

Advancements in Breast Cancer Detection and Treatment: Targeted Therapies, Imaging Innovations, and Deep Learning Analysis

1st Aakriti Kumari

Computer Science and Engineering
Indira Gandhi Delhi Technical University For Women
Delhi, India
aakritikumari247@gmail.com

2nd Himanshi

Computer Science and Engineering
Indira Gandhi Delhi Technical University For Women
Delhi, India
himanshi26082004@gmail.com

3rd Ishika Chhoker

Computer Science and Engineering
Indira Gandhi Delhi Technical University For Women
Delhi, India
chhokerishika95@gmail.com

ABSTRACT

Cancer starts when healthy cells of the breast change and growth happens in them uncontrollably thereby forming mass or various sheets of cells further hence known as a tumor. A tumor can be cancerous or noncancerous, which is also called benign. A cancerous tumor is malignant, meaning it can grow and spread to various different parts of the patient's body. Breast Cancer begins when healthy cells in the breast change and growth happens in them uncontrollably. Invasive ductal carcinoma (IDC or infiltrating ductal carcinoma) is a type of breast cancer that starts in the milk ducts of the breast and moves to nearby tissue. IDC is typically diagnosed after detection through routine mammogram. While IDC does not always produce typical breast cancer symptoms, we might notice a new lump or unusual changes in the breast. If the physician suspects breast cancer, additional tests and a breast biopsy may confirm a diagnosis of IDC. The size of the initial tumor and the extent of cancer spread determines the stage of breast cancer. This paper describes deep learning model i.e. Convolution neural network (CNN), for detecting and visualising malignant regions of cells in the breast. The accuracy of our project has come out to be 93.13 percent. Our major aim of this sophisticated project is to encourage the research aspect in the medical field and technology development by using the concepts of deep learning. We have used a large and varying data set for the aspect of our research which consists 90,000 breast histopathology images (64,583 IDC -ve and 25,417 IDC +ve). Breast cancer is a major topic right now, and it must not be overlooked at any cost.

1. INTRODUCTION

Invasive Ductal Carcinoma of the Breast, the most prevalent kind of breast cancer, remains a huge medical problem, driving

researchers worldwide to explore creative ways for early diagnosis and successful treatment. Scientists at the University of Liverpool have developed a promising biomedical compound that could potentially halt the spread of breast cancer by targeting proteins involved in metastasis. [9] Through their research, they identified a novel compound capable of blocking the interaction between the metastasis-inducing protein S100A4 and its cellular target. This compound was further enhanced by connecting it to a molecule that triggers the cell's natural protein degradation process. [10]

Breast Magnetic Resonance Imaging (MRI) has gained prominence as a modern technique for detecting breast cancer with improved contrast. It plays a vital role in identifying additional tumor sites and assessing disease extent, especially beneficial for invasive lobular carcinoma cases that are challenging to detect through mammography. [7] While Breast MRI is successful in uncovering the primary tumor in 70% to 86% of hidden breast cancer cases, its ability to detect residual cancer post-surgery is limited due to surgical changes. Mammograms are life-saving tools for early breast cancer detection, traditionally recommended annually for individuals with an average risk, starting at age 40. [5] However, in 2023, the U.S. Preventive Services Task Force (USPSTF) proposed revised recommendations suggesting biennial mammograms starting at age 40 for those at average risk. Three-dimensional (3-D) mammography, also known as digital breast tomosynthesis (DBT), is a variant of digital mammography. It captures thin breast "slices" from various angles and reconstructs images with computer software, similar to CT scanning [6]. While 3-D mammography employs low-dose x-rays, combining it with standard two-dimensional (2-D) digital mammography increases radiation exposure. Recent tomosynthesis advancements offer standalone DBT, potentially reducing radiation exposure closer to standard mammography levels. [14] Screening

aids early cancer detection, sparing many women (nearly 98%) from breast cancer diagnoses if mammograms and follow-up exams reveal no abnormalities.

