# Multi-Institutional Assessment and Crowdsourcing Evaluation of Deep Learning for Automated Classification of Breast Density



Ken Chang, PhD<sup>a</sup>, Andrew L. Beers, AB<sup>a</sup>, Laura Brink, BS<sup>b</sup>, Jay B. Patel, BS<sup>a</sup>, Praveer Singh, PhD<sup>a</sup>, Nishanth T. Arun, BTech<sup>a</sup>, Katharina V. Hoebel, MD<sup>a</sup>, Nathan Gaw, PhD<sup>a</sup>, Meesam Shah, BS<sup>b</sup>, Etta D. Pisano, MD<sup>c,d</sup>, Mike Tilkin, MS<sup>e</sup>, Laura P. Coombs, PhD<sup>f</sup>, Keith J. Dreyer, DO, PhD<sup>g,h,i</sup>, Bibb Allen, MD<sup>j,k,l</sup>, Sheela Agarwal, MD, MBA<sup>m</sup>, Jayashree Kalpathy-Cramer, PhD<sup>a,n</sup>

Credits awarded for this enduring activity are designated "SA-CME" by the American Board of Radiology (ABR) and qualify toward fulfilling requirements for Maintenance of Certification (MOC) Part II: Lifelong Learning and Self-assessment. To access the SA-CME activity visit https://cortex.acr.org/Presenters/CaseScript/CaseView?Info=6VA3oEyni0LFSnDR5TF78POxFum8xUX0f8lQB91nLQc %253d. SA-CME credit for this article expires July 31, 2020.

#### Abstract

Objective: We developed deep learning algorithms to automatically assess BI-RADS breast density.

**Methods:** Using a large multi-institution patient cohort of 108,230 digital screening mammograms from the Digital Mammographic Imaging Screening Trial, we investigated the effect of data, model, and training parameters on overall model performance and provided crowdsourcing evaluation from the attendees of the ACR 2019 Annual Meeting.

<sup>b</sup>American College of Radiology, Reston, Virginia.

<sup>e</sup>Chief Information Officer and EVP for Technology (ACR), Reston, Virginia.

<sup>h</sup>Associate Professor of Radiology,Harvard Medical School, Boston, Massachusetts.

<sup>i</sup>Chief Science Officer, ACR Data Science Institute, Reston, Virginia.

<sup>k</sup>Secretary General, International Society of Radiology, Reston, Virginia.

<sup>&</sup>lt;sup>a</sup>Athinoula A. Martinos Center for Biomedical Imaging, Department of Radiology, Massachusetts General Hospital, Boston, Massachusetts.

<sup>&</sup>lt;sup>c</sup>Chief Research Officer (ACR), Reston, Virginia.

<sup>&</sup>lt;sup>d</sup>Professor in Residence, Beth Israel Lahey/Harvard Medical School, Boston, Massachusetts.

<sup>&</sup>lt;sup>f</sup>Vice President (ACR), Reston, Virginia.

<sup>&</sup>lt;sup>g</sup>Chief Data Science Officer, Chief Imaging Information Officer, Massachussetts General Hospital and Brigham Women's Hospital (MGH & BWH), Chief Executive, MGH & BWH Center for Clinical Data Science; Vice Chairman of Radiology – Informatics, MGH & BWH, Boston, Massachusetts.

<sup>&</sup>lt;sup>1</sup>Chief Medical Office, ACR Data Science Institute, Reston, Virginia.

<sup>&</sup>lt;sup>1</sup>Partner, Grandview Medical Center, Birmingham, Alabama.

<sup>&</sup>lt;sup>m</sup>Lennox Hill Radiology, New York, New York.

<sup>&</sup>lt;sup>n</sup>Scientific Director (CCDS), Director (QTIM lab and the Center for Machine Learning), Associate Professor of Radiology, MGH/Harvard Medical School, Boston, Massachusetts.

Corresponding author and reprints: Jayashree Kalpathy-Cramer, PhD, Athinoula A. Martinos Center for Biomedical Imaging, 149 13th Street, Charlestown, MA 02129; e-mail: kalpathy@nmr.mgh.harvard.edu.

J. Patel reports grants from National Institute of Biomedical Imaging and Bioengineering, outside the submitted work. Dr Kalpathy-Cramer reports grants from National Cancer Institute, personal fees from Infotech, Soft, outside the submitted work. Dr Hoebel reports grants from Martinos Scholars fund, outside the submitted work. Dr Chang reports grants from National Institute of Biomedical Imaging and Bioengineering, grants from National Cancer Institute, outside the submitted work. Dr Agarwal reports personal fees from Bayer HealthCare, outside the submitted work. The other authors state that they have no conflict of interest related to the material discussed in this article.

