

A Novel Transparency Strategy-based Data Augmentation Approach for BI-RADS Classification of Mammograms

Sam B. Tran[†]
Vingroup Big Data Institute
10000, Hanoi, Vietnam
samsamtranbao@gmail.com

Huyen T. X. Nguyen[†]
Vingroup Big Data Institute
10000, Hanoi, Vietnam
nguyenhuyen.csp.bka@gmail.com

Chi Phan[†]
VinUni-Illinois Smart Health Center,
VinUniversity, 10000 Hanoi, Vietnam
21chi.pth@vinuni.edu.vn

Ha Q. Nguyen
Vingroup Big Data Institute,
VinBigData JSC,
10000 Hanoi, Vietnam
v.hanq3@vinbigdata.org

Hieu H. Pham
VinUni-Illinois Smart Health Center,
College of Engineering and Computer Science,
VinUniversity, 10000 Hanoi, Vietnam
Correspondence: hieu.ph@vinuni.edu.vn

[†] These authors contribute equally to the work and share the first authorship.

Abstract—Image augmentation techniques have been widely investigated to improve the performance of deep learning (DL) algorithms on mammography classification tasks. Recent methods have proved the efficiency of image augmentation on data deficiency or data imbalance issues. In this paper, we propose a novel transparency strategy to boost the Breast Imaging Reporting and Data System (BI-RADS) scores of mammogram classifiers. The proposed approach utilizes the Region of Interest (ROI) information to generate more high-risk training examples for breast cancer (BI-RADS 3, 4, 5) from original images. Our extensive experiments on three different datasets show that the proposed approach significantly improves the mammogram classification performance and surpasses a state-of-the-art data augmentation technique called CutMix. This study also highlights that our transparency method is more effective than other augmentation strategies for BI-RADS classification and can be widely applied to other computer vision tasks.

Index Terms—Mammogram, deep learning, data augmentation, abnormality detection, BI-RADS classification.

I. INTRODUCTION

Breast cancer has currently become the most common cancer, based on statistics from the International Agency for Research on Cancer (IARC) in December 2020 [11]. The American Cancer Society (ACS) stated that the average hazard proportion of a woman in the United States developing breast cancer during her life was about 13% [1]. WHO estimated that 2.3 million women were diagnosed with breast cancer, and there were about 685,000 deaths worldwide in 2020 [18]. The expert recommendation for women at high risk of breast cancer is to take diagnostic screening annually to detect cancer earlier and receive effective treatments beforehand. Mammography is a prevalent X-ray examination for breasts and is employed in computer-aided diagnosis (CADx) systems to assist radiologists in assessing breast cancer risk. In

particular, the BI-RADS score is used as a risk evaluation and quality assurance tool that supplies a widely accepted lexicon and reporting schema for breast imaging [5]. This standard consists of seven assessment levels: BI-RADS 0 (incomplete), BI-RADS 1 (negative), BI-RADS 2 (benign), BI-RADS 3 (probably benign), BI-RADS 4 (suspicious for malignancy), BI-RADS 5 (highly suggestive of malignancy), and BI-RADS 6 (known biopsy-proven malignancy). There are several recent DL-based studies which focus on BI-RADS classification [10], [17] rather than benign/malignant classification. However, malignancy cases are often much less than benign cases, leading to data issues consisting of data deficiency and data imbalance.

Recently, several data augmentation strategies have been proposed to resolve this problem and boost further training efficiency [2], [4], [15], [19], [20]. There are some deep learning techniques for medical image synthesis, but mainly using generative adversarial networks (GANs) [2], [15]. The potential of GANs for image processing issues is enormous because they can be trained to mimic any dataset. However, several systematic reviews of GANs [3], [15] reveal that there are still many challenging issues in using this technique for medical imaging, including training instability, computational cost, and concerns regarding the quality of the generated samples. Whereas, some basic augmentation techniques such as cropping, filtering, flipping, and noise injection are simple to apply but still adequate for model performance improvement [2]. Two combination methods including [19], [20] could achieve remarkable results because a new image is generated by integrating some original images. Additionally, we observed that lesion areas in medical imaging, especially in mammograms, play an essential role, but have not been focused on by most current image augmentation techniques [12].

This has motivated us to propose a new basic augmentation technique, called Transparency, for mammography BI-RADS classification. We conducted extensive experiments to show that the proposed approach has higher performance compared to a state-of-the-art data augmentation technique called CutMix [19] in mammograms.

