoothing images.
     - o Preprocessing for edge detection algorithms.
   - **Advantages:**
     - o Preserves edges better than a simple averaging filter.



Figure 3.37 Example
of original Mammogram
image CALC [13].



Figure 3.38 Example of
original Mammogram image
CIRC [13].



Figure 3.39 Example of
original Mammogram image
NORM [13].

42

Figure 3.40 Example of Mammogram image CALC after applying Gaussian Blurring.

Figure 3.41 Example of Mammogram CIRC image after applying Gaussian Blurring.

Figure 3.42 Example of Mammogram NORM image after applying Gaussian Blurring.

**3. Median Filtering**

Median Filtering is a non-linear filtering technique that replaces each pixel value with the median value of the surrounding pixel neighborhood. This method is especially effective for removing salt-and-pepper noise from an image.

- **Applications:**
  - Noise removal while preserving edges.
  - Enhancing image quality in medical imaging.
- **Advantages:**
  - Excellent for preserving edge details while removing noise.



Figure 3.43 Example of Mammogram image CALC after applying Median Filtering.

Figure 3.44 Example of Mammogram CIRC image after applying Median Filtering.

Figure 3.45 Example of Mammogram NORM image after applying Median Filtering.

43

4. **Histogram Equalization**

Histogram Equalization enhances the contrast of an image by redistributing the intensity values of pixels. It works by flattening the histogram of the image so that the intensity levels are spread more evenly. This technique is particularly useful for improving the visibility of details in low-contrast images.

- **Applications:**
  - Enhancing low-contrast images.
  - Improving image quality in satellite and medical imaging.
- **Advantages:**
  - Simple and effective for global contrast enhancement.
- **Disadvantages:**
  - May over-enhance the noise or artifacts in the image.



Figure 3.46 Example of Mammogram image CALC after applying Histogram Equalization.

Figure 3.47 Example of Mammogram CIRC image after applying Histogram. Equalization.

Figure 3.48 Example of Mammogram NORM image after applying Histogram Equalization.

5. **CLAHE Equalization (Contrast Limited Adaptive Histogram Equalization)**

CLAHE is an advanced form of histogram equalization that works on small regions (tiles) of the image. By limiting the contrast amplification, it prevents noise from being over-enhanced. CLAHE is particularly useful for improving local contrast in images.

- **Applications:**
  - Enhancing medical images like X-rays and CT scans.
  - Improving visibility in images with varying lighting conditions.
- **Advantages:**
  - Enhances contrast adaptively in localized regions.
  - Reduces noise amplification compared to standard histogram equalization.



| Figure 3.49 Example of Mammogram image CALC after applying CLAHE Equalization. | Figure 3.50 Example of Mammogram CIRC image after applying CLAHE Equalization. | Figure 3.51 Example of Mammogram NORM image after applying CLAHE Equalization. |

6. **Image Inversion**

Image inversion is a process that transforms an image by replacing each pixel value with its complement. For grayscale images, this involves subtracting each pixel value from the maximum intensity (e.g., 255 for 8-bit images), effectively swapping white and black regions.

- **Applications:**
  - Creating negative images.
  - Highlighting features in medical imaging.
  - Artistic effects in photography.

Figure 3.52 Example of Mammogram image CALC after applying Image Inversion.

Figure 3.53 Example of Mammogram CIRC image after applying Image Inversion.

Figure 3.54 Example of Mammogram NORM image after applying Image Inversion.

We tried all these techniques in the end we chose the combination that gave the best result

| epoch | metrics/precision(B) | metrics/recall(B) | metrics/mAP50(B) | metrics/mAP50-95(B) | model |
|---|---|---|---|---|---|
| 100 | 0.94667 | 0.79405 | 0.88127 | 0.65362 | Gaussian blurring and Inverted Image |
| 100 | 0.86954 | 0.78755 | 0.87447 | 0.65287 | Inverted Image |
| 100 | 0.95931 | 0.82221 | 0.88869 | 0.64495 | Gaussian Blurring |
| 100 | 0.93984 | 0.75838 | 0.84525 | 0.62567 | CLAHE |
| 100 | 0.87009 | 0.72754 | 0.82172 | 0.61399 | Median Filtering |
| 100 | 0.88919 | 0.71725 | 0.82805 | 0.59907 | Inverted Image and Gaussian Blurring |
| 100 | 0.87336 | 0.74369 | 0.79725 | 0.58654 | Histogram Equalization |

Table 3 Best filter combination

*Note: these filters were only used in the deep learning model, none of these filters was applied in the machine learning model.*

In the end we applied Gaussian blurring then Inverted Image techniques because they gave the best model result which will be later discussed in chapter 5.

### 3.1.3 Data Integration and Transformation

The integration of disparate data sources and transformation into a unified format were critical steps in preparing the dataset for model training. These steps involved organizing data, standardizing variables, and ensuring consistency across all sources.

1. **Dataset Organization**

To ensure a coherent and accessible structure, the **MIAS** dataset was organized as follows:

- o **Image Repository:** All images were stored in a dedicated folder (images/), maintaining a clear and consistent directory structure.
- o **Annotations:** For each image, a corresponding annotation file (.txt) was created. These files were named identically to their respective images (e.g., image1.txt for image1.jpg) to establish a direct mapping.

The overall dataset structure was standardized as:

**/dataset/**

**/images/**

**image1.jpg**

**image2.jpg**

**...**

**/labels/**

**image1.txt**

**image2.txt**

**...**

*Note: written in BASH*

2. **Label Transformation**

Each object in the images was labeled using bounding box coordinates in the YOLO annotation format:

**class_id center_x center_y width height**

*Note: written in Arduino*

The variables were normalized to ensure uniformity and interoperability between different datasets:

- o **class_id:** Numerical identifier for the object class, indexed from 0.
- o **center_x:** Horizontal center of the bounding box, normalized by the image width (range: 0 to 1).
- o **center_y:** Vertical center of the bounding box, normalized by the image height (range: 0 to 1).
- o **width:** Bounding box width, normalized by the image width (range: 0 to 1).
- o **height:** Bounding box height, normalized by the image height (range: 0 to 1).