Doctors employ a specific classification system called TNM (Tumor, Node, Metastasis) to determine the various stages of breast cancer. The "T" in TNM represents the tumor size, measured in millimeters or centimeters. Breast cancer can be categorized into stages 0 through 4, determined not only by tumor size but also by several other factors. [12] Breast cancer staging extends beyond the TNM system to consider additional factors like hormone markers, given the distinct characteristics of breast tissue. A doctor can provide explanations for assigning a particular stage to a cancer diagnosis. [3]

Invasive ductal carcinoma, or IDC, is the predominant form of breast cancer, constituting 80% of diagnoses. Early detection leads to a favorable prognosis, boasting a 5-year survival rate of nearly 100% for localized cases. Surgical approaches vary based on tumor characteristics, including location, size, stage, and patient health, with options ranging from lumpectomy to mastectomy. [8]

A CNN, or Convolutional Neural Network, is a deep learning architecture built primarily for image and visual data processing applications. CNNs have shown to be extremely effective in a variety of computer vision applications, including picture classification, object identification, image segmentation, and others. Their capacity to acquire hierarchical features from raw data automatically makes them a critical tool in the fields of deep learning and computer vision. [4]

2. METHODOLOGY

Artificial intelligence is a broad term that is used to describe the system and the machines that resembles human intelligence. Machine learning is a subdomain of Artificial Intelligence which focuses on the development of the system and dedicated towards improving the performance of those system by acquiring knowledge from the input data.

Deep Learning is a branch of machine learning which uses artificial neural networks(ANN). It possesses the ability of grasping complex patterns and correlations within data. [1] It's popularity surged in recent years due to the advancement in processing power and availability of large datasets because it primarily based on Artificial Neural network also referred as Deep neural network. These neural network are based on the structure and functioning of neurons in human brain. DNN requires huge amount of data to learn. In DNN there is an input layer and more than one hidden layer. Every previous layer provides an input for next layer of neurons means output of one neuron becomes input for other and this continues till final layer is reached and produce an output. This network learns complex representations of input data through non-linear transformations of input data.

CNN are similar to traditional ANNs in the way they both consists of neurons that refine there performance through learning. On the basis of countless ANNs every neuron will receive an input and perform a operation(scalar product). The only difference between CNNs and traditional ANNs is

that CNNs is used primarily for the recognition of pattern within images. One of the disadvantage of traditional ANNs is difficulty in handling computational complexity associated with processing image data. [13]

2.1 CNN Architecture

One of the neural networks CNN(Convolutional Neural networks) have been substantiated to be very useful in the diagnosis and categorization of breast cancer. Named ConvNets due to the hidden layers it includes and suggest inscriptions. Convolution Neural Networks comprises of convolution , pooling and fully linked layer [2]. In CNN the convolution layer groups the properties of the original picture , pooling reduces the dimensions of the picture and fully linked layer returns the desired output. CNN is very useful as brighter pixels can be used to represent the picture borders and lay stress on the attributes of the local picture for more processing. Every layer of CNN includes an activation function , also convolution and pooling function. It takes an image and one filter as an input and gives the required images an output. For example, all the parameters that have been used in the data like image size($128 \times 128 \times 3$), image height ,width etc affects the performance of neural networks. The image channel have RGB image channel(3 represents RGB) and the height of 128×128 pixels so processing size becomes 49,152 bytes [11]. CNN is a typical form of neural network which may overfit very small images, waste memory and need comprehensive image processing.

$$C = \sum_{i=1}^i \sum_{j=1}^j I_{ij} F_{ij} \quad (1)$$

In equation 1 ,F represents the filter or convolution kernel and i,j represents rows and columns. Figure 1 shows the input image, filter and the new 2 dimensional output in which the image is multiplied with the kernal or filter. Convolution splits images into single layer neural network known as perceptrons which are subsequently flattened on (y)X axis(z). Every layer has x filters for the location of characteristics. X-sized feature maps are generated by layer L which are illustrated as follows:

$$C_i^L = B_i^L + \sum_{j=1}^{x^{(L-1)}} F_{i,j}^L * C_j^{(L-1)} \quad (2)$$

Where B_i^L and $F_{i,j}^L$ represent bias matrix and the filter that connects the j^{th} feature map in the layer respectively.

