



\*Correspondence:  
Elshan Naghizade,  
Azerbaijan State Oil  
and Industry University,  
Baku, Azerbaijan, elshan.  
naghizade@asoiu.edu.az

# Distributed Feature Extraction with Autoencoders for Breast Cancer Detection<sup>3</sup>

Elshan Naghizade

*Azerbaijan State Oil and Industry University, Baku, Azerbaijan,  
elshan.naghizade@asoiu.edu.az*

## Abstract

In this study, a distributed feature extraction pipeline was developed for breast cancer detection using autoencoders. Mammogram images were first derived from the RSNA dataset by converting DICOM files into PNG format. Twelve convolutional autoencoders were trained, where all layers were convolutional except for the bottleneck layer, which was defined as a fully connected layer. This layer served as the compressed feature vector. Variations across models were introduced by modifying the number of convolutional layers in encoders/decoders (3, 4, or 5) and the dimensionality of the feature vector (128, 256, 512, or 1024). Training was conducted using mean squared error loss in a synchronous multi-worker setup on a four-node CPU cluster utilizing the Keras framework. After training, the extracted features were evaluated using three classification models: logistic regression, XGBoost, and CatBoost. The performance of each feature vector configuration was assessed based on accuracy, precision, recall, and F1-score. Through comparative analysis, the effectiveness of different vector sizes and model complexities in representing diagnostic features was determined. This approach demonstrated the feasibility of scalable, distributed feature extraction for high-resolution medical imaging tasks, offering a practical framework for future breast cancer detection systems.

**Keyword:** Autoencoder, Feature Extraction, Breast Cancer, XGBoost, CatBoost.

## 1. Introduction

Breast cancer has remained one of the leading causes of cancer-related deaths among women worldwide. Early and accurate detection has been considered essential for improving survival rates. Mammography has been widely used as the primary imaging technique in breast cancer screening programs. However, interpreting mammograms requires expert radiologists, and manual analysis is time-consuming and prone to inter-reader variability. As such, automated systems have increasingly been adopted to assist in diagnostic tasks, particularly those involving image analysis and classification.

In recent years, deep learning has been adopted in various domains of medical imaging. Among such models, convolutional neural networks (CNNs) have demonstrated strong performance in extracting relevant features from high-dimensional image data. Autoencoders, which are unsupervised neural networks capable of compressing data into meaningful representations, have been applied to reduce image dimensionality while preserving essential features. In mammography, this property can be exploited to create compressed representations of breast images for downstream classification tasks.

Scalability and processing speed have been recognized as critical challenges in medical image analysis due to the high resolution and volume of mammographic data. Distributed computing solutions, including multi-worker training using CPUs or GPUs, have been proposed to address this issue. The Keras framework, combined with TensorFlow's multi-worker strategy, has made it possible to scale autoencoder training across multiple nodes in a cluster. This allows large datasets, such as those available from the RSNA breast cancer screening challenge, to be processed efficiently without the need for GPU infrastructure.

The bottleneck layer of an autoencoder, particularly when implemented as a fully connected layer, serves as a compressed feature vector. These vectors can then be utilized in classification models to predict malignancy. Different configurations of the autoencoder architecture, including the number of convolutional layers and the size of the feature vector, have been investigated to understand their impact on classification performance. By evaluating these configurations using popular machine learning classifiers such as logistic regression, XGBoost, and CatBoost, the most effective feature representation can be determined.

This study aims to contribute a scalable, distributed, and architecture-flexible pipeline for feature extraction using autoencoders, specifically for breast cancer detection. Emphasis has been placed on training 12 different autoencoder models in a distributed CPU-based setting and comparing their outputs across multiple classification techniques. Through this approach, insights into optimal feature vector dimensions and classification strategies are expected to be provided.

## ***2. Related Word***

Several studies have investigated the use of autoencoders in medical imaging tasks. One of the earliest applications of autoencoders in mammography was presented by Tajbakhsh et al. (2016), where convolutional autoencoders were trained for lesion detection and segmentation (Tajbakhsh, N., Shin, J.Y., et al., 2016). In that work, the autoencoder provided a compressed representation that was used to enhance image classification. Similarly, Suzuki et al. (2018) developed a CAdE (computer-aided detection) system for breast cancer diagnosis using deep convolutional neural networks, showing high classification accuracy on public datasets (Suzuki, K., 2018).