An example annotation:

47

<div align="center">

0  0.5  0.5  0.2  0.3

1  0.7  0.6  0.1  0.2

</div>

3. **Unified Configuration**

To further ensure consistency, a data.yaml configuration file was created to define the dataset's structure and metadata:

<div align="center">

***# Paths to datasets***
**train:** /path/to/train/images
**val:** /path/to/val/images

***# Number of classes in the dataset***
**nc:** <number_of_classes>

***# Class names***
**names:**
0: 'class_name_1'
1: 'class_name_2'
2: 'class_name_3'

</div>

*Note: written in yaml*

This unified structure facilitated seamless integration and preprocessing of data while ensuring compatibility with the YOLO model training pipeline. The rigorous standardization of variables and directory structure enhanced reproducibility and robustness of the training process.

## 3.2 Feature Engineering and Selection

For the **CBIS-DDSM** dataset, all categorical features in the dataset were encoded using one-hot encoding techniques to convert them into a numerical representation suitable for analysis. This encoding process ensures that categorical variables are transformed into a format that machine learning algorithms can effectively utilize for training and prediction.

By applying one-hot encoding, each categorical feature is represented by a binary vector, enabling the model to interpret and leverage these variables accurately during the learning process.

For the **MIAS** dataset, in this step, the raw data from the dataset was preprocessed to enhance model performance. The key processes included:

48

1) **Handling missing values:**

Missing values in the dataset were identified and replaced with zeros (df.fillna(0)). This ensured the consistency and usability of the dataset.

2) **Feature Creation**:

New features, such as HEIGHT and WIDTH, were calculated based on the radius of the detected masses:

- WIDTH = 2 × RADIUS
- HEIGHT = 2 × RADIUS

3) **Encoding Categorical Variables**:

The CLASS column, which contained categorical labels, was encoded numerically for use in machine learning models. The mapping included:

- CALC → 0
- CIRC → 1
- SPIC → 2
- MISC → 2
- ARCH → 3
- ASYM → 4

4) **Directory Creation for Data Preparation**:

Organized the dataset by creating directories for images and labels using Python's Path module to facilitate smooth integration with downstream tasks.

5) **Feature Scaling**:

The image dimensions (width = 640, etc.) were set for uniformity, ensuring compatibility with the model's input requirements.

→By transforming raw data into these meaningful features, the foundation was laid for improved model performance and interpretability.

## 3.3 Ethical Considerations and Data Privacy

This subsection investigates the ethical implications of the project, including considerations related to data privacy, security, and fairness. Ensuring the ethical use of data is crucial, especially when dealing with sensitive medical information such as mammograms. Anonymizing sensitive information is a critical step to protect patient privacy. This involves removing or encoding Personally Identifiable Information (PII) to prevent the identification of individuals from the data.

The datasets used in this project include MIAS and CBIS-DDSM. The Mammographic Image Analysis Society (MIAS) dataset, as described by its original source [12], contains a reduced set of digital mammograms, all anonymized, with detailed metadata including class labels, image dimensions, and lesion characteristics. MIAS was specifically designed for academic and non-commercial research purposes, ensuring compliance with ethical standards and promoting advancements in mammogram analysis. By anonymizing patient details, the dataset safeguards privacy while enabling research in detecting and classifying breast abnormalities.

Similarly, the CBIS-DDSM dataset ensures patient anonymity while still providing valuable metadata for research, such as patient age and breast density [11]. The original developers of CBIS-DDSM obtained the necessary patient consents for the dataset [11]. Moreover, compliance with relevant regulations like the General Data Protection Regulation (GDPR) and the Health Insurance Portability and Accountability Act (HIPAA) is essential. These regulations set strict guidelines on data protection and privacy, requiring that any data used in research must be securely stored and handled to prevent unauthorized access.

By adhering to these ethical considerations, the project ensures the responsible use of sensitive medical information. This compliance not only protects patients but also ensures that the research aligns with legal standards, fostering trust and integrity in the research process. Through such measures, this work contributes to the responsible advancement of medical technologies while upholding the highest standards of ethics.

**CHAPTER 4: MODEL DEVELOPMENT AND ARCHITECTURAL DESIGN**

**4.1 System Architecture Overview**



Figure 4.1 Data Flow Diagram for Breast Cancer Detection System

This diagram represents the data flow in a breast cancer detection system, from data collection to user interface display. The system includes preprocessing techniques for image enhancement, data augmentation for model robustness, YOLO model training and testing, and a web-based interface for user interaction. The workflow incorporates automated inference and database storage of results, ensuring accurate predictions and efficient data management.

**4.2 Model Architecture and Design**

In our initial exploration, we employed several traditional machine learning models, including Random Forest, Logistic Regression, and Decision Trees, to address the task at hand. Random Forest, an ensemble learning method, was utilized for its robustness and ability to handle complex datasets through multiple decision trees. Logistic Regression, a linear model, was chosen for its simplicity and effectiveness in binary classification tasks. Decision Trees were also implemented due to their interpretability and capability to model non-linear relationships. While these models provided valuable insights and baseline performance, we ultimately decided to transition to YOLOv8, a deep learning

51

model, due to its superior precision and accuracy in object detection tasks. YOLOv8's advanced architecture and ability to deliver detailed, class-wise performance metrics aligned more closely with our requirements, making it the optimal choice for achieving higher-quality results.

This subsection focuses on the architecture and design of the YOLOv8 models and machine learning models used in the project. It includes the selection of appropriate algorithms, network architectures, and model complexity. Key considerations such as scalability, interpretability, and the choice of frameworks or libraries are also addressed.

### Overview of YOLOv8 Architectures

YOLO (You Only Look Once), a groundbreaking object detection and image segmentation model, was originally developed by Joseph Redmon and Ali Farhadi at the University of Washington. Since its launch in 2015, YOLO has gained widespread recognition for its exceptional speed and accuracy **[14]**.

YOLOv8, released in 2023 by Ultralytics, builds on this legacy by introducing new features and improvements that enhance performance, flexibility, and efficiency. It supports a comprehensive range of vision AI tasks, making it a versatile choice for various applications **[14]**.