**Results:** Our best-performing algorithm achieved good agreement with radiologists who were qualified interpreters of mammograms, with a four-class  $\kappa$  of 0.667. When training was performed with randomly sampled images from the data set versus sampling equal number of images from each density category, the model predictions were biased away from the low-prevalence categories such as extremely dense breasts. The net result was an increase in sensitivity and a decrease in specificity for predicting dense breasts for equal class compared with random sampling. We also found that the performance of the model degrades when we evaluate on digital mammography data formats that differ from the one that we trained on, emphasizing the importance of multi-institutional training sets. Lastly, we showed that crowdsourced annotations, including those from attendees who routinely read mammograms, had higher agreement with our algorithm than with the original interpreting radiologists.

**Conclusion:** We demonstrated the possible parameters that can influence the performance of the model and how crowdsourcing can be used for evaluation. This study was performed in tandem with the development of the ACR AI-LAB, a platform for democratizing artificial intelligence.

Key Words: ACR AI-LAB, artificial intelligence, BI-RADS, breast density, deep learning, DMIST, generalizability, mammogram, neural networks

J Am Coll Radiol 2020;17:1653-1662. Copyright © 2020 Published by Elsevier Inc. on behalf of American College of Radiology

### INTRODUCTION

Breast cancer is one of the leading causes of death among women in the United States, with the number of deaths expected to be over 41,000 in 2019 [1]. Early mammographic screening has resulted in a decrease in breast cancer mortality [2,3]. The correct mammographic interpretation of breast density, which measures extent of fibroglandular tissue, is important in the assessment of breast cancer risk because there is increased risk with increased density [4,5]. Furthermore, the identification of dense breast may stratify patients who may have masked cancers and may benefit from additional ultrasound or MRI. As such, there is now legislation in many states that patients must be notified of their breast density after mammography [6].

Qualitative assessment by means of the widely used BI-RADS include four categories: (a) almost entirely fatty, (b) scattered fibroglandular densities, (c) heterogeneously dense, or (d) extremely dense [7]. These criteria are subjective, resulting in interrater variability among radiologists. A study by Sprague et al showed that the likelihood of any given mammogram being rated as dense (heterogeneously dense and extremely dense) is highly dependent on the interpreting radiologist, with the percentage ranging from 6.3% to 84.5% [8]. Other studies have reported intrareader variability to be  $\kappa = 0.58$  (among 34 community radiologists) and the interrater variability to be  $\kappa = 0.643$ (between a consensus of five breast radiologists and the original interpreting breast radiologist) [6,9]. Similarly, commercially available software shows a wide range of agreement with clinical experts, and the probability of dense classification is dependent on the specific software used [10,11]. This intra- and interrater variability, and even

intersoftware variability, may confer undue patient anxiety and potential harm to the patient (ie, possible unnecessary supplemental screening examinations).

As such, there has been interest in using automated approaches to improve accuracy and consistency of breast density assessment. Commercial software uses quantitative imaging features to assess breast density, with mixed agreement with radiologist interpretation [11]. Deep learning methods have yielded state-of-the-art results in a wide range of computer vision tasks without the need for domaininspired handcrafted imaging features. Moreover, recent studies have shown the potential of deep learning in medical fields such as dermatology, ophthalmology, and radiology [12-14]. A recent study from Lehman et al demonstrates the utility of deep learning for mammographic density assessment in clinical practice at a single institution and mammography system [6]. Here, we further this work by validating the deep learning approach on a multiinstitutional imaging cohort with a variety of digital mammography systems. Furthermore, we provide an indepth analysis of how choice of data, model, and training parameters affects algorithm performance. In addition to that, we investigate the generalizability of models across different digital mammography data formats. Lastly, we deploy our system at the ACR 2019 Annual Meeting for a crowdsourced evaluation.