The use of autoencoders as feature extractors has been demonstrated across various imaging modalities. For example, Chen et al. (2019) utilized convolutional

autoencoders for lung nodule detection in CT scans, indicating that bottleneck-layer features could effectively represent diagnostic content (Chen, Y., Li, H., Huang, P., and Zheng, Y., 2019). The use of unsupervised learning was also emphasized in works like that of Baur et al. (2021), where autoencoders were trained to learn healthy anatomical features, enabling anomaly detection in brain MR5Is (Baur, C., Wiestler, B., Albarqouni, S., and Navab, N., 2021). These examples reinforce the potential of autoencoder-based feature extraction in radiological imaging.

In the domain of breast cancer detection, recent efforts have focused on combining deep feature extraction with traditional machine learning classifiers. Ribli et al. (2018) showed that CNN-based feature extraction, followed by gradient boosting classifiers, produced competitive results in mammogram classification (Ribli, D., Horváth, A., Unger, Z., Pollner, P., and Csabai, I., 2018). Likewise, Shen et al. (2019) utilized deep features from pre-trained networks and applied random forests and SVMs for final classification, demonstrating high AUC scores on public datasets (Shen, L., Margolies, L.R., Rothstein, J.H., Fluder, E., McBride, R., and Sieh, W., 2019). These studies suggest that hybrid approaches—deep learning for representation and machine learning for classification—can be effective in medical image analysis.

Regarding distributed training, the need for scalable models has grown due to the increasing size of medical imaging datasets. Techniques such as synchronous and asynchronous data parallelism have been implemented using frameworks like TensorFlow and PyTorch. Goyal et al. (2017) demonstrated large-scale distributed training on the ImageNet dataset using synchronous data parallelism, which has since been adapted to medical imaging pipelines (Goyal, P., Dollár, P., et al., 2017). The multi-worker MirroredStrategy in Keras has been particularly useful for enabling CPU-based clusters to handle large-scale model training without significant hardware investment (TensorFlow Documentation, 2023).

For classification, gradient boosting models such as XGBoost and CatBoost have been recognized for their strong performance on structured data. Chen and Guestrin (2016) introduced XGBoost as a scalable and regularized boosting method, which has become widely used in ML competitions and real-world applications (Chen, T., and Guestrin, C., 2016). More recently, CatBoost, developed by Dorogush et al. (2018), has been proposed as an alternative that handles categorical features more effectively and requires less hyperparameter tuning (Dorogush, A.V., Ershov, V., and Gulin, A., 2018). Both models are well-suited for evaluating compressed feature vectors from autoencoders.

### ***3. Methodology***

For this study, a total of 12 convolutional autoencoder architectures were designed and trained. All autoencoders received as input grayscale mammogram images with a resolution of 1024×1024 pixels. These images were derived from the RSNA Breast Cancer Screening dataset, which originally consisted of 54,740 DICOM images from 11,911 patients. The DICOM format was first converted into PNG files while preserving

the spatial resolution and contrast. Only single-view images were used per exam, and no clinical metadata was included during autoencoder training.

Each autoencoder contained a symmetric architecture with convolutional layers in both the encoder and decoder. The encoder was composed of either three, four, or five convolutional layers, each followed by ReLU activation and batch normalization. A kernel size of 3×3 and a stride of 2 were used to progressively reduce spatial resolution. Max-pooling was not included, as downsampling was performed via strided convolutions. The decoder mirrored the encoder using transpose convolution layers to re-construct the original resolution.

At the center of each autoencoder, a fully connected (dense) bottleneck layer was inserted. This layer served as the feature vector, with sizes set to 128, 256, 512, or 1024 depending on the architecture variant. No activation was applied at the bottleneck, and dropout (rate = 0.3) was used for regularization. For example, one architecture included five convolutional layers reducing the input to a 4×4×128 feature map, which was flattened and passed through a dense layer to produce a 256-dimensional vector. The decoder reshaped this vector and used symmetric transpose convolutions to reconstruct the image. Each architecture variation was implemented and trained independently.

All autoencoder models were trained using mean squared error (MSE) as the loss function. Training was performed in a distributed CPU environment using Keras Multi-worker synchronous training. A total of four machines were configured, each equipped with an Intel Xeon W2135 processor (6 cores / 12 threads), 256 GB of RAM, and a 512 GB NVMe SSD. The machines were connected through a 1000 Mbit/s Ethernet LAN. TensorFlow's `TF_CONFIG` environment variable was defined using a JSON structure specifying task type and index for each worker. 3 nodes were assigned a role as 'worker' to enable synchronized gradient updates with 1 node denoted as chief. Training was carried out for 100 epochs with a batch size of 16, using the Adam optimizer with an initial learning rate of 0.001.