One of the key reasons we chose to use the YOLOv8 model is its ability to provide a detailed class-wise breakdown of performance metrics. This granular insight is crucial for understanding how well the model performs on individual classes, particularly in datasets with a diverse array of object categories. YOLOv8's architecture, which includes CSP (Cross-Stage Partial connections), the CSPDarknet backbone, and the PANet head, ensures a balance between speed, accuracy, and complexity. This makes it well-suited for our needs, as it allows us to evaluate precision, recall, mAP50, and mAP50-95 for each class, ensuring robust and reliable object detection across varying levels of difficulty.

YOLOv8 offers a range of model sizes to cater to diverse performance requirements, enabling us to select the optimal model for our specific use case. Its combination of advanced features, efficiency, and detailed performance analysis capabilities makes it an ideal choice for our object detection tasks.

Below are the details for each model size:

**1. YOLOv8-Nano (yolov8n)**

- **Algorithm Selection and Network Architecture:**
  - Backbone: CSPDarknet
  - Input layer: 640x640 image
  - Convolutional layers (kernel sizes: 3x3, 1x1)
  - MaxPooling layers
  - Cross-Stage Partial (CSP) layers for reduced complexity
  - Neck: PANet (Path Aggregation Network) with upsampling and downsampling layers, and feature map concatenation
  - Head: YOLO Head with detection layers for multiscale object detection

- **Model Complexity:**
  - Parameter Count: ~1.5M
  - Output: Bounding boxes, class scores, and objectness score
- **Use Case:** Real-time applications prioritizing speed, such as mobile and edge devices
- **Scalability and Interpretability:** Optimized for environments with limited computational resources

**2. YOLOv8-Small (yolov8s)**

- **Algorithm Selection and Network Architecture:**

  - Backbone: CSPDarknet
  - CSP layers for computational efficiency and accuracy
  - Neck: PANet for multiscale feature fusion
  - Head: YOLO Head with three detection layers

- **Model Complexity:**
  - Parameter Count: ~3M
  - Output: Bounding boxes, class scores, and objectness score
- **Use Case:** Small-scale applications requiring higher accuracy than Nano models
- **Scalability and Interpretability:** Maintains a balance between speed and accuracy for edge devices.

**3. YOLOv8-Medium (yolov8m)**

- **Algorithm Selection and Network Architecture:**

53

- o Backbone: CSPDarknet with increased depth and width
- o Enhanced CSP layers for complex feature extraction
- o Neck: PANet with refined multiscale feature aggregation
- o Head: YOLO Head with three detection layers
- **Model Complexity:**
  - o Parameter Count: ~7M
  - o Output: Bounding boxes, class scores, and objectness score
- **Use Case:** General-purpose detection tasks balancing speed and accuracy
- **Scalability and Interpretability:** Suitable for moderate-scale applications

### Frameworks and Libraries

YOLOv8 models are implemented using state-of-the-art machine learning frameworks and libraries, including PyTorch and Ultralytics' YOLO implementation. These tools ensure robust model training, scalability, and efficient deployment.

### Summary Table of YOLOv8 Architectures

| Model Size | Backbone (CSPDarknet) | Layers Overview | Parameter Count |
|---|---|---|---|
| **YOLOv8-Nano** | Shallow | Few layers, reduced CSP depth, simple PANet | ~1.5M |
| **YOLOv8-Small** | Shallow | More CSP layers, medium PANet complexity | ~3M |
| **YOLOv8-Medium** | Medium | Deeper CSP layers, more complex PANet and detection head | ~7M |
| **YOLOv8-Large** | Deeper | Deeper CSP, advanced PANet and detection head | ~20M |
| **YOLOv8-Xlarge** | Deepest | Deepest CSP, advanced PANet and detection layers | ~40M |

Table 4 Summary Table of YOLOv8 Architectures

## 4.3 Training Process

The training process for the deep learning models was meticulously designed to ensure optimal performance and generalization. Below is a detailed breakdown of the steps and methodologies employed:

### 1. Data Preparation and Splitting

The dataset was organized and preprocessed to facilitate effective training and evaluation:

- **Dataset Setup**: The dataset was loaded from the specified directory (dataset_edit), and labels were extracted from the corresponding text files.
- **YAML Configuration**: A YAML file (data.yaml) was used to define the dataset structure, including class names and directories.
- **Data Splitting**: The dataset was split into **5-fold cross-validation** sets to ensure robust evaluation. Each fold was divided into training and validation subsets, with directories created for images and labels accordingly.

### 2. Model Initialization

The YOLO model was initialized using pre-trained weights to leverage transfer learning:

- **Model Setup**: The YOLO model was initialized with the weights_path and configured for the object detection task.
- **Training Configuration**: The model was trained using the Adam optimizer, which combines adaptive learning rates and momentum for efficient convergence.

### 3. Training Process

Multiple training runs were conducted to fine-tune the model and identify the best hyperparameters:

- **Epochs**: The training process was conducted with epoch limits ranging between **100 and 150** to balance training time and model performance.
- **Batch Size**: Different batch sizes were experimented with to optimize computational efficiency and memory usage.

55

- **Early Stopping**: A patience value of **10 epochs** was used to halt training if the validation loss did not improve, preventing unnecessary computation.

### 4. Performance and Challenges

- **Overfitting and Underfitting**: No significant issues with overfitting or underfitting were observed, indicating that the model generalized well to unseen data.
- **Convergence**: The model consistently achieved convergence within the specified epoch range, demonstrating the effectiveness of the Adam optimizer and training strategy.

### 5. Results and Output

- **Training Output**: The trained model was saved along with evaluation metrics for each fold, ensuring comprehensive analysis and comparison.
- **Cross-Validation**: The 5-fold cross-validation approach provided a reliable estimate of the model's performance across different subsets of the data.

This structured training process ensured the development of a robust and high-performing object detection model, ready for deployment in real-world applications.

## 4.4 Hyperparameter Tuning

Hyperparameter tuning is a critical step in optimizing the performance of machine learning models. It involves adjusting parameters that are not learned during training, such as the learning rate, batch size, number of epochs, and regularization strength. Below is a detailed explanation of the hyperparameter tuning process employed in this

56

project, incorporating insights from the referenced resource and focusing on the best-performing model.

### 1. Hyperparameter Selection

The following hyperparameters were selected for tuning based on their significant impact on model performance:

- **Batch Sizes**: Values of 8 and 16 were tested to balance memory usage and computational efficiency.

