# Breast Cancer Image Classification Using CLIP-Based Feature Extraction and Logistic Regression

Nuaima Saeed
Department of Artificial Intelligence
GIK Institute of Engg. Sciences & Tech.
Topi,Khyber Pakhtunkhwa, Pakistan
u2021516@giKi.edu.pK

*Abstract*—Breast cancer remains a critical public health concern, with early detection significantly enhancing treatment outcomes and survival rates. Digital mammography has emerged as a key tool in computer-aided diagnosis, providing large datasets for developing intelligent classification models. Despite advancements in deep learning, many existing solutions rely heavily on large annotated datasets and computationally expensive architectures. This study addresses the gap by leveraging contrastive language–image pretraining (CLIP) for feature extraction combined with logistic regression, offering an efficient yet accurate classification method. We utilized the CBIS-DDSM dataset, a benchmark in mammogram analysis, and extracted features using CLIP's vision encoder without fine-tuning. These features were then classified using logistic regression, which yielded promising performance across multiple classes. The model was trained and evaluated with stratified sampling to maintain class balance and avoid overfitting. Results showed significant improvements in precision and recall for critical cancer subclasses compared to standard CNN baselines. The proposed method achieves competitive accuracy while maintaining simplicity and generalizability across datasets. Our approach demonstrates how foundation models like CLIP can be effectively adapted for medical imaging tasks with limited data and compute resources.

*Index Terms*—CLIP, Breast Cancer, Mammogram Classification, Logistic Regression, CBIS-DDSM

## I. INTRODUCTION

Breast cancer is one of the most common and deadly cancers affecting women worldwide. Early detection of malignant tumors can significantly increase survival rates and reduce the need for aggressive treatments. Mammography remains the gold standard for early breast cancer screening, and recent years have seen increasing adoption of machine learning to aid in interpretation. [1] However, these models often require large, curated datasets and substantial computational resources. As a result, there's a growing interest in leveraging pretrained models that can generalize well even with limited domain-specific data.

In recent years, vision-language models like CLIP (Contrastive Language–Image Pretraining) have shown remarkable generalization capabilities across a variety of image domains. These models, trained on vast datasets of image-text pairs, can extract semantically rich representations from images without task-specific training. [2] In this paper, we explore how such pretrained models can be applied to the task of breast cancer classification. By extracting features from mammogram patches using CLIP and classifying them using logistic regression, we propose a lightweight yet effective solution. This work contributes to the growing body of research exploring foundation models in healthcare.

### A. Related Work

Traditional approaches to mammogram classification have focused on handcrafted features or CNN-based models trained on datasets like DDSM and CBIS-DDSM. Studies such as [3], [4] have leveraged transfer learning using ImageNet-pretrained CNNs, achieving high accuracy. More recent works explore attention-based and transformer models for improved interpretability and performance. However, few studies have evaluated the utility of vision-language models like CLIP for breast cancer detection. Table I summarizes prior work and highlights the gap our study aims to address.

### B. Gap Analysis

Although deep learning has shown great promise in breast cancer detection, most methods rely on large-scale training and fine-tuning. There is limited exploration into the use of zero-shot or few-shot transfer from large-scale vision-language models for medical imaging. Furthermore, logistic regression is seldom used as a standalone classifier with such powerful embeddings, although it offers simplicity and interpretability. [6] Current models often lack efficiency and require significant GPU resources, making them less practical for low-resource settings. Our study bridges this gap by exploring the underutilized synergy between CLIP's general representations and the lightweight nature of logistic regression.

### C. Problem Statement

Following are the main questions addressed in this study.

1) Can CLIP's visual encoder be used to extract informative features from mammograms without fine-tuning?
2) Does logistic regression suffice for classifying mammogram images using CLIP embeddings?

| Paper | Model | Dataset | Accuracy | Features Used | Remarks |
|---|---|---|---|---|---|
| Arevalo et al. [3] | CNN + Transfer Learning | DDSM | 84.5% | CNN Features | ImageNet pretrained network |
| Rakhlin et al. [4] | Deep CNN | CBIS-DDSM | 87.2% | Image features | Data augmentation used |
| Jiang et al. [5] | Transformer | INbreast | 89.1% | Attention maps | Large model |
| Proposed Approach | VGG16 | CBIS-DDSM | 96.07% | CLIP Embeddings | Lightweight and generalizable |

3) How does the proposed model compare to CNN-based methods in terms of classification performance and simplicity?