Suppose that blue channel is at -1, the other channels might have values of +1 or 0. One should use the dot product method , to determine the convolution value. By using convolution the images get deformed. The picture size is same (F_s) if I_s and I_s where $I_s=1$ in this example.

$$C_s = ((I_s - F_s)/S + 1) \quad (3)$$

If I_s , F_s , S have values 6, 3, 1 respectively the C_s will be equal to $4(6-3+1) = 4$. The equation illustrates the method used to calculate the output dimension:

$$[(Width - F_s + 2P)/S] + 1 \quad (4)$$

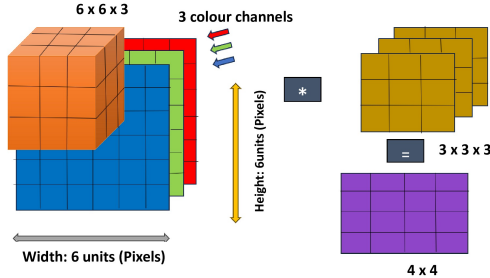


Fig. 1. Convolutional Operation

Fs represents filter size, P represents padding, S represents stride, and equation (4) represents the floor value.

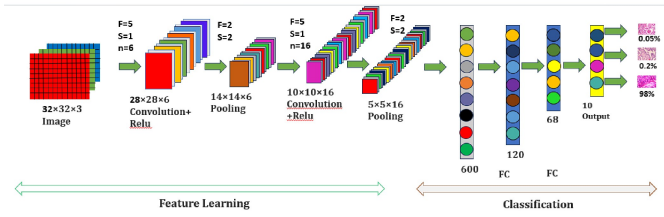


Fig. 2. Convolutional Neural Network

Figure 2 shows that how convolution method applied on every image before pooling. F.C. layer uses the activation function to classify the cancer images. There are various named architecture in CNN that are used for the detection and classification of breast cancer. At basic level Convolutional neural network consists of series of layers which process the input data and produce a classification score. In building a CNN model , convolution layer, pooling , Relu and fully connected layers are used [11].

3. EXPERIMENT ANALYSIS

3.1 Dataset Description

The dataset used in our study comprises 162 whole mount slide images of breast cancer (BCa) specimens that were scanned at 40x magnification. From these images, we extracted a total of 277,524 patches, each measuring 50 x 50 pixels. These patches consist of 198,738 samples that are IDC negative and 78,786 samples that are IDC positive. The filenames for each patch follow a specific format, such as 10253_idx5_x1351_y1101_class0.png where "u" represents the patient ID (e.g., 10253_idx5), "X" denotes the x-coordinate of the patch's cropping location, "Y" represents the y-coordinate of the cropping location, and "C" indicates the class, with 0 indicating non-IDC and 1 indicating IDC.

3.2 Data Preprocessing

3.2.1 Data Collection and Classification: We applied the glob function to collect PNG image file paths from sub-directories of the /content/IDC_regular_ps50_idx5/ directory.

This is a recursive procedure that includes photos from every subdirectory. Images are divided into two categories: '0' for IDC negative images and '1' for IDC positive images. Based on their filenames, we divided the images into two lists: non_img and can_img. We visualised a random selection of photos from both the 'IDC -ve' and 'IDC +ve' categories in a 6x6 grid using TensorFlow and Keras.(figure 3)