- **Epochs**: Values of 100 and 150 were evaluated to ensure sufficient training time without overfitting.

- **Learning Rate**: Although not explicitly tuned in this process, the learning rate was managed by the Adam optimizer, which adapts it dynamically during training.

### 2. Grid Search Strategy

A **grid search** approach was used to systematically evaluate all possible combinations of hyperparameters:

- **Grid Creation**: All combinations of batch sizes and epochs were generated and tested.

- **Evaluation**: Each combination was evaluated using **5-fold cross-validation** to ensure robust and reliable results. This approach helps mitigate the risk of overfitting to a specific training-validation split.

57

## 3. Training Process

For each hyperparameter combination, the YOLO model was trained and evaluated:

- **Model Initialization**: The model was initialized with pre-trained weights (yolov8l.pt) to leverage transfer learning, which significantly reduces training time and improves performance.

- **Training Configuration**: The Adam optimizer was used for training, as it combines the benefits of adaptive learning rates and momentum for efficient convergence. Early stopping (patience=10) was implemented to halt training if the validation loss did not improve, preventing overfitting and saving computational resources.

- **Fold-wise Training**: Each hyperparameter combination was tested across all folds to ensure consistency in performance and to identify the most generalizable model.

## 4. Insights from Referenced Resource

The referenced resource highlights additional strategies for hyperparameter tuning, such as:

- **Learning Rate Scheduling**: Dynamically adjusting the learning rate during training to improve convergence.

- **Data Augmentation**: Enhancing the dataset with transformations (e.g., rotation, scaling) to improve model robustness.

- **Advanced Optimization Techniques**: Exploring methods like Bayesian optimization for more efficient hyperparameter search.

While these techniques were not explicitly applied in this project, they provide valuable avenues for future optimization.

**5. Best Performing Model**

The grid search identified the best-performing combination of hyperparameters:

- **Batch Size**: 8

- **Epochs**: 150

- **Average mAP**: 0.864

- **Training Time**: 12.5 hours

This configuration achieved the highest mean Average Precision (mAP) across all folds, demonstrating excellent accuracy and generalization. The use of early stopping and cross-validation ensured efficient use of computational resources, and no significant overfitting or underfitting was observed.

**6. Conclusion**

The systematic approach to hyperparameter tuning, combined with insights from advanced techniques, ensured the development of a high-performing and robust object detection model, ready for deployment in real-world applications. The best-performing model, trained with a batch size of 8 and 150 epochs, serves as a strong foundation for future optimizations and scalability.

**4.5 Validation Strategy**

This subsection delineates the validation strategy employed to evaluate the performance of the machine learning model. The primary methodology utilized was **5-fold cross-validation**, a rigorous technique designed to assess the model across multiple subsets of the dataset, thereby providing a more reliable estimation of its generalization capabilities.

To augment the training process, **data augmentation techniques** were applied to the images, enriching the dataset's diversity and enhancing the model's robustness.

The model's performance was quantified using **YOLOv8 metrics**, including **precision**, **recall**, **mAP (mean Average Precision)**, and **F1-score**. These metrics offer a comprehensive evaluation of the model's accuracy in detection and classification, as well as the balance between precision and recall.

The cross-validation approach effectively addressed the **bias-variance tradeoff**, ensuring the model was neither underfitting nor overfitting. Additionally, it facilitated the identification of the optimal model configuration by rigorously testing its performance on unseen data.

This validation strategy underscores the critical importance of assessing generalization performance, thereby ensuring the model's reliability and efficacy when deployed in real-world applications.

## 4.6 Code Implementation

The project was implemented using a well-structured technology stack to ensure efficiency and seamless integration of components. The key technologies and tools utilized are as follows:

- **Programming Language**: The project was developed using **Python version 3.10.12**, leveraging its extensive libraries and frameworks for machine learning, web development, and data processing.

- **Web Framework**: **Flask** was employed as the backend framework to create a RESTful API, enabling communication between the machine learning model and the web application.

- **Frontend Development**: The user interface was designed using **HTML**, **CSS**, and **JavaScript**, ensuring a responsive and intuitive experience for end-users.

- **Database**: **PostgreSQL** was used as the relational database management system to store and manage structured data, providing reliability and efficient data retrieval.

All development and testing were conducted on a local device, ensuring a controlled environment for implementation and evaluation. This combination of technologies facilitated the development of a cohesive system, integrating machine learning capabilities with a user-friendly web application and a robust database backend.

60

# CHAPTER 5: MODEL/SYSTEM EVALUATION AND TESTING PLAN

Regarding the system, we did not conduct tests for Scalability Assessment, Resource Utilization, Responsiveness, or Latency because the project was implemented and executed on a local device, in the future we will consider these options.

## 5.1 Evaluation Metrics

This subsection details the specific metrics employed to assess the performance of the models and the system. These metrics were essential for evaluating both the effectiveness of the models and the overall performance of the system.

### 5.1.1 Model Performance Metrics

A key component of the output is the detailed class-wise analysis of performance metrics. This level of detail is particularly valuable for assessing how effectively the model performs on individual classes, especially in datasets encompassing a wide variety of object categories. For every class in the dataset, the following metrics are provided:

- **Class**: This specifies the name of the object category, such as "CALC," "NORM," or "CIRC."

- **Images**: This indicates the number of images in the validation set that include the specified class.

- **Instances**: This reflects the total number of occurrences of the class across all images in the validation set.

- **Box(P, R, mAP50, mAP50-95)**: This set of metrics evaluates the model's object detection capabilities:

  - **P (Precision)**: Measures the accuracy of the model's detections by showing the proportion of correct identifications.

  - **R (Recall)**: Assesses the model's ability to detect all instances of the class within the images.

o **mAP50**: Represents the mean average precision at an intersection over union (IoU) threshold of 0.50, focusing on the model's performance for "easier" detections.

o **mAP50-95**: Averages the mean average precision across IoU thresholds ranging from 0.50 to 0.95, offering a holistic evaluation of the model's accuracy across varying levels of detection complexity.

This breakdown helps in pinpointing strengths and weaknesses in the model's performance for each specific class.