Our main contribution in this work is developing a transparency data augmentation technique that can generate new compelling images based on original images. A new image focuses on lesion areas without losing global image context by blurring the original image except for lesion areas. The new image would still have the same distribution as the original one and a deep focus on lesions. The proposed method is easy to apply to any medical dataset whose images have lesion bounding boxes and can be widely applied to other computer vision tasks. Experimental results indicate that our transparency technique outperforms the typical augmentation techniques by 2.7%, 4.3%, and 6.9% on the MIAS, our private dataset and VinDr-Mammo, respectively. It also surpasses the state-of-the-art CutMix [19] on our private mammography dataset and VinDr-Mammo. The rest of the paper is organized as follows. The methodology is provided in Section II. Our experiments are presented in Section III. Finally, Section IV discusses the experimental results and concludes the paper.

II. METHODOLOGY

A. Preprocessing

One of the essential methods to improve the performance of machine learning models is data pre-processing, which is a group of different techniques to generate the most informative dataset for the training process. This section provides a brief detail of pre-processing data methods for our BI-RADS classification research. Four strategies Cropping, Flipping, CutMix, and Transparency will be discussed in this section.

1) *Image cropping and flipping*: One limitation of mammograms is that they contain a large black area without valuable information. Hence, we implement a breast detection model based on YOLOv5 [7] to crop the breast out of the original image. First, we used the labeling tool to label 2,000 images from the internal dataset. Then a dataset is built with 1,600 samples for training and 400 samples for validation. The medium YOLOv5 [7] achieved the mAP of 0.995 for breast detection during testing. The cropping step can reduce input size before entering the DL model, avoiding wasting resources and shortening the computational time. An example of this process is illustrated in Figure 1. We then flipped all the right mammograms along the vertical axis. By this way, we obtained a dataset where all the breast images have the same orientation.

2) *CutMix algorithm*: The CutMix augmentation strategy [19] was applied to mammograms to increase the number of unusual samples. The key idea behind the CutMix algorithm is to create a new pattern by merging the interpolation of both

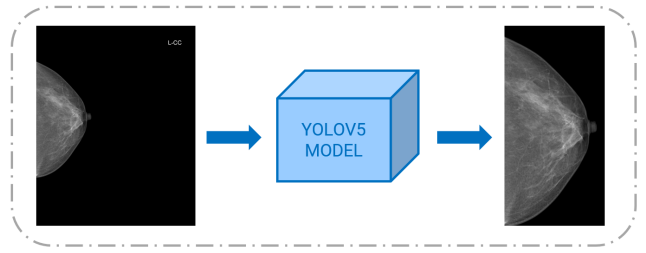


Fig. 1: Image cropping by YOLOv5 [7].

images and two labels. In this study, we simply moved the lesion region from high-risk scans to low-risk ones by changing pixel values. We denote that $x \in \mathbb{R}^{W \times H \times C}$ and y correspond to a mammography sample and its label where W, H, C are width, height, and channels of this sample, respectively. First, we extract the lesion area from image A (x_A, y_A) with abnormal bounding boxes $Box = (x_{min}, x_{max}, y_{min}, y_{max})$ and ground truth labels $y_A \in \{\text{“BIRADS 3”}, \text{“BIRADS 4”}, \text{“BIRADS 5”}\}$. A mask $M \in \{0; 1\}^{W \times H}$ is created by reducing the background pixel value to 0 and keeping the box pixels to 1 as

$$M_{ij} = \begin{cases} 1; & \text{if } (x_{min} \leq i \leq x_{max}) \& (y_{min} \leq j \leq y_{max}) \\ 0; & \text{if } (x_{min} > i \mid i > x_{max} \mid y_{min} > j \mid j > y_{max}). \end{cases} \quad (1)$$

We then generate the new sample (x', y') through the equation

$$\begin{aligned} x' &= M \odot x_A + (1 - M) \odot x_B, \\ y' &= y_A, \end{aligned} \quad (2)$$

where \odot denotes element-wise multiplication, image B (x_B, y_B) with $y_B \in \{\text{“BIRADS 1”}, \text{“BIRADS 2”}\}$.