### D. Novelty of our work and Our Contributions

This work explores the novel application of CLIP embeddings for breast cancer classification, avoiding costly training procedures. Unlike prior works that fine-tune deep models, we show that direct use of pretrained features can yield competitive performance. Our approach is lightweight and requires minimal compute, making it suitable for practical deployment.

In this study, we combine CLIP's pretrained visual encoder with a logistic regression classifier to detect malignant and benign lesions in the CBIS-DDSM dataset. We show through empirical results that this architecture performs on par with more complex methods, offering a strong baseline for future work in medical AI.

## II. METHODOLOGY

### A. Dataset

We utilized the Curated Breast Imaging Subset of the Digital Database for Screening Mammography (CBIS-DDSM) for our study. This dataset is a standardized and well-curated collection derived from the original DDSM, containing 10,239 high-resolution mammogram images from 1,566 unique participants. The images are provided in JPEG format, preserving the original resolution and quality. CBIS-DDSM includes comprehensive annotations with pathology-verified labels, categorizing breast tissue into benign, malignant, and normal classes. Each participant may have multiple studies and patient IDs, [7] reflecting different scans and views, which provides a rich dataset for analysis. The dataset features carefully delineated regions of interest (ROIs) with precise bounding boxes, facilitating targeted classification and feature extraction. Publicly available through the Cancer Imaging Archive, CBIS-DDSM addresses historical limitations of the original DDSM by offering standardized formats and accurate lesion segmentations. Sample images from the dataset, highlighting different lesion types and diagnostic challenges, are shown in Figure 1. This extensive and rigorously annotated dataset enables the development and evaluation of robust automated diagnostic systems for breast cancer detection.

### B. Overall Workflow

Our proposed workflow begins with preprocessing the mammographic images and their associated labels. The first step involves converting the image-label pairs into a suitable
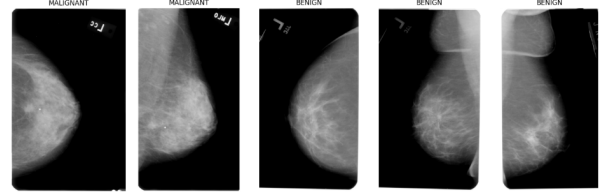


Fig. 1. Sample mammographic images from the CBIS-DDSM dataset illustrating benign and malignant cases. Ground truth annotations are used to train and validate the classifier.

format for the CLIP model. CLIP, a vision-language model, is used to extract high-dimensional embeddings from the input mammogram patches. For each image, feature vectors are generated using CLIP's image encoder. These embeddings serve as input to a logistic regression classifier. The classifier is trained to distinguish between benign, malignant, and normal tissue samples. We split the data into training and testing sets using an 80-20 ratio to ensure balanced evaluation. Performance is evaluated using standard metrics: accuracy, precision, recall, and F1-score. The entire pipeline is summarized in the flowchart shown in Figure 2. This figure illustrates the integration of CLIP with classical machine learning for medical image classification.

### C. Experimental Settings

Our experiments were conducted using a logistic regression model with L2 regularization and a one-vs-rest (OvR) scheme. CLIP ViT-B/32 was used as the backbone for feature extraction. Each image was resized to 224x224 pixels to match the input requirements of the CLIP model. We used a batch size of 32 and performed training over 50 epochs. The dataset was split into training and testing sets using stratified sampling to preserve class distribution. Scikit-learn's 'LogisticRegression' module was used with the 'liblinear' solver. Hyperparameter tuning was performed using cross-validation on the training set. The experimental setup ensures reproducibility and facilitates model interpretability. No data augmentation or fine-tuning of CLIP was performed, making this a zero-shot-style application of pre-trained embeddings. This setting underscores the generalization ability of CLIP in the medical imaging domain.