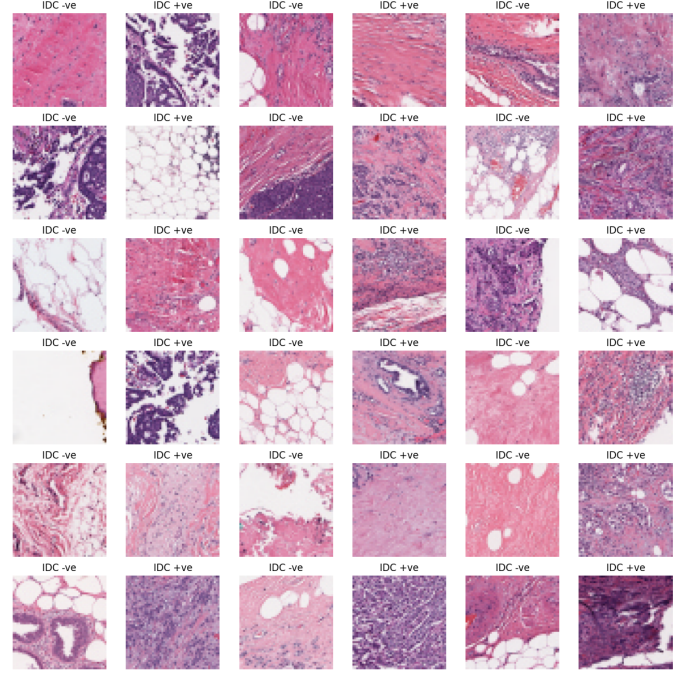


Fig. 3. Classification of Breast Cancer Images

3.2.2 Image Visualization: In the realm of visualizing breast cancer images, a comprehensive approach is embraced. This entails the creation of two distinct subplots, each structured as a 5x10 grid, dedicated to showcasing images categorized as positive and negative cases. Through the iteration over selected indices representing these samples, the corresponding images are retrieved and depicted using the 'imshow' function. This phase, which relies more on visual presentation rather than intricate code, provides a tangible portrayal of a subset of the dataset, effectively highlighting both positive and negative instances. Furthermore, the data preprocessing pipeline encompasses the formulation of essential coordinate extraction functions, specifically '*extract_coords*' and '*get_cancer_dataframe*.' These functions play a pivotal role in the extraction and refinement of coordinate data embedded within image filenames, facilitating subsequent in-depth analyses. Additionally, within the pipeline, advanced visualization functions like '*visualise_breast_tissue_base*' and '*visualise_breast_tissue*' are crafted to enable detailed examination of breast tissue and the identification of cancerous regions within the images(figure 3, figure 4) Lastly, for the visualization of binary classification, the '*visualise_breast_tissue_binary*' function is employed to generate scatter plots of breast tissue images.(figure 5) Within

these plots, data points are color-coded based on binary classification labels (0 or 1), offering valuable insights into the distribution of cancerous and non-cancerous regions. Collectively, these visualization and coordinate extraction elements enrich the comprehension and exploration of the breast cancer image dataset, establishing a robust foundation for subsequent analyses and model development.

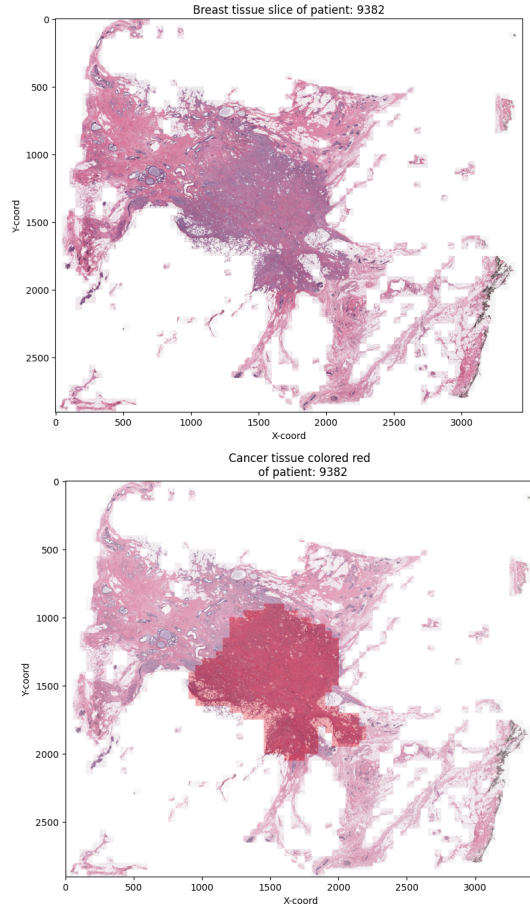


Fig. 4. Identification of the cancerous region (red color region)

3.2.3 Dataset Splitting: The dataset is divided into two distinct sets: training and testing, as a crucial step in the data preprocessing process. This division serves the purpose of setting aside a specific portion of the data for later model assessment. Within this context, X encapsulates the image data, whereas y contains the associated labels, representing either 0 or 1. Notably, the `test_size` parameter, configured at 0.3, designates that 30% of the dataset is allocated for testing, while the remaining 70% forms the training dataset. This strategic partitioning ensures that the model can be rigorously evaluated while being trained on a substantial portion of the data.