3) *Transparency algorithm*: Due to the lack of data, especially BI-RADS 3, BI-RADS 4, BI-RADS 5 mammograms, an image augmentation strategy called “Transparency Strategy” was proposed to generate more high-risk samples from original images. The minimum condition for implementing this algorithm is that the training dataset contains images with unusual bounding boxes. A new training sample (x', y') is created by transforming pixel values from the original sample (x, y) . The transformation operation will be described in detail below

$$\begin{aligned} x' &= M \odot x, \\ y' &= y. \end{aligned} \quad (3)$$

This operator is a label-preserving transformation, BI-RADS of the new sample is the original label. The sample increment algorithm artificially increases the training dataset size by preserving the pixels at the lesion and reducing the background pixel values of the bounding box image. For abnormal images, $Box = (x_{min}, x_{max}, y_{min}, y_{max})$ indicates the bounding box location. In the mask, $M \in \{\alpha; 1\}^{W \times H}$, α is a random number that ranges from 0.1 to 0.9, and 1 is the value of pixels that is inside the lesion bounding box. The range of alpha was chosen empirically. Then the mask has been created with the following formulas

$$M_{ij} = \begin{cases} 1; & \text{if } (x_{min} \leq i \leq x_{max}) \& (y_{min} \leq j \leq y_{max}) \\ \alpha; & \text{if } (x_{min} > i \mid i > x_{max} \mid y_{min} > j \mid j > y_{max}). \end{cases} \quad (4)$$

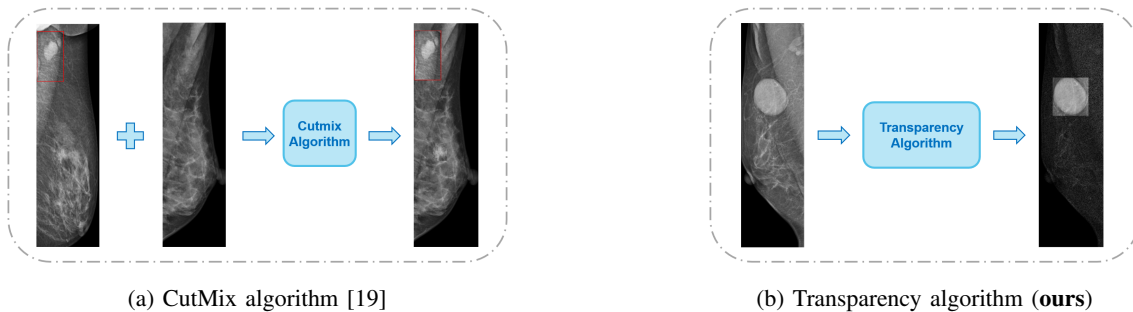


Fig. 2: Illustration of the CutMix augmentation algorithm and the proposed Transparency algorithm.

Before putting the data into the models, we define a data loader that loads the data and generates more samples based on the Transparency algorithm. In comparing several data augmentation techniques that use information from abnormal bounding boxes, such as CutMix [19], Mixup [20], Cutout [4] the advantages of the proposed approach are using the whole breast image, utilizing informative lesion location, and adjusting the weight of the region needed to focus. See Figure 2 for more details.

B. Deep Learning BI-RADS Classification

We implemented a state-of-the-art deep learning model to classify BI-RADS on mammograms in this work. Our network consists of two main components: (i) an extraction feature based on the Efficientnet-B2 [16] architecture, the output of which is a feature representation for each sample image, and (ii) a fully connected layer as a classifier to predict results from computed representations. We refer the readers to Tan *et al.* [16] for more detail about the network architecture.

III. EXPERIMENTS

A. Datasets and experimental setup

We evaluate the effectiveness of the proposed approach on three datasets: our private dataset, VinDr-Mammo [9], and MIAS [6].

Our private dataset was collected from Hospital 108 (H108) and Hanoi Medical University Hospital (HMHU) from 2018 to 2020. The dataset includes 36,614 screening mammogram images that come with their BI-RADS classification; each image was annotated by a team of three radiologists specializing in breast imaging for global labels (BI-RADS 1–5). These images were divided into three groups by the global label stratification method [14]: training set (25,373), validation set (5,398), and test set (5,393). In the dataset, 2,503 images were remarked on three local labels (lesions), including Discrete Mass, Spiculated Mass, and Stellate Mass with bounding boxes. The number of lesions on BI-RADS 3 & 4 & 5 occupies almost a total of mass images that contains local labels. Descriptions of three sets on our internal dataset are provided in Table I.

VinDr-Mammo dataset is a large-scale benchmark full-field digital mammography dataset consisting of 5,000 four-view

exams with breast-level BI-RADS and findings annotations. This is currently the largest public dataset with 20,000 scans providing BI-RADS assessment and suspicious/probably benign findings. Mammography images were acquired retrospectively from H108 and HMUH. We divided this dataset by the given stratification method in [14]. Therefore, one-fifth (4,000 images) of the VinDr-Mammo dataset was used for testing and the rest part (16,000 images) of dataset was used for training. **MIAS dataset** was collected from the United Kingdom National Breast Screening Program and then labeled by specialists in 1994. This dataset contains 322 images that were divided into three classes: 209 images for normal, 61 images for benign, and 52 images for malignant class. We randomly stratified the dataset into training and test sets with a ratio of 0.8/0.2. As a result, there are 265 images used for training and 67 images for testing the classification algorithms.