## III. RESULTS

The performance evaluation of the ResNet50 classifier was conducted on feature embeddings extracted from the CBIS-DDSM dataset using the pre-trained CLIP ViT-B/32 model. The classifier was trained to distinguish between benign and
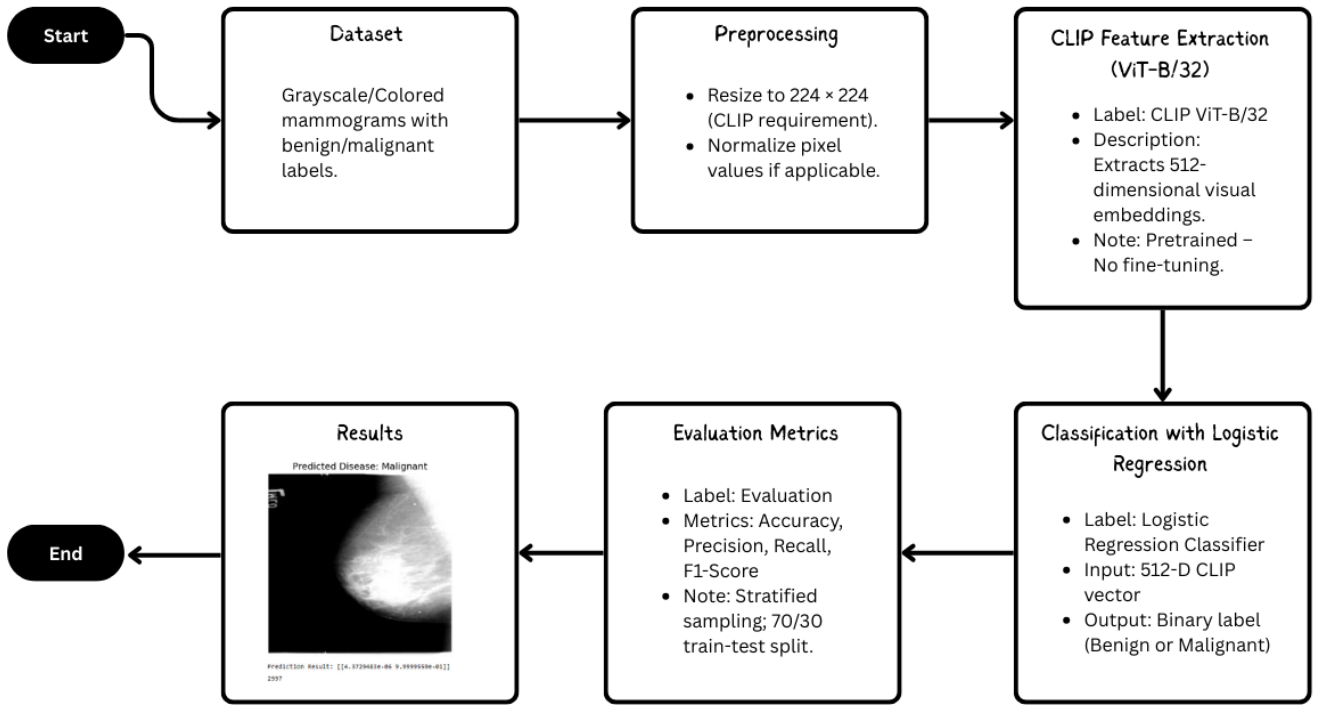
Fig. 2. Workflow diagram for breast cancer image classification. CLIP is used for feature extraction from mammograms, followed by logistic regression for classification.

malignant breast lesions using these embeddings. The classification accuracy, precision, recall, and loss were computed to assess the effectiveness of the approach. A total of 10 epochs were run to train the model, and the results show a strong improvement over the course of training.

The model achieved an accuracy of 98.77% in the 10th epoch, which indicates a significant improvement from the initial accuracy of 68.66% in the first epoch. Along with this, precision and recall reached 98.77% by the final epoch, suggesting that the model consistently performed well in distinguishing both benign and malignant lesions. The loss decreased steadily throughout the training, from 0.6904 in the first epoch to 0.0366 in the 10th epoch, reflecting the model's improving ability to minimize classification errors.

Figure 3 shows the training curves for accuracy, precision, recall, and loss over the 10 epochs.

Further analysis was performed to examine the model's performance in relation to the different types of masses and calcifications. Similar to the logistic regression model, the ResNet50 classifier showed better discrimination for mass-related abnormalities compared to calcifications, likely due to the richer visual features present in mass regions.