*5.1.2 Other Relevant Metrics*

No additional or domain-specific metrics were used beyond those provided by YOLOv8. The predefined metrics effectively captured the necessary performance insights, aligning with the project's objectives and ensuring a thorough evaluation of the model's capabilities.

## 5.2 Evaluation Methodology

The evaluation methodology employed in this project utilized **5-fold cross-validation** as the primary technique to assess model performance and generalization. The dataset was partitioned into five equal subsets, with the model trained on four subsets and validated on the remaining one in each iteration. This approach ensured that all data points were used for both training and validation across different folds, providing a robust and unbiased evaluation of the model's capabilities. The **holdout method** and separate **train-test splits** were not employed, as cross-validation alone offered a comprehensive and balanced assessment of the model's performance.

## 5.3 Model Performance Evaluation

This section focuses on evaluating the performance of the models across different datasets, ensuring they generalize well and perform effectively in real-world scenarios.

*5.3.1 Performance on Training Data*

The performance on the training data was assessed to determine the model's ability to capture underlying patterns and relationships. The primary objective was to measure how well the model fits the training data and learns from the patterns within it. Evaluation metrics such as **accuracy**, **precision**, **recall**, **F1 score**, and **loss** were

utilized to understand the model's performance in terms of both prediction accuracy and learning effectiveness. Visualization tools, including YOLOv8's default charts (e.g., loss and accuracy curves), provided insights into the model's convergence during training.

*5.3.2 Performance on Validation Data*

The model's generalization performance was evaluated on unseen validation data to ensure robustness and prevent overfitting. The objective was to assess how well the model performs on data it had not encountered during training, which is critical for evaluating its ability to generalize. Metrics such as **validation loss** and **accuracy** were used, with a focus on generalization. **K-fold cross-validation** was employed, splitting the data into five subsets and rotating the training and validation sets to ensure robust performance across different portions of the dataset. Additionally, **early stopping** (patience=10) was implemented to halt training if the model's performance on the validation set did not improve after 10 consecutive epochs, mitigating the risk of overfitting.

| epoch | Batch Size | metrics/precision(B) | metrics/recall(B) | Avg mAP50 | Avg mAP50-95 | model |
|---|---|---|---|---|---|---|
| 100 | BATCH16 | 0.936264 | 0.874996 | 0.915228 | 0.702222 | YOLOv8n |
| 100 | BATCH8 | 0.963418 | 0.895576 | 0.929538 | 0.73526 | YOLOv8n |
| 150 | BATCH16 | 0.968632 | 0.94474 | 0.965226 | 0.76441 | YOLOv8n |
| 150 | BATCH8 | 0.97572 | 0.92428 | 0.951576 | 0.74697 | YOLOv8n |
| 100 | BATCH16 | 0.941932 | 0.82766 | 0.898058 | 0.682484 | YOLOv8n |
| 100 | BATCH8 | 0.93976 | 0.85937 | 0.917864 | 0.675212 | YOLOv8n |
| 150 | BATCH16 | 0.98079 | 0.924754 | 0.960214 | 0.757896 | YOLOv8n |
| 150 | BATCH8 | 0.973944 | 0.913524 | 0.949862 | 0.764594 | YOLOv8n |
| 100 | BATCH16 | 0.97405 | 0.87908 | 0.930717 | 0.75601 | YOLOv8s |
| 100 | BATCH8 | 0.954071 | 0.880967 | 0.933696 | 0.719581 | YOLOv8s |
| 150 | BATCH16 | 0.979043 | 0.938867 | 0.968229 | 0.769173 | YOLOv8s |
| 150 | BATCH8 | 0.947817 | 0.886443 | 0.941343 | 0.69773 | YOLOv8s |
| 100 | BATCH16 | 0.939912 | 0.845018 | 0.906272 | 0.659844 | YOLOv8m |
| 100 | BATCH8 | 0.935334 | 0.882754 | 0.93016 | 0.67228 | YOLOv8m |
| 150 | BATCH16 | 0.968644 | 0.91308 | 0.944042 | 0.74176 | YOLOv8m |
| 150 | BATCH8 | 0.945214 | 0.88727 | 0.93602 | 0.699118 | YOLOv8m |

Table 5 average for each combination of epochs and batch size for each different model

The table presents a comparison of model performance **before** and **after preprocessing**. The first four models were evaluated using the raw, unprocessed dataset, establishing a baseline for performance. The remaining models were assessed after applying preprocessing techniques, specifically **Gaussian Blurring** and **Image Inversion**. This structured comparison highlights the impact of preprocessing on model

63

accuracy and generalization, providing insights into the effectiveness of these techniques in enhancing dataset quality and improving overall performance.

The best model was

| epoch | Batch Size | metrics/precision(B) | metrics/recall(B) | metrics/mAP50(B) | metrics/mAP50-95(B) | model |
|---|---|---|---|---|---|---|
| 150 | BATCH8 | 0.98629 | 0.94979 | 0.98985 | 0.81547 | YOLOv8n |
| 150 | BATCH8 | 0.99427 | 0.9469 | 0.98205 | 0.78781 | YOLOv8n |
| 150 | BATCH16 | 0.99403 | 0.97458 | 0.9874 | 0.83265 | YOLOv8s |
| 150 | BATCH16 | 0.98424 | 0.93647 | 0.95631 | 0.79197 | YOLOv8m |

Table 6 best model information

*5.3.3 Performance on Test Data*

The model's final performance was tested on a separate test dataset to evaluate its ability to generalize to completely unseen data. The objective was to assess how well the model performs in real-world scenarios, using metrics such as **accuracy**, **precision**, **recall**, **F1 score**, and **loss** to provide a comprehensive evaluation. This final testing phase ensured that the model would be effective and reliable when deployed in practical applications.