## MATERIALS AND METHODS Patient Cohort

Digital screening mammograms from 33 clinical sites were retrospectively obtained through the Digital Mammographic Imaging Screening Trial, the details of which were

mammograph syste	em	
12-Bit Monochrome 1	12-Bit Monochrome 2	14-Bit Monochrome 1
Senoscan, Fischer Medical, Wheat Ridge, CO (99.9%) Kodak Lumiscan 75, Rochester, NY (.1%)	Senographe, General Electric Medical Systems, Waukesha, WI (93.8%) Other (6.1%) Mammo-Clinical (.1%)	Senographe, General Electric Medical Systems, Waukesha, WI (94.1%) Mammo-Clinical (5.9%)

Table 1. Breakdown of data format by digitalmammograph system

previously published [15]. Each examination was interpreted by a single radiologist from a cohort of radiologists using ACR BI-RADS breast density lexicon (category a: fatty, category b: scattered, category c: heterogeneously dense, category d: extremely dense) [7]. A total of 92 radiologists read the examinations. Readers in the United States were all qualified interpreters of mammograms under federal law. Canadian readers met equivalent standards. Each site's lead radiologist received training to read for Digital Mammographic Imaging Screening Trial and in turn trained the site's other readers. All images were previously de-identified before this study. The mammograms were saved in DICOM format with four different image data formats, corresponding to different digital mammography systems or different versions of the same system (Table 1): 12-bit monochrome 1 (30.3%), 12-bit monochrome 2 (11.2%), 14-bit monochrome 1 (58.0%), and 14-bit monochrome 2 (0.5%). The 14-bit monochrome 2 images were excluded to ensure that each image data format included in our study had adequate representation for training of our deep learning model. Our final patient cohort consisted of 108,230 digital screening images from 21,759 patients (Table 2), which was divided into training, validation, and testing sets on the patient level. The training set was used to develop the model, and the validation set was used to assess model performance during training to prevent overfitting. The test set was unseen until the model training was complete.

## Experiments on Data, Model, and Training Parameters

Neural network models were implemented in DeepNeuro with Keras/TensorFlow backend [16]. We investigated the effect of data, model, and training parameters on algorithm performance. A schematic of the various experiments investigating data, model, and training parameters are summarized in Figure 1A. To investigate the effect of training set size, we used various different training set sizes and assessed the resulting performance on the test set. We tested four commonly used neural network architectures, each of which differ in number of layers and design: ResNet50, DenseNet121, InceptionV3, and VGG16 [17-21]. We also investigated the benefit of pretraining by comparing ImageNet (a large computer

Table 2. Summary of demographics in the patient cohort with regard to age, sex, race, and breast density

Demographic	Training (n = 62,316)	Validation (n = 6,978)	Testing (n = 38,936)
Age (median years, IQR)	46 (53-61)	46 (53-61)	47 (53-61)
Female (%)	100	100	100
Race White Black or African American Hispanic or Latino Asian American Indian or Alaska Other or unknown	50,414 8,389 2,273 819 63 358	5,622 925 289 62 8 72	30,845 5,733 1,416 633 19 290
Radiologist-assessed breast density Fatty Scattered Heterogeneously dense Extremely dense	6,980 (11.2%) 27,733 (44.5%) 23,987 (38.5%) 3,616 (5.8%)	873 (12.5%) 2,985 (42.8%) 2,753 (39.5%) 367 (5.3%)	4,575 (11.8%) 17,191 (44.2%) 14,585 (37.5%) 2,585 (6.6%)

IQR = interquartile range.



Fig. 1. (A) A summary of all the data, model, and training parameter experiments performed. (B) Performance on the testing set (measured by four-class  $\kappa$  agreement with radiologist interpretation) increased as the percentage of training set used. The 95% confidence interval is plotted in light green. (C) Effect of model and training parameters on testing set four-class  $\kappa$  agreement with radiologist interpretation. Black lines denote 95% confidence interval. \*P < .05, \*\*P < .01, \*\*\*\*P < .001.

vision data set of natural images) pretrained versus random initialization [22]. A variety of cost functions were also used (categorical cross-entropy, mean absolute error, mean squared error, and ordinal regression) to assess the effect of objective function (and their underlying assumptions of the nature of the labels) on performance [23]. The training set was augmented in real time by means of random horizontal or vertical flips (50% probability of each) and random rotations ( $0^{\circ}$ -45°). At each minibatch, images from each breast density class were sampled with either random (weighting in the empirical density class distribution) or equal class (weighting each density class equally) probability to assess the effect of class weighting on performance. We also evaluated the effect of model ensembling by averaging the output of two to four trained models of the same architecture (ResNet50). Model ensembling describes the process by which several independently trained models are combined to improve performance [24]. The default model used 100% of the training set, ImageNet pretrained weights, ResNet50 architecture, no ensembling, categorical cross-entropy loss function, augmentation, and equal class sampling. Only one parameter was modified at a time in the experiments, keeping all other parameters the same as the default model (ceteris paribus).