B. Training Methodology & Evaluation Metrics

Our experiments were built on PyTorch and using a PC with an Nvidia GTX 1080 GPU. We trained the feature extractors using SGD optimizer [13] with a momentum of 0.9 and cosine annealing learning rate [8]. The cross-entropy function was used to calculate the error. For model evaluation, we used *F1*-score on the 5-class BI-RADS level. *F1*-score is the harmonic mean of precision and recall. For multi-class problems, the *F1*-score macro, which is defined as the mean of class-wise *F1*-scores could be used. In our experiments, the results are appraised on image-level for BI-RADS classification. The classification models for different augmentation methods are trained with the same network architecture (EfficientNet-B2 [16]) and a fixed image size of 1024×768 . The number of epochs was set to 50, and the training process stopped in case there was no improvement in the *F1*-score of the validation set after 15 consecutive epochs by an early stopping callback. The performance of different techniques is assessed on the test set with the same evaluation metrics and network architecture.

C. Experimental Results

This section reports the effectiveness of the proposed Transparency approach and compares its performance with two other data augmentation techniques: baseline and CutMix [19]. We observed that our method consistently outperforms the baseline and the state-of-the-art method by a

TABLE I: Description of the private dataset used for model development and validation.

Data		Training set					Validation set					Test set				
BI-RADS		1	2	3	4	5	1	2	3	4	5	1	2	3	4	5
Lesions																
Discrete mass		2	31	1,077	779	153	0	3	245	193	25	2	7	237	168	35
Spiculated mass		0	1	1	49	226	0	0	0	11	41	0	0	0	18	48
Stellate mass		0	0	1	70	95	0	0	0	8	21	0	0	3	17	21
Total mass images		2	13	901	571	222	0	2	217	134	38	2	4	201	138	58
Total images		16,203	5,435	1,699	1,514	522	3,385	1,208	372	325	108	3,376	1,191	369	338	119

TABLE II: Quantitative results (F1-score) of baseline and Transparency technique on the private dataset, VinDr-Mammo, and MIAS datasets. The highest score are highlighted in **bold**.

Dataset	Class	Baseline	Transparency (ours)
Private dataset	BI-RADS 1	0.85	0.86
	BI-RADS 2	0.59	0.61
	BI-RADS 3	0.21	0.31
	BI-RADS 4	0.52	0.57
	BI-RADS 5	0.60	0.63
	Macro - F1	0.552	0.595
VinDr-Mammo	BI-RADS 1	0.89	0.89
	BI-RADS 2	0.59	0.61
	BI-RADS 3	0.55	0.68
	BI-RADS 4	0.51	0.50
	BI-RADS 5	0.50	0.69
	Macro - F1	0.607	0.676
MIAS	Normal	0.70	0.71
	Benign	0.32	0.28
	Malignant	0.27	0.37
	Macro-F1	0.428	0.455

large margin in different settings.

1) *Comparison with the baseline*: The baseline model is trained on our private dataset, VinDr-Mammo, and MIAS datasets, in which all images are cropped with the breast detector and flipped before fitting into the model. Compared to the baseline, the proposed transparency approach achieved better classification performance of about 4.3% higher in F1-scores on the private dataset, 6.9% and 2.7% higher on the VinDr-Mammo and MIAS dataset, respectively. Table II reports the experimental result on these three datasets.

2) *Comparison with a state-of-the-art technique*: Experiments on the private dataset and VinDr-Mammo showed that both CutMix [19] and the proposed transparency technique considerably enhance the performance of the classifier. In particular, the proposed method performed significantly better than the CutMix [19] algorithm for all classes on the VinDr-Mammo dataset, and for three classes such as BI-RADS 1, BI-RADS 2, and BI-RADS 3 on our in-house dataset. In particular, the Macro-F1 of the proposed method surpassed the CutMix algorithm by around 1% on the private dataset and by 6.5% on VinDr-Mammo. The experimental result on these two datasets was illustrated in Table III.

IV. CONCLUSION

Resolving data imbalance is one of the most challenging problems in machine learning, especially in medical image analysis, where anomalous patterns are rare and hard to

TABLE III: Quantitative results (F1-score) using CutMix and Transparency technique on the private dataset and VinDr-Mammo dataset. The highest score are highlighted in **bold**.