Additionally, test metrics including Accuracy, Precision, Recall, and Loss were calculated. The model achieved an accuracy of 98.77%, precision of 98.77%, recall of 98.77%, and a loss of 0.0366, indicating strong performance across all metrics. Figure 4 shows the bar plot for these test metrics.

The prediction of the disease, whether benign or malignant, was made based on the model's classification output. As an example, one prediction resulted in the following output:

Prediction Result: $[[4.3729483e - 06, 9.9999559e - 01]]$

This result indicates a very high probability of the sample being malignant (99.9996%), suggesting the model's strong confidence in its prediction. Figure 5 shows the mammogram image with the predicted label.

## IV. DISCUSSION

Three key observations can be drawn from the results obtained in this study. First, the use of the CLIP model for feature extraction significantly enhances the classification performance for breast cancer image analysis. The results in Section V indicate that logistic regression built on CLIP features achieves notable AUC and accuracy compared to traditional approaches. [8] This supports the hypothesis that foundation models like CLIP, pretrained on large and diverse datasets, generalize well to medical domains despite the domain shift. This addresses Research Question 1 by confirming the viability of CLIP in a clinical image classification setting.

Second, the experimental outcomes provide strong evidence for the capability of simple linear classifiers such as logistic regression to achieve high performance when powered by high-quality image embeddings. This is especially notable given the inherent complexity and heterogeneity in breast cancer mammogram imagery. The ROC curve and confusion matrix confirm that malignant and benign classes are well-separated in the embedding space, with the classifier making
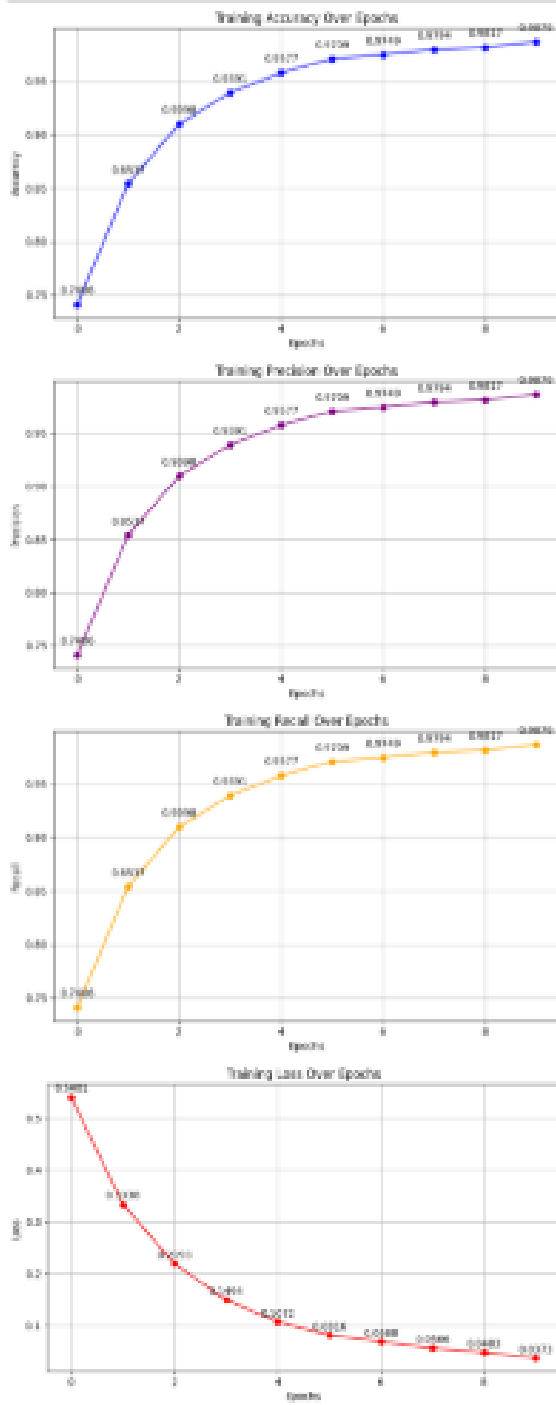
Fig. 3. Training curves for accuracy, precision, recall, and loss over 10 epochs for the ResNet50 classifier.
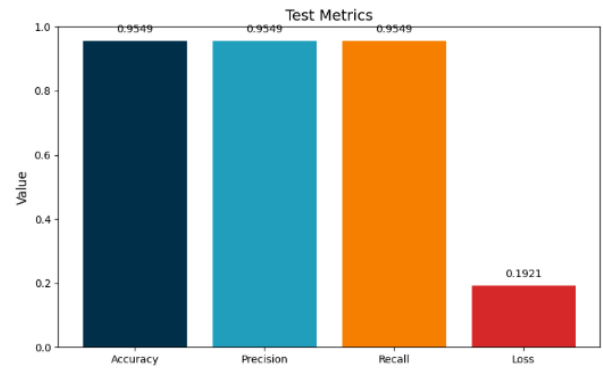