3.2.4 Model Selection: The architecture and configuration of the machine learning or deep learning model that will

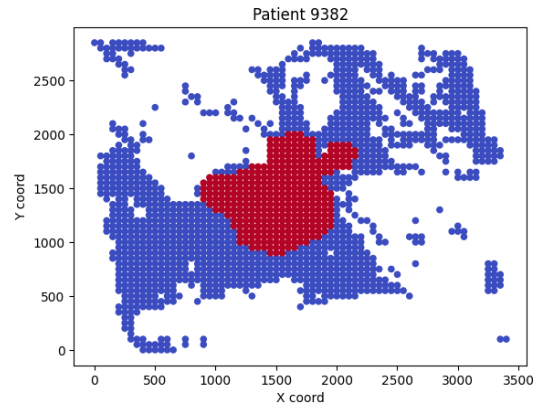


Fig. 5. Scatter plot: Identification of the cancerous region

be used for a particular job are chosen through the model selection process. The Keras Sequential API is used to specify the model architecture. In this, the type and number of layers, activation mechanisms, and other model hyperparameters are specified. For the goal of classifying breast cancer, a convolutional neural network (CNN) architecture is chosen. The specific layers, such as convolutional layers, pooling layers, and dense layers, are defined as part of the model selection process.

3.2.5 Model Training Process: The model is trained across 20 epochs with a batch size of 50. The training data is used to update the model's internal parameters, allowing it to learn how to categorise breast cancer images. The history object stores information about the training process, including accuracy and loss values for each epoch.

3.2.6 Result Visualization: Figure 6 shows the accuracy of training and validation over epochs. Figure 7 depicts the loss in training and validation throughout epochs. Figure 8 i.e. a bar plot to show the number of true and false predictions. It displays the number of true predictions (`true_predict`) and the number of false predictions (`false_predict`). For the single image prediction we had defined a function called `predict_single_image(idx)` that made predictions on a single image specified by its index (`idx`) in the dataset. We loaded and preprocess the image, and then used the trained model to predict its labelled. (figure 9)

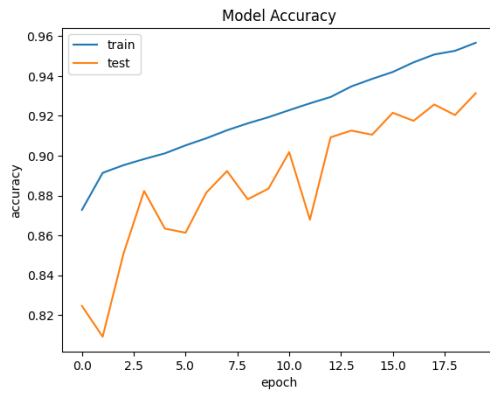


Fig. 6. Accuracy vs Epoch

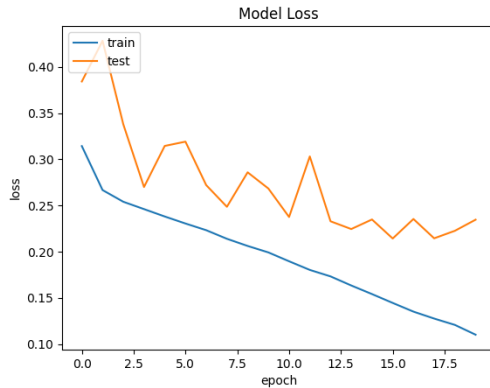


Fig. 7. Loss vs Epoch

4. CONCLUSION

Our study on breast histopathology images has yielded valuable insights and outcomes that contribute to the ongoing battle against this devastating disease. Through a combination of cutting-edge technology, data analysis, and collaborative efforts, we have made significant strides in the early detection and diagnosis of breast cancer. The accuracy of our model has come out to be 93.13%. The application of the CNN technique to the analysis of a wide range of data, such as malignant regions of breast cell, is a key outcome of this study. Our research also contributes to the advancement targeted therapies by enabling early and accurate diagnosis, identifying biomarkers, and tracking therapy responses. This model enable healthcare providers to create highly personalised treatment regimens and choose the most efficient treatments for individual patients based on their unique cancer profiles. This accuracy increases not just treatment efficacy, but also minimises unnecessary therapies and improves overall patient outcomes in the battle against breast cancer. cancer remains a leading cause of mortality across nations in contemporary times, underscoring the urgent need for continued research, prevention, and improved treatments to combat this formidable global health challenge. As a result, we hope to make a significant contribution to this continuing scientific endeavour.

