

UNIVERSITY OF SINDH

Laar Campus

PROJECT REPORT

Breast Cancer Detection Using Deep Learning

Submitted by: Raimal Raja Kolhi

Supervisor: Assistant Professor Dileep Kumar

Date: 31 January 2026

ABSTRACT

Breast cancer is one of the leading causes of cancer-related deaths among women worldwide, with early detection being critical for successful treatment outcomes. Manual diagnosis of breast cancer from histopathology slides is time-consuming, labor-intensive, and prone to human error, creating a significant need for automated diagnostic solutions. This project addresses this challenge by developing an automated breast cancer detection system using deep learning techniques to classify Invasive Ductal Carcinoma (IDC) in breast histopathology images. The system leverages the EfficientNetB0 architecture through transfer learning, trained on the Breast Histopathology Images dataset containing over 277,524 images from Kaggle. The model was developed using a two-phase training approach: first training with a frozen base model for 6 epochs, followed by fine-tuning the last 30 layers with a reduced learning rate. Data augmentation techniques including rotation, flips, and zooming were applied to improve model generalization. The project was implemented in a cloud-based Google Colab environment for training and deployed locally using a Flask web framework. The final model achieved a test accuracy of 78.95% with an AUC score of 0.8552, demonstrating the potential for automated breast cancer detection. The web-based interface allows medical professionals to upload histopathology images and receive instant predictions with confidence levels, providing a practical tool for assisting in clinical diagnosis.

TABLE OF CONTENTS

List of Figures and Tables	4
1. Project Overview	5
2. Introduction	5
2.1 Problem Statement	5
2.2 Solution	5
3. Environment Setup	6
3.1 Kaggle API	6
3.2 Cloud Development (Google Colab)	6
3.3 Local Development (Flask App)	7
4. Implementation	8
4.1 Data Preprocessing and Training	8
4.2 Local Deployment Structure	9
5. Results and Dashboard	9
5.1 Performance Metrics	9
5.2 Web Interface	10
6. Component Versions	10

LIST OF FIGURES AND TABLES

Tables

Table 1:	Project Information	5
Table 2:	Cloud Storage Structure	7
Table 3:	Local File Directory Structure	7
Table 4:	Two-Phase Training Approach	8
Table 5:	Performance Metrics	10
Table 6:	Component Versions	10

1. PROJECT OVERVIEW

This project, titled "Breast Cancer Detection", was developed using Deep Learning AutoML techniques in the domain of medical imaging and computer vision. The dataset used is the "Breast Histopathology Images" obtained online from Kaggle using the Kaggle API, containing over 277,524 images. The primary objective of this project is to develop an automated system for detecting Invasive Ductal Carcinoma (IDC) in breast histopathology images with a target accuracy of 90% or higher.

Table 1: Project Information

Project Title	Breast Cancer Detection Using Deep Learning
Domain	Medical Imaging and Computer Vision
Dataset	Breast Histopathology Images (Kaggle)
Dataset Size	277,524 images
Target Accuracy	90%+
Architecture	EfficientNetB0 (Transfer Learning)

2. INTRODUCTION

2.1 Problem Statement

Manual diagnosis of breast cancer from histopathology slides is time-consuming and prone to human error. The traditional approach requires pathologists to examine tissue samples under microscopes, a process that is not only labor-intensive but also subject to inter-observer variability. This creates a critical need for automated, consistent, and accurate diagnostic tools that can assist medical professionals in making timely and reliable diagnoses.

2.2 Solution

A Deep Learning model based on the EfficientNetB0 architecture is employed to classify tissue patches as IDC positive or negative. This solution leverages transfer learning to benefit from pre-trained weights while adapting to the specific characteristics of breast histopathology images. The model is designed to provide rapid, consistent predictions that can serve as a second opinion or screening tool for pathologists.

3. ENVIRONMENT SETUP

The project followed a three-phase development lifecycle: creating Kaggle API integration, cloud-based training, and local-based deployment. This approach ensures efficient development while maintaining reproducibility and scalability.