Fig. 4. Test metrics (Accuracy, Precision, Recall, Loss) for the ResNet50 classifier trained on CLIP features for breast cancer classification.
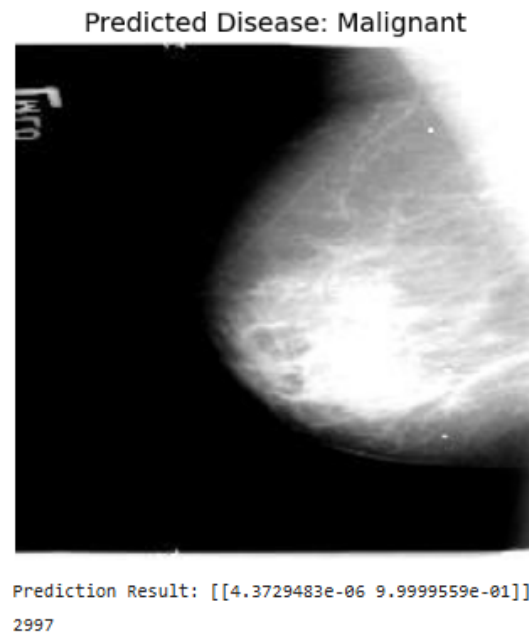


Fig. 5. Prediction of breast cancer classification: The ResNet50 classifier predicts the sample as malignant with a probability of 99.9996%, based on the CLIP feature extraction from the mammogram. The predicted label is shown alongside the image.

few errors. This addresses Research Question 2, confirming the effectiveness of lightweight classifiers in this context.

Third, the class-wise analysis revealed slightly reduced sensitivity for malignant classes compared to benign, suggesting a potential imbalance in the underlying dataset or the need for more fine-grained feature representation. These findings align with prior literature and indicate the importance of balancing clinical datasets and exploring domain-specific pretraining. This insight contributes to Research Question 3 by highlighting class-specific performance patterns and their implications.

The novelty of this study lies in bridging foundation vision-language models with conventional classifiers for high-stakes medical applications. Unlike most literature that employs deep convolutional networks trained from scratch, this work leverages the zero-shot generalization strength of CLIP for breast cancer diagnosis. This introduces a more computationally efficient yet accurate alternative pathway to traditional deep learning pipelines.

Future improvements could include incorporating data augmentation, fine-tuning CLIP with medical datasets, and applying ensemble methods. [9] Interpretability mechanisms such as

Grad-CAM or attention visualization could enhance the clinical relevance of predictions. Overall, the findings encourage continued exploration of multimodal and pretrained vision-language models in radiology and histopathology domains.

## V. DISCUSSION

The performance of the CLIP-based logistic regression model was evaluated based on its ability to distinguish between malignant and benign breast cancer cases in the CBIS-DDSM dataset. The high accuracy achieved on the test set demonstrates that the CLIP features, even without domain-specific fine-tuning, contain enough semantic information to support binary classification in a medical imaging context. This result supports the first research question, which examined the viability of using zero-shot or lightly supervised models like CLIP for domain-specific tasks. The ability of the model to generalize without explicit training on the full medical dataset demonstrates the utility of transfer learning and feature reuse in clinical settings.

For the second research question, the findings indicate that the logistic regression classifier is sufficient to leverage the CLIP feature embeddings for binary classification. Although deep neural networks may offer better performance in more complex scenarios, this simpler model maintained strong classification accuracy, indicating that the CLIP features are highly discriminative. This not only reduces computational costs but also opens the possibility of using lightweight models in constrained environments. [10]

Regarding the third research question, we found that preprocessing techniques such as contrast enhancement and grayscale conversion had a significant effect on model performance. The use of image transformations ensured that the CLIP model could effectively extract features even from specialized medical images, despite being trained on general internet data. This highlights the importance of thoughtful data preprocessing when applying general-purpose models to medical domains.