Class	Private dataset		VinDr-Mammo	
	CutMix	Transparency (ours)	CutMix	Transparency (ours)
BI-RADS 1	0.85	0.86	0.89	0.89
BI-RADS 2	0.59	0.61	0.59	0.61
BI-RADS 3	0.25	0.31	0.53	0.68
BI-RADS 4	0.57	0.57	0.48	0.50
BI-RADS 5	0.68	0.63	0.56	0.69
Macro-F1	0.589	0.595	0.611	0.676

collect. This study introduced a Transparency strategy-based technique that creates diseased samples by adjusting the pixel values. Experimentally, the proposed approach shows strong evidence that it could improve the BI-RADS classification task on mammogram exams. Our approach is simple and can be applied to various tasks in medical imaging, especially for lesion detection and classification problems.

Compliance with Ethical Standards. Our work follows all applicable ethical research standards and laws. The study has been reviewed and approved by the hospital’s institutional review board (IRB). The need for obtaining informed patient consent was waived because this work did not impact clinical care.

Acknowledgements. This study was supported by Smart Health Center, VinBigData JSC. We would like to acknowledge Hanoi Medical University Hospital for providing access to their image databases. In particular, we thank all of our radiologists who participated in this project.

REFERENCES

- [1] American Cancer Society. How Common Is Breast Cancer? <https://www.cancer.org/cancer/breast-cancer/about/how-common-is-breast-cancer.html>. Accessed Aug. 11, 2021 [Online].
- [2] P. Chlap, M. Hang, N. Vandenberg, J. A. Dowling, L. Holloway, and A. Haworth. A review of medical image data augmentation techniques for deep learning applications. In *Journal of Medical Imaging and Radiation Oncology*, 2016.
- [3] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath. Generative adversarial networks: An overview. volume 35, pages 53–65, 2018.
- [4] T. DeVries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- [5] Dr Daniel J Bell and Dr Yuranga Weerakkody et al. Breast imaging-reporting and data system (BI-RADS). <https://radiopaedia.org/articles/breast-imaging-reporting-and-data-system-bi-rads>. Accessed Aug. 11, 2021 [Online].
- [6] J. S. et al. Mammographic image analysis society (mias) database v1.21. 2015.

- [7] G. Jocher. Yolov5. <https://github.com/ultralytics/yolov5>, 2020.
- [8] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- [9] H. T. Nguyen, H. Q. Nguyen, H. H. Pham, K. Lam, L. T. Le, M. Dao, and V. Vu. Vindr-mammo: A large-scale benchmark dataset for computer-aided diagnosis in full-field digital mammography. *medRxiv*, 2022.
- [10] H. T. X. Nguyen, S. B. Tran, D. B. Nguyen, H. H. Pham, and H. Q. Nguyen. A novel multi-view deep learning approach for bi-rads and density assessment of mammograms. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 2144–2148, 2022.
- [11] W. H. Organization. Latest global cancer data: Cancer burden rises to 19.3 million new cases and 10.0 million cancer deaths in 2020. International Agency for Research on Cancer, Lyon, France, Dec. 15, 2020 [Online].
- [12] P. Oza, P. Sharma, S. Patel, F. Adedoyin, and A. Bruno. Image augmentation techniques for mammogram analysis. In *Multidisciplinary Digital Publishing Institute (MDPI)*, 2022.
- [13] S. Ruder. An overview of gradient descent optimization algorithms. In *arXiv preprint arXiv:1609.04747*, 2016.
- [14] K. Sechidis, G. Tsoumakas, and I. Vlahavas. On the stratification of multi-label data. In D. Gunopulos, T. Hofmann, D. Malerba, and M. Vazirgiannis, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 145–158, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [15] Y. Skandarani, P.-M. Jodoin, and A. Lalonde. Gans for medical image synthesis: An empirical study. In *Preprint arXiv*, 2021.
- [16] M. Tan and Q. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning (ICML)*, pages 6105–6114, 2019.
- [17] K.-J. Tsai, M.-C. Chou, H.-M. Li, S.-T. Liu, J.-H. Hsu, W.-C. Yeh, C.-M. Hung, C.-Y. Yeh, and S.-H. Hwang. A high-performance deep neural network model for bi-rads classification of screening mammography. *Sensors*, 22(3), 2022.
- [18] World Health Organization. Breast Cancer. <https://www.who.int/news-room/fact-sheets/detail/breast-cancer>. Accessed Aug. 11, 2021 [Online].
- [19] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. J. Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6022–6031, 2019.
- [20] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. Mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.