3.1 Kaggle API

To access the dataset programmatically, a Kaggle account was created and configured. The process involves creating an account on Kaggle, navigating to settings, scrolling down, and generating a new API token named "BreastCancerDetection". The API key and username are then specified in the "breast_cancer_colab_training.ipynb" notebook to enable online dataset access, eliminating the need for manual downloading. When the new token is generated, Kaggle automatically creates a kaggle.json file which can be placed directly into the Colab environment for seamless dataset download.

3.2 Cloud Development (Google Colab)

The model training was performed in Google Colab, leveraging cloud-based GPU resources for faster computation. Google Drive was mounted to maintain a persistent file structure for saving trained models and results, ensuring that training outputs are preserved across sessions.

Table 2: Cloud Storage Structure

Directory	Contents
/models/	best_model.keras, model_config.pkl
/results/	Confusion matrices, training history plots

3.3 Local Development (Flask App)

After training, the model files were downloaded to a local machine for deployment within a Flask web framework. This enables the creation of a user-friendly web interface for making predictions on new histopathology images.

Local File Directory Structure:

Table 3: Local File Directory Structure

Directory/File	Description
----------------	-------------

BreastCancerDetection/	Root project directory
flask_app/	Flask application directory
app.py	Main Flask application file
templates/	HTML templates directory
index.html	Main web interface
static/	Static files directory
css/style.css	Stylesheet
js/script.js	JavaScript file
models/	Trained model files
breast_cancer_model.h5	Model from Colab
model_config.pkl	Configuration from Colab
uploads/	Auto-generated upload directory
requirements.txt	Python dependencies

4. IMPLEMENTATION

4.1 Data Preprocessing and Training (Colab)

Augmentation:

Data augmentation techniques such as rotation, horizontal and vertical flips, and zooming were applied using ImageDataGenerator to improve model generalization and reduce overfitting. These transformations help the model learn robust features that are invariant to minor variations in image orientation and scale.

Model Architecture:

A transfer learning approach was employed using EfficientNetB0 as the base model, followed by Global Average Pooling, Batch Normalization, and Dropout layers (with rates of 0.5 and 0.3) to prevent overfitting. EfficientNetB0 was chosen for its balance between model complexity and performance, making it suitable for deployment in resource-constrained environments.

Two-Phase Training:

Table 4: Two-Phase Training Approach

Phase	Description	Details
Phase 1	Frozen Base Model	Training for 6 epochs with frozen EfficientNetB0 base
Phase 2	Fine-tuning	Fine-tuning last 30 layers with reduced learning rate

4.2 Local Deployment Structure

The local environment was configured to match the training environment's dependencies to avoid compatibility errors. A virtual environment was created within the folder structure, and specific versions of required packages were installed using either a dedicated setup script (requirements.ps1) or requirements.txt. This ensures that the model performs as expected in the deployment environment, maintaining consistency between training and inference.

5. RESULTS AND DASHBOARD

5.1 Performance Metrics

The model achieved high performance, with metrics recorded in the model configuration. The initial test accuracy of 78.95% demonstrates promising results, with the potential to reach the target accuracy of 90%+ by increasing the training epochs to 20. The AUC score of 0.8552 indicates strong discriminative ability between IDC positive and negative tissue samples.

Table 5: Performance Metrics

Metric	Value	Notes
Test Accuracy (Initial)	78.95%	With 6 epochs training
Target Accuracy	90%+	Achievable with 20 epochs
AUC Score	0.8552	Strong discriminative ability

5.2 Web Interface

A responsive dashboard was developed to allow users to drag and drop histopathology images for analysis. The Flask backend processes the uploaded image using OpenCV, runs it through the EfficientNet model, and returns the probability of IDC presence along with a confidence level. This user-friendly interface makes the technology accessible to medical professionals without requiring technical expertise in machine learning or programming.

6. COMPONENT VERSIONS

The project was built using specific versions of libraries and frameworks to ensure compatibility and reproducibility. The following table lists all major components and their versions:

Table 6: Component Versions

Component	Version
Python	3.10.0
TensorFlow	2.16.1
Keras	3.0

Flask	3.0.0
OpenCV	4.8.1.78
NumPy	2.0