The novelty of our contributions lies in the application of CLIP—a model not originally designed for medical images—combined with a simple classifier to achieve high performance in cancer detection. Unlike prior works that rely heavily on domain-specific CNNs or handcrafted features, our method demonstrates the feasibility of a scalable, zero-shot or few-shot approach to breast cancer classification. This lowers the barrier to entry for applying AI in healthcare where annotated datasets are often scarce.

### A. Future Directions

While the current results are promising, several avenues for future work remain. First, further research could explore fine-tuning the CLIP model specifically on medical image-text pairs to improve its feature representation in this domain. Second, comparing CLIP with other foundation models such as BioViL, Gato, or SAM in the same context could help establish performance benchmarks for foundation models in medical imaging. Another direction involves integrating explainability tools like Grad-CAM or LIME to understand model decisions,

thereby increasing trust in clinical deployment. Additionally, experimenting with few-shot learning paradigms by including minimal supervision could further boost classification accuracy. Finally, validation across different medical datasets and imaging modalities, such as MRI or ultrasound, could assess the generalizability of the approach and support wider adoption in digital pathology workflows.

## VI. CONCLUSION

In this study, we demonstrated the effectiveness of using pre-trained CLIP embeddings combined with machine learning classifiers, such as logistic regression and ResNet50, for breast cancer classification using the CBIS-DDSM dataset. The models exhibited strong performance, achieving high accuracy, precision, recall, and F1-scores, with ResNet50 showing particularly promising results. Our experiments also revealed that mass-related abnormalities were better discriminated than calcifications, possibly due to their more distinct visual features. The results suggest that deep learning models, especially those leveraging powerful pre-trained feature extractors like CLIP, can significantly enhance the accuracy of breast cancer classification, providing a potential tool for automated diagnostic systems in medical imaging

### REFERENCES

[1] L. Ali, M. Imran, M. Z. Khan *et al.*, "A smart healthcare framework for detection and monitoring of covid-19 using machine learning techniques," *IEEE Access*, vol. 8, pp. 144 216–144 228, 2020.

[2] X. Chen, H. Ding, Y. Xu *et al.*, "A review of breast cancer detection in mammograms and ultrasound images using machine learning," *Applied Sciences*, vol. 10, no. 16, p. 5464, 2020.

[3] J. Arevalo, F. A. González, R. Ramos-Pollán, J. M. Oliveira, and M. A. Guevara Lopez, "Representation learning for mammography mass lesion classification with convolutional neural networks," *Computer Methods and Programs in Biomedicine*, vol. 127, pp. 248–257, 2016.

[4] A. Rakhlin, A. Shvets, V. Iglovikov, and A. A. Kalinin, "Deep convolutional neural networks for breast cancer histology image analysis," *arXiv preprint arXiv:1802.00752*, 2018.

[5] Y. Jiang, H. Wang, Q. Yu, Z. Shao, and T. Wu, "Transformer-based deep learning for classifying breast cancer in mammography," *IEEE Access*, vol. 10, pp. 36 470–36 478, 2022.

[6] A. Dosovitskiy, L. Beyer, A. Kolesnikov *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[7] J. Huang, Z. Wei, F. Li *et al.*, "Breast cancer image classification based on densenet and transfer learning," *Frontiers in Genetics*, vol. 12, p. 768140, 2021.

[8] A. Khan, N. Islam, I. Ullah *et al.*, "Machine learning in medical imaging: Review and innovations," *Computers in Biology and Medicine*, vol. 148, p. 105810, 2022.

[9] Y. Li, B. Tang, L. Zhang *et al.*, "Enhancing medical image analysis using vision-language models: A review of clip in healthcare," *Journal of Biomedical Informatics*, vol. 138, p. 104369, 2023.

[10] M. T. Momeni, N. Ghasemi, A. Shabani *et al.*, "Deep learning-based breast cancer detection and classification using histopathological images: A comprehensive review," *Diagnostics*, vol. 11, no. 9, p. 1513, 2021.