Here are the results of the machine learning model

Random Forest Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| **0** | 0.82 | 0.91 | 0.86 | 113 |
| **1** | 0.91 | 0.82 | 0.86 | 126 |
| **Accuracy** | | | 0.86 | 239 |
| **macro avg** | 0.86 | 0.86 | 0.86 | 239 |
| **weighted avg** | 0.87 | 0.86 | 0.86 | 239 |

Table 7 Random Forest Classification Report

Logistic Regression Classification Report:

| | | | | |
|---|---|---|---|---|
| **0** | **0.81** | **0.85** | **0.83** | **113** |
| **1** | 0.86 | 0.82 | 0.84 | 126 |
| **Accuracy** | | | 0.83 | 239 |
| **macro avg** | 0.83 | 0.83 | 0.83 | 239 |
| **weighted avg** | 0.83 | 0.83 | 0.83 | 239 |

Table 8 Logistic Regression Classification Report

Decision Tree Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| **0** | 0.77 | 0.86 | 0.81 | 113 |
| **1** | 0.86 | 0.77 | 0.81 | 126 |
| **Accuracy** | | | 0.81 | 239 |
| **macro avg** | 0.81 | 0.81 | 0.81 | 239 |
| **weighted avg** | 0.82 | 0.81 | 0.81 | 239 |

Table 9 Decision Tree Classification Report

## 5.4 GUI Acceptance Testing Evaluation

As part of the system evaluation, the Graphical User Interface (GUI) underwent rigorous acceptance testing to ensure its functionality and usability. One key aspect of the testing involved verifying the system's behavior when users uploaded incorrect or non-medical images. The GUI was designed to handle such cases by predicting these images as **normal**, ensuring robustness and preventing misleading results. Additionally, the interface was tested for user-friendliness, with feedback from users indicating that the GUI was intuitive, easy to navigate, and visually appealing. The overall design and functionality of the UI were well-received, meeting the project's usability standards and enhancing the user experience.

**CHAPTER 6: MODEL DEPLOYMENT AND INTEGRATION**

**6.1 Technology Stack for AI Systems**

The project employs a comprehensive and robust technology stack, incorporating the following key components:

- **Backend Framework**: Flask is utilized to serve the AI model and manage request handling.
- **Frontend Technologies**: HTML and CSS are employed for the development of user interfaces.
- **Database**: PostgreSQL is implemented for the storage of processed image data and associated predictions.
- **AI Model**: YOLO (You Only Look Once), integrated via the Ultralytics library, is leveraged for image detection and classification tasks.
- **Programming Language**: Python is used for backend development and AI model operations.

This integrated stack ensures a scalable, efficient, and reliable system for the project's objectives.

**6.2 Data Acquisition, Preprocessing, and Database Integration**

This section outlines the processes involved in acquiring and preprocessing the data, as well as integrating the results into the database. The steps ensure that the data is properly prepared for model training and that all outcomes, including predictions and annotations, are systematically stored for future reference and analysis.

**1. Data Input**

- **Image Upload**: Images are uploaded through a user-friendly web interface, ensuring a seamless data input process.

**2. Data Preprocessing**

- **Resizing**: Images are resized to match the input dimensions required by the YOLO model.

66

- **Enhancement Techniques**: Inversion and Gaussian blurring are applied to enhance image quality and improve processing accuracy.
- **Compatibility Checks**: Data compatibility and consistency are ensured to facilitate accurate predictions by the model.

**3. Database Integration**

- **Database Setup**: A **PostgreSQL** database is configured to store all relevant data, including uploaded images, preprocessed images, and prediction results.
- **Schema Design**: The database schema includes the following fields:
    - **ID (Primary Key)**: A unique identifier for each image.
    - **Image_name_with_id**: The image name with its id.
    - **Original_Image_data**: The image file as uploaded by the user.
    - **Preprocessed_Image_data**: The image after preprocessing, ready for model input.
    - **Prediction**: The predicted class label with the confidence value (e.g., "normal" or "abnormal").
    - **Label_index**: The numerical index corresponding to the predicted class.
    - **Label_name**: The predicted class label
    - **Annotations**: If the image is predicted as "abnormal," the bounding box coordinates (x, y, width, height) of the tumor are stored. For "normal" predictions, this field is set to null.
    - **Timestamp**: The date and time when the image was processed and stored.
- **Data Insertion**: After preprocessing and model inference, the original image, preprocessed image, prediction results, annotations (if applicable), and timestamp are saved to the database using [specify methods, e.g., "SQL queries" or "ORM (Object-Relational Mapping) tools"].

**4. Tools and Libraries**

- **Data Processing**: Libraries such as **Pandas**, **NumPy**, and **OpenCV** are used for data manipulation and preprocessing.

67

- **Database Interaction**: The **Psycopg2** library is utilized to connect Python with PostgreSQL and manage database operations.

## 6.3 Model Deployment Strategies

The deployment of the model was executed on a local device to ensure efficient testing and development. The following strategies were implemented:

- **Deployment Framework**: Flask was utilized as the deployment platform, serving as the host for the model API. This lightweight and flexible framework enabled seamless integration and interaction with the AI model.
- **Local Deployment**: The system was deployed locally, running on port 8080, to facilitate testing and development on the local device. This approach allowed for rapid iteration and validation of the model's performance in a controlled environment.
- **Inference Pipeline**:
  - **Image Processing**: Input images were processed and passed to the YOLO model for prediction.
  - **Result Generation**: The model generated predictions, including bounding boxes and labels, which were then returned to the user in a structured and interpretable format.

This local deployment strategy provided a robust foundation for development and testing, ensuring the model's reliability and accuracy before potential scaling to cloud-based environments.

## 6.4 User Interface Design

The user interface (UI) of the application is designed to be intuitive, user-friendly, and visually appealing. The following components and features were implemented:

- **Web Pages**:
  - **Home Page**: Serves as the landing page, introducing the application and its functionality to users.
  - **Upload Page**: Provides a straightforward interface for users to upload images for processing.

- o **Results Page**: Displays the prediction outcomes, including annotated images with bounding boxes, labels, and confidence scores, in a clear and organized manner.

- **Styling**:
  - o CSS is employed to enhance the usability and aesthetics of the interface, ensuring a clean and responsive layout.
  - o The design prioritizes a seamless user experience, with intuitive navigation and visually consistent elements across all pages.