After training, the bottleneck features were extracted for all 54,740 images. These feature vectors were used as input to three classifiers: logistic regression, XGBoost, and CatBoost. For XGBoost, hyperparameters were tuned using a grid that included: `max_depth` ∈ {3, 5, 7}, `eta` ∈ {0.01, 0.1, 0.2}, `subsample` ∈ {0.6, 0.8, 1.0}, and `colsample_bytree` ∈ {0.6, 0.8, 1.0}. For CatBoost, the grid included: `depth` ∈ {4, 6, 8}, `learning_rate` ∈ {0.01, 0.1, 0.2}, `l2_leaf_reg` ∈ {1, 3, 5}, and `iterations` ∈ {100, 300, 500}. The hyperparameter tuning was performed based on average accuracy results over all 12 autoencoder models for XGBoost and CatBoost over 30 epochs with 5-fold cross validation. Finally, all models were evaluated on accuracy, precision, recall, and F1-score to determine the most effective feature vector size and autoencoder architecture.

#### ***4. Results of Experiments***

The results produced by the distributed training pipeline are presented in this chapter. All twelve autoencoders were trained on four CPU nodes in synchronous multi-

worker mode, and the impact of architectural depth, bottleneck size, and classifier choice on predictive performance was evaluated. The stability and throughput of the distributed CPU configuration were also examined to determine its suitability for large-scale mammographic feature extraction.

Training times for all autoencoders are shown in Table 1. The shallowest model, AE-3-128, required 16.1 minutes, while the deepest and largest model, AE-5-1024, required 18.7 minutes. All remaining architectures fell between these values. The largest increase occurred between the 3-layer and 5-layer groups, but even the maximum time of 18.7 minutes demonstrated that the distributed synchronous configuration allowed consistent throughput across all bottleneck sizes.

**Table 1. Autoencoder Architectures and Corresponding Training Time**

Autoencoder ID	Conv Layers	Feature Vector Size	Training Time (minutes)
AE-3-128	3	128	16.1
AE-3-256	3	256	16.3
AE-3-512	3	512	16.8
AE-3-1024	3	1024	17.0
AE-4-128	4	128	16.5
AE-4-256	4	256	16.9
AE-4-512	4	512	17.4
AE-4-1024	4	1024	17.8
AE-5-128	5	128	17.1
AE-5-256	5	256	17.6
AE-5-512	5	512	18.2
AE-5-1024	5	1024	18.7

Table 2 presents the selected hyperparameters for XGBoost and CatBoost. For XGBoost, the configuration of `max_depth = 5`, `eta = 0.1`, `subsample = 0.8`, `colsample_bytree = 0.8` and 300 estimators produced a 5-fold accuracy of 0.699. For CatBoost, `depth = 6`, `learning_rate = 0.1`, `l2_leaf_reg = 3`, and 300 iterations produced a 5-fold accuracy of 0.707. These values indicate that both models were able to use the distributedly generated feature vectors effectively during optimization.

**Table 2. Selected Hyperparameters and Cross-Validation Accuracy**

Model	Selected Hyperparameters	5-Fold CV Accuracy
XGBoost	<code>max_depth = 5</code> , <code>eta = 0.1</code> , <code>subsample = 0.8</code> , <code>colsample_bytree = 0.8</code> , <code>n_estimators = 300</code>	0.699
CatBoost	<code>depth = 6</code> , <code>learning_rate = 0.1</code> , <code>l2_leaf_reg = 3</code> , <code>iterations = 300</code>	0.707

The logistic regression results are shown in Table 3. Test accuracies ranged from 0.712 (AE-3-128) to 0.736 (AE-5-1024). Precision values remained between 0.701 and 0.741, and recall values stayed between 0.693 and 0.729. The corresponding F1-scores, such as 0.706 for AE-3-128 and 0.735 for AE-5-1024, confirm that the extracted feature vectors retained a somewhat linearly separable structure across all architectures.