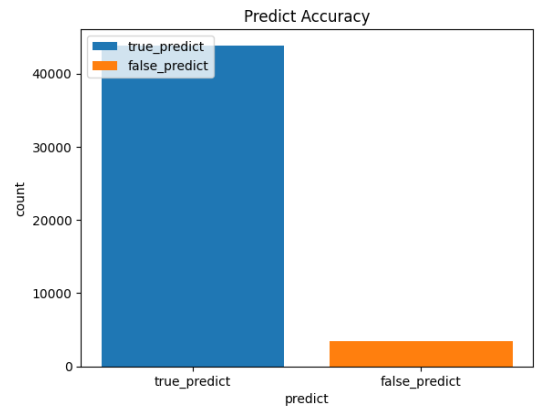


Fig. 8. True and False prediction

```
Cancer : 1 No Cancer : 0
True Label : 0
1/1 [=====] - 0s 29ms/step
Predicted Label : 0
```

Fig. 9. Single Image Predictions

DATA AVAILABILITY STATEMENT

<https://www.kaggle.Com/paultimothymooney/breast-histopathology-images>.(accessed on 27 July 2023).

REFERENCES

- [1] Theyazn HH Aldhyani, Rajit Nair, Elham Alzain, Hasan Alkahtani, and Deepika Koundal. Deep learning model for the detection of real time breast cancer images using improved dilation-based method. *Diagnostics*, 12(10):2505, 2022.
- [2] Theyazn HH Aldhyani, Rajit Nair, Elham Alzain, Hasan Alkahtani, and Deepika Koundal. Deep learning model for the detection of real time breast cancer images using improved dilation-based method. *Diagnostics*, 12(10):2505, 2022.
- [3] Gábor Cserni, Ewa Chmielik, Bálint Cserni, and Tibor Tot. The new tnm-based staging of breast cancer. *Virchows Archiv*, 472:697–703, 2018.
- [4] Eytan Gilboa. The cnn effect: The search for a communication theory of international relations. *Political communication*, 22(1):27–44, 2005.
- [5] Peter C Gøtzsche and Karsten Juhl Jørgensen. Screening for breast cancer with mammography. *Cochrane database of systematic reviews*, (6), 2013.
- [6] Mark A Helvie. Digital mammography imaging: breast tomosynthesis and advanced applications. *Radiologic Clinics*, 48(5):917–929, 2010.
- [7] Margaret Houser, David Barreto, Anita Mehta, and Rachel F Brem. Current and future directions of breast mri. *Journal of Clinical Medicine*, 10(23):5668, 2021.
- [8] Sylvie Ménard, Stefania Fortis, Fabio Castiglioni, Roberto Agresti, and Andrea Balsari. Her2 as a prognostic factor in breast cancer. *Oncology*, 61(Suppl. 2):67–72, 2001.
- [9] Elwin A Morgan, Shiva S Forootan, Janet Adamson, Christopher S Foster, Hiroshi Fujii, Michihiro Igarashi, Carol Beesley, Paul H Smith, and Youqiang Ke. Expression of cutaneous fatty acid-binding protein (c-fabp) in prostate cancer: potential prognostic marker and target for tumourigenicity-suppression. *International journal of oncology*, 32(4):767–775, 2008.
- [10] Carlo Palmieri, Alison Musson, Catherine Harper-Wynne, Duncan Wheatley, Gianfilippo Bertelli, Iain R Macpherson, Mark Nathan, Ellie McDowall, Ajay Bhojwani, Mark Verrill, et al. A real-world study of the first use of palbociclib for the treatment of advanced breast cancer within the uk national health service as part of the novel ibrance® patient program. *British journal of cancer*, pages 1–9, 2023.

- [11] Faezehsadat Shahidi, Salwani Mohd Daud, Hafiza Abas, Noor Azurati Ahmad, and Nurazeen Maarop. Breast cancer classification using deep learning approaches and histopathology image: a comparison study. *IEEE Access*, 8:187531–187552, 2020.
- [12] S Eva Singletary and James L Connolly. Breast cancer staging: working with the sixth edition of the ajcc cancer staging manual. *CA: a cancer journal for clinicians*, 56(1):37–47, 2006.
- [13] Mohammad Mustafa Taye. Understanding of machine learning with deep learning: Architectures, workflow, applications and future directions. *Computers*, 12(5):91, 2023.
- [14] Srinivasan Vedantham, Andrew Karellas, Gopal R Vijayaraghavan, and Daniel B Kopans. Digital breast tomosynthesis: state of the art. *Radiology*, 277(3):663–684, 2015.