This UI design ensures that users can interact with the application effortlessly, making it accessible and engaging for a wide range of audiences.
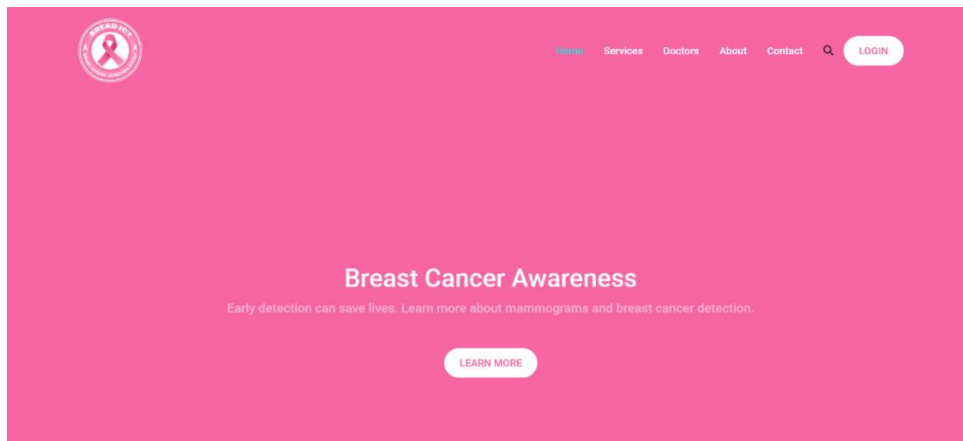

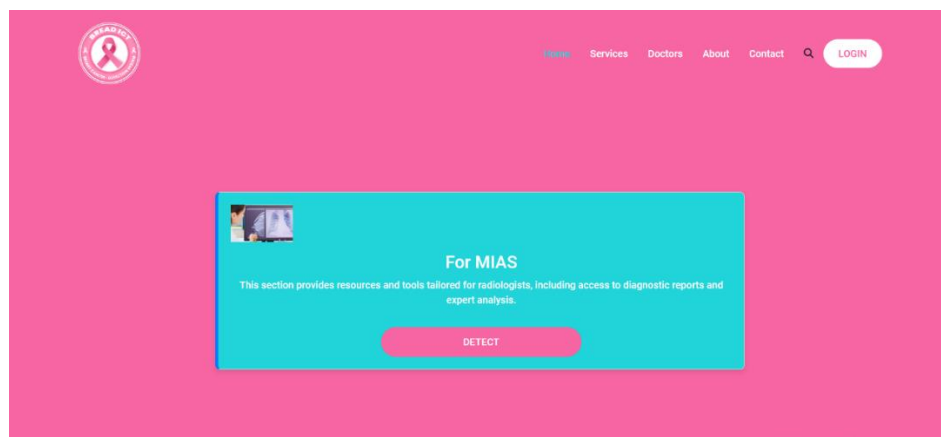
Figure 6.1 UI design main page
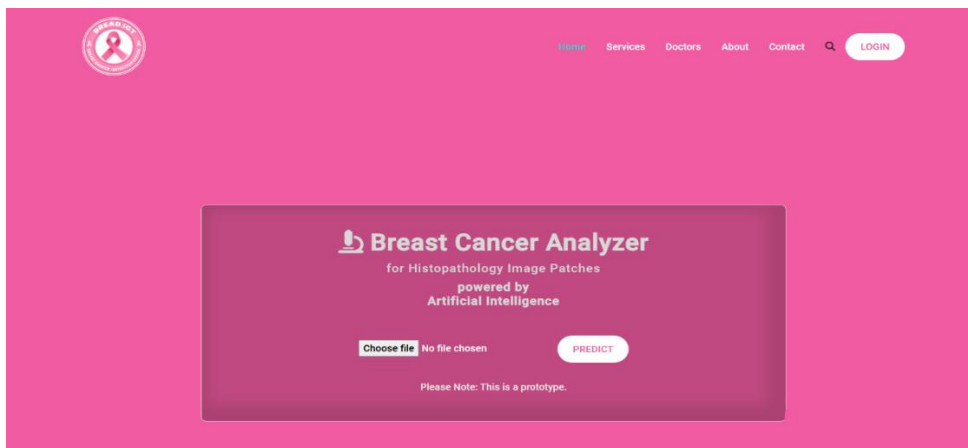


Figure 6.2 UI design detect page
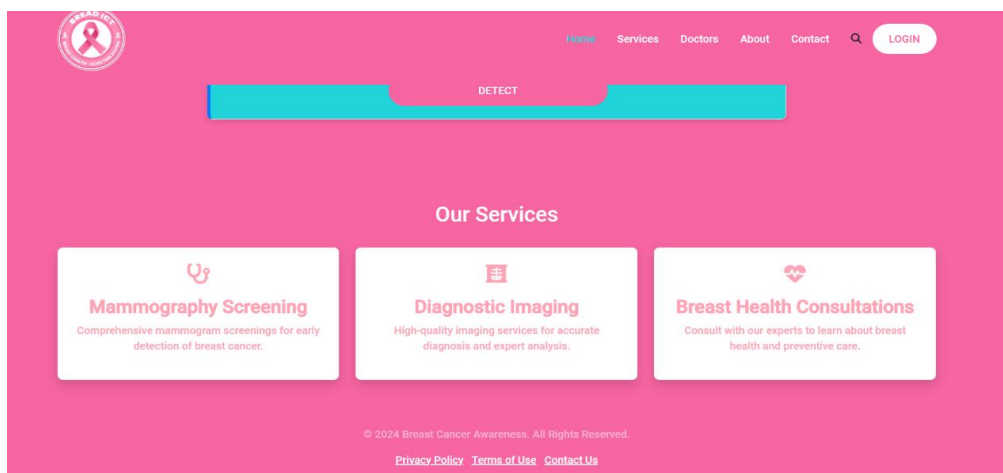
Figure 6.3 UI design predict page
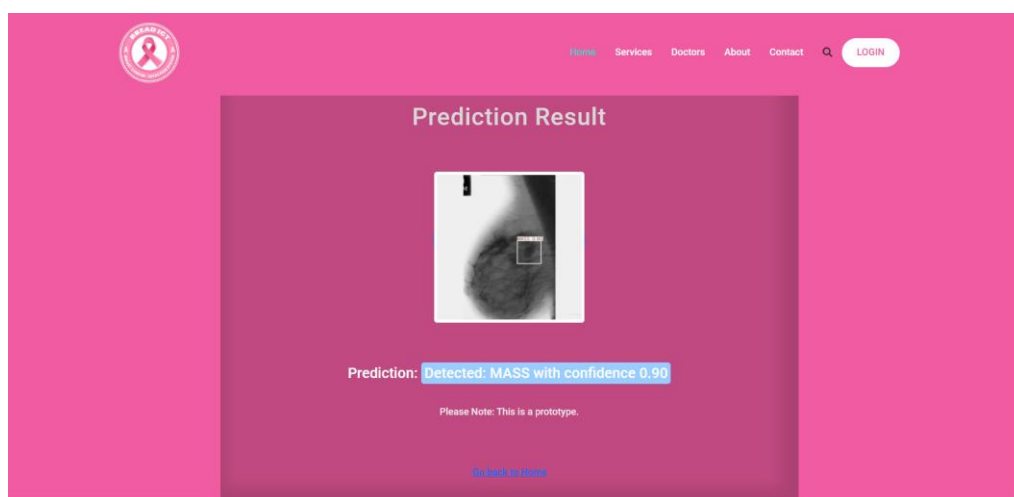