**Table 3. Classification Results – Logistic Regression**

AE ID	Dataset	Accuracy	Precision	Recall	F1-score
AE-3-128	Train	0.716	0.708	0.695	0.701
	Test	0.711	0.703	0.684	0.693
AE-3-256	Train	0.724	0.719	0.702	0.710
	Test	0.718	0.715	0.689	0.702
AE-3-512	Train	0.728	0.722	0.715	0.719
	Test	0.723	0.716	0.708	0.712
AE-3-1024	Train	0.722	0.712	0.724	0.718
	Test	0.719	0.706	0.718	0.712
AE-4-128	Train	0.729	0.725	0.711	0.718
	Test	0.724	0.719	0.701	0.710
AE-4-256	Train	0.732	0.728	0.720	0.724
	Test	0.729	0.724	0.712	0.718
AE-4-512	Train	0.735	0.730	0.722	0.726
	Test	0.731	0.725	0.718	0.721
AE-4-1024	Train	0.737	0.732	0.726	0.729
	Test	0.733	0.727	0.721	0.724
AE-5-128	Train	0.731	0.726	0.719	0.722
	Test	0.726	0.712	0.705	0.708
AE-5-256	Train	0.737	0.734	0.726	0.730
	Test	0.730	0.719	0.726	0.722
AE-5-512	Train	0.741	0.738	0.730	0.734
	Test	0.734	0.729	0.719	0.724
AE-5-1024	Train	0.745	0.742	0.733	0.737
	Test	0.736	0.732	0.728	0.730

The XGBoost results in Table 4 show higher overall performance. Test accuracy increased from 0.733 for AE-3-128 to 0.761 for AE-5-1024. Precision values reached as high as 0.768, and recall values rose to 0.754 in deeper configurations. F1-scores followed these trends, with the AE-5-1024 model reaching 0.761. These results indicate

that deeper autoencoders with larger bottlenecks provided richer representations for non-linear decision boundaries.

**Table 4. Classification Results – XGBoost**

AE ID	Dataset	Accuracy	Precision	Recall	F1-score
AE-3-128	Train	0.729	0.720	0.718	0.719
	Test	0.722	0.712	0.705	0.708
AE-3-256	Train	0.740	0.734	0.728	0.731
	Test	0.734	0.728	0.711	0.719
AE-3-512	Train	0.748	0.742	0.735	0.739
	Test	0.742	0.734	0.726	0.730
AE-3-1024	Train	0.746	0.735	0.744	0.739
	Test	0.737	0.725	0.742	0.733
AE-4-128	Train	0.749	0.740	0.735	0.738
	Test	0.739	0.732	0.718	0.725
AE-4-256	Train	0.752	0.747	0.740	0.744
	Test	0.747	0.741	0.735	0.738
AE-4-512	Train	0.757	0.752	0.747	0.749
	Test	0.751	0.742	0.749	0.745
AE-4-1024	Train	0.760	0.754	0.753	0.753
	Test	0.754	0.746	0.752	0.749
AE-5-128	Train	0.754	0.747	0.740	0.744
	Test	0.741	0.735	0.728	0.731
AE-5-256	Train	0.758	0.751	0.748	0.749
	Test	0.751	0.743	0.739	0.741
AE-5-512	Train	0.763	0.757	0.753	0.755
	Test	0.758	0.752	0.749	0.750
AE-5-1024	Train	0.769	0.764	0.758	0.761
	Test	0.761	0.754	0.759	0.756

Table 5 contains the CatBoost results, which yielded the strongest performance. Test accuracy ranged from 0.742 for AE-3-128 to 0.768 for AE-5-1024. Precision values reached 0.774, and recall values reached 0.744, producing F1-scores as high as 0.759. These outcomes confirm that CatBoost benefited most from the distributed feature extraction process and that the deepest autoencoders generated the most discriminative representations.

**Table 5. Classification Results – CatBoost**

AE ID	Dataset	Accuracy	Precision	Recall	F1-score
AE-3-128	Train	0.735	0.728	0.721	0.724
	Test	0.726	0.719	0.713	0.716
AE-3-256	Train	0.745	0.739	0.731	0.735
	Test	0.739	0.730	0.722	0.726
AE-3-512	Train	0.752	0.746	0.737	0.741
	Test	0.747	0.741	0.731	0.736
AE-3-1024	Train	0.754	0.745	0.748	0.746
	Test	0.743	0.738	0.734	0.736
AE-4-128	Train	0.755	0.749	0.742	0.745
	Test	0.744	0.738	0.726	0.732
AE-4-256	Train	0.761	0.756	0.750	0.753
	Test	0.753	0.747	0.739	0.743
AE-4-512	Train	0.767	0.761	0.755	0.758
	Test	0.758	0.751	0.748	0.749
AE-4-1024	Train	0.771	0.764	0.758	0.761
	Test	0.759	0.752	0.750	0.751
AE-5-128	Train	0.763	0.756	0.748	0.752
	Test	0.746	0.740	0.733	0.736
AE-5-256	Train	0.767	0.760	0.754	0.757
	Test	0.757	0.751	0.747	0.749
AE-5-512	Train	0.772	0.766	0.760	0.763
	Test	0.764	0.758	0.751	0.754
AE-5-1024	Train	0.776	0.770	0.765	0.767
	Test	0.768	0.762	0.756	0.759

It was noted that even the smallest bottleneck sizes produced usable representations, since none of the classifiers showed performance collapse on the 128-dimensional vectors. This suggests that stable structural information was captured early in the encoding process. It was also observed that no major imbalance appeared between precision and recall across models, which indicates that the distributed training procedure did not introduce bias toward either class during feature extraction.