Figure 6.4 UI design services page



Figure 6.5 UI design prediction result page

# CHAPTER 7: CONCLUSION AND FUTURE WORK

This chapter offers a detailed summary of the project's results and key takeaways. It highlights the accomplishments, assesses the outcomes, and reflects on the insights obtained. The chapter wraps up by examining the lessons learned and suggesting possible avenues for future development and improvement.

## 7.1 Summary and Project Evaluation

This section summarizes the project's accomplishments and evaluates its success against initial objectives. The project developed a breast cancer detection system using the YOLO model, trained on the CBIS-DDSM and MIAS datasets. Preprocessing techniques like Gaussian Blurring, CLAHE, and Image Inversion, along with data augmentation methods such as RandomResizedCrop and HorizontalFlip, were applied to enhance model performance.

A **Flask-based web application** enabled users to upload images and view predictions, while results were stored in a **PostgreSQL database** for traceability. Key achievements include high model accuracy, a seamless workflow, and a robust database system. The project successfully met its goals, demonstrating the potential of AI in medical imaging.

## 7.2 Future Directions

This section reflects on the lessons learned during the project's implementation, addressing the challenges encountered, successes achieved, and areas for improvement. Furthermore, it outlines future directions and recommendations for enhancing and expanding the project.

To build upon the current accomplishments, the following future initiatives are proposed:

1. **Dataset Expansion**: Incorporate additional datasets to enhance the model's robustness and accuracy, ensuring its effectiveness across diverse patient demographics and imaging conditions.
2. **Accessibility**: Expand the system's availability to all individuals, not just radiologists and doctors, making it accessible to patients and the general public for early detection and awareness.

71

3. **Scalability**: Improve the system's scalability to accommodate larger datasets and increased user traffic, ensuring its feasibility for widespread use.

4. **Comprehensive Diagnosis System**: Develop the system into a **breast cancer diagnosis system** by integrating advanced diagnostic features, such as risk assessment and treatment recommendations.

5. **User Interface Enhancements**: Refine the web application's interface to ensure it is intuitive and user-friendly for individuals with varying levels of technical expertise.

6. **Model Optimization**: Explore advanced optimization techniques, including hyperparameter tuning and transfer learning, to further enhance the model's performance and reliability.

7. **Real-Time Processing**: Implement real-time image processing and prediction capabilities to deliver immediate results to users.

8. **Public Awareness and Education**: Collaborate with healthcare organizations to promote the system's use for public awareness and early detection of breast cancer.

By pursuing these future directions, the project can evolve into a more comprehensive and impactful tool for breast cancer detection and diagnosis, ultimately contributing to improved healthcare outcomes and advancing the field of AI-driven medical diagnostics.

# REFERENCES

1. Breast Cancer Research Foundation, https://www.bcrf.org/breast-cancer-statistics-and-resources/ Accessed date: 15/3/2024.

2. International Agency For Research On Cancer-World Health Organization, https://www.iarc.who.int/news-events/current-and-future-burden-of-breast-cancer-global-statistics-for-2020-and-2040/, Accessed date: 15/3/2024

3. World Health Organization , https://www.who.int/news/item/03-02-2023-who-launches-new-roadmap-on-breast-cancer , Accessed date: 15/3/2024.

4. King Hussien Cancer Center, https://www.khcc.jo/en/cancer-types/breast-cancer , Accessed date: 17/3/2024

5. Oregon Health & Science University ,
https://news.ohsu.edu/2020/06/30/understanding-the-aggressive-breast-cancers-missed-by-mammogram-screening, Accessed date: 20/3/2024

6. Karssemeijer, N. & te Brake, G. M. Detection of stellate distortions in mammograms. *IEEE Trans. Med. Imaging* **15**, 611–619 (1996), Accessed date: 8/4/2024

7. Yang, S. K. et al. Screening mammography—detected cancers : Sensitivity of a computer-aided detection system applied to full-field digital mammograms. *Radiology* **244**, 104–111 (2007), Accessed date: 8/4/2024

8. Rangayyan, R. M., Mudigonda, N. R. & Desautels, J. E. Boundary modelling and shape analysis methods for classification of mammographic masses. *Med. Biol. Eng. Comput.* **38**, 487–496 (2000), Accessed date: 8/4/2024.

9. Mudigonda, N. R., Rangayyan, R. M. & Desautels, J. E. Gradient and texture analysis for the classification of mammographic masses. *IEEE Trans. Med. Imaging* **19**, 1032–1043 (2000), Accessed date: 8/4/2024.

10. Görgel, P., Sertbas, A. & Uçan, O. N. Computer-aided classification of breast masses in mammogram images based on spherical wavelet transform and support vector machines. *Expert Syst* **32**, 155–164 (2015), Accessed date: 8/4/2024.

11. Cancer-Imaging-Archive,

    https://www.cancerimagingarchive.net/collection/cbis-ddsm/ , Accessed date: 8/4/2024.

12. The mini-MIAS database of mammograms,

    http://peipa.essex.ac.uk/info/mias.html , Accessed date: 15/11/2024.

13. The mini-MIAS database of mammograms,

    https://www.kaggle.com/datasets/kmader/mias-mammography, Accessed date: 15/11/2024.

14. YOLO model documentation, https://docs.ultralytics.com/, Accessed date: 30/11/2024.