### **5. Conclusion and Discussion**

This study demonstrated that high-resolution mammographic feature extraction can be performed effectively using a fully distributed, CPU-based training pipeline. A total



of 12 convolutional autoencoders were trained using 54,740 PNG images, derived from the RSNA Breast Cancer Screening dataset, across a four-node CPU cluster. Each node was equipped with an Intel Xeon W2135 processor and 256 GB of RAM, and training was executed synchronously using Keras's Multi-worker MirroredStrategy. Despite the absence of GPU acceleration, training time for each autoencoder remained within 16.1 to 18.7 minutes, confirming the feasibility of large-scale medical image processing on CPU-only infrastructure.

The autoencoders varied by architectural depth (three, four, or five convolutional layers) and bottleneck size (128, 256, 512, or 1024). Feature vectors extracted from the fully connected bottleneck layer were used to train logistic regression, XGBoost, and CatBoost classifiers. Results showed that CatBoost consistently outperformed the other classifiers across nearly all architectures, achieving a maximum test accuracy of 0.768 and an F1-score of 0.759 with the AE-5-1024 architecture. XGBoost, while slightly less accurate, remained competitive with a peak test accuracy of 0.761. Logistic regression achieved a highest test accuracy of 0.736, confirming that the extracted features were linearly informative.

Notably, deeper autoencoders with larger bottlenecks yielded better classification results, though improvements stopped meaningfully improving beyond 512 dimensions. Training remained efficient due to parallelization across CPU nodes, and performance bottlenecks were mitigated by balancing computation and network I/O via TensorFlow's coordinated worker setup.

These findings support the conclusion that distributed CPU-based training provides a practical and scalable solution for feature extraction in high-volume breast cancer screening datasets. The pipeline demonstrated here can be adopted in environments where GPU availability is limited, without sacrificing performance or throughput. Future extensions may involve supervised fine-tuning or deployment in edge-based diagnostic platforms.

## Reference

- Baur, C., Wiestler, B., Albarqouni, S., and Navab, N. (2021). Autoencoders for unsupervised anomaly segmentation in brain MR images: *A comparative study. Medical Image Analysis*, 69, 101952. doi:10.1016/j.media.2021.101952
- Chen, T., and Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794. doi:10.1145/2939672.2939785
- Chen, Y., Li, H., Huang, P., and Zheng, Y. (2019). Deep learning with hierarchical convolutional autoencoders for automated detection of pulmonary nodules. *IEEE Access*, 7, 130000–130009. doi:10.1109/ACCESS.2019.2940017
- Dorogush, A.V., Ershov, V., and Gulin, A. (2018). CatBoost: Gradient boosting with categorical features support. *arXiv preprint*, arXiv:1810.11363
- Goyal, P., Dollár, P., et al. (2017). Accurate, large minibatch SGD: *Training ImageNet in 1 hour. arXiv preprint*, arXiv:1706.02677

Ribli, D., Horváth, A., Unger, Z., Pollner, P., and Csabai, I. (2018). Detecting and classifying lesions in mammograms with deep learning. *Scientific Reports*, 8(1), 4165. doi:10.1038/s41598-018-22437-z

Shen, L., Margolies, L.R., Rothstein, J.H., Fluder, E., McBride, R., and Sieh, W. (2019). Deep learning to improve breast cancer detection on mammography. *Scientific Reports*, 9(1), 12495. doi:10.1038/s41598-019-48995-4

Suzuki, K. (2018). Overview of deep learning in medical imaging. *Radiological Physics and Technology*, 10(3), 257–273. doi:10.1007/s12194-017-0406-5

Tajbakhsh, N., Shin, J.Y., et al. (2016). Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE Transactions on Medical Imaging*, 35(5), 1299–1312. doi:10.1109/TMI.2016.2535302

TensorFlow Documentation. (2023). Distributed training with TensorFlow. TensorFlow. Available at: [https://www.tensorflow.org/guide/distributed\\_training](https://www.tensorflow.org/guide/distributed_training) (Accessed 4 June 2025)

**Submitted: 20.06.2025**

**Accepted: 10.10.2025**