## **Experiments on Image Data Formats**

To visualize the differences in intensity distributions across image formats, histograms of preprocessed images from the testing set were generated. The dimensionality of histograms was then reduced to a two-dimensional projection and plotted to inspect for similarity across image formats [25]. The effect of image format of training images on generalizability of models was investigated. We trained ResNet50 models using 12-bit monochrome 1 images only, 12-bit monochrome 2 images only, 14-bit monochrome 1 images only, and all images. The performance of these models for each image format was then assessed. Projections of the intermediate output of the penultimate layer of the neural network were also plotted for images in the testing set using a model trained on all images to evaluate the learned features learned by the deep learning model. Further information about the patient cohort and experiments is available in the e-only Supplementary Information.

## Crowdsourcing Assessment

As further evaluation of our breast density algorithm, we deployed an annotation workstation at the ACR 2019 Annual Meeting. Attendees of all levels (researchers, medical students, residents, radiologists) were invited to

Table 3. Demographics of participants of thecrowdsourcing assessment	
Demographic	n
Experience Radiologist (breast) Radiologist (other) Resident Student	3 10 2 2
Read mammograms No Yes	10 7

perform annotations on a subset of images within our patient cohort. Representative images of all breast density classifications from the BI-RADS manual were provided to attendees during annotation. Attendees were able to inspect all images (all views available) from a given patient study and were asked to provide a BI-RADS breast density assessment. In total, 3,649 annotations were performed on 1,083 patient studies by 17 raters (demographics summarized in Table 3). On average, there were three annotations per patient study, and each rater performed 215 annotations. Consensus of the crowd was determined by majority vote, with random tiebreak. In our analysis, we looked at agreement between crowd and original interpreting radiologist annotation as well as crowd and algorithm (ResNet50). The ResNet50 model was chosen because it was the best-performing architecture among those tested. Only a single model (as opposed to an ensembled model) was used to reflect the common scenario in which only a single model is deployed for computation efficiency.

## **Statistical Analysis**

Agreement between raters was assessed via linear  $\kappa$  coefficient across the four breast density categories in the testing set (four-class  $\kappa$ ). For reference, a  $\kappa$  of 0.21 to 0.40, 0.41 to 0.60, and 0.61 to 0.80 represents fair, moderate, and substantial agreement, respectively [26].

## RESULTS

## Effect of Data Parameters on Performance

The performance of training set size on testing set performance was investigated, showing that K coefficient increases as the training set size increases. When 2% (n = 1,247 images) of the training set was used, the mean four-class K was 0.563 (95% confidence interval [CI], 0.551-0.576). In contrast, when using 100% (n = 62,316 images) of the training set, the mean four-class K was 0.660 (95% CI 0.657-0.664) (Fig. 1B). There was a statistically significant difference between the performance of using 2% to 60% and 100% of the training set (*t* test P < .05). There was no difference in the performance of using 80% and 100% of the training set (P = .291).

## Effect of Model Parameters on Performance

The mean four-class  $\kappa$  of models with randomly initialized weights was 0.327 (95% CI 0.273-0.384), compared with ImageNet pretrained weights 0.660 (95% CI 0.657-0.664, P < .001) when using the full training set (Fig. 1C). In the experiments assessing model architecture, the mean four-class  $\kappa$  of ResNet50, DenseNet121, InceptionV3, and VGG16 was 0.660 (95% CI 0.657-0.664), 0.650 (95% CI

0.640-0.659), 0.644 (95% CI 0.635-0.652), and 0.660 (95% CI 0.658-0.664), respectively. There was no statistically significant difference between the performance of the various architectures. The mean four-class K of no ensembling, ensembling two models, ensembling three models, and ensembling four models was 0.660 (95% CI 0.657-0.664), 0.665 (95% CI 0.664-0.666), 0.666 (95% CI 0.666-0.667), 0.667 (95% CI 0.666-0.668), respectively. The performance of ensembling four models and three models was greater than that of no ensembling (P = .041 and .036, respectively).

## Effect of Training Parameters on Performance

For categorical cross-entropy, mean absolute error, mean squared error, and ordinal regression, the mean four-class κ was 0.660 (95% CI 0.657-0.664), 0.649 (95% CI 0.644-0.653), 0.654 (95% CI 0.646-0.661), and 0.664 (95% CI 0.659-0.669), respectively. The performance of categorical cross-entropy and ordinal regression was significantly greater than mean absolute error (P = .011and P = .004, respectively). The mean four-class  $\kappa$  with no augmentation was 0.658 (95% CI 0.646-0.666), compared with augmentation 0.660 (95% CI 0.657-0.664; P = .675). The mean four-class  $\kappa$  with random and equal class sampling at each minibatch was 0.665 (95% CI 0.662-0.669) and 0.660 (95% CI 0.657-0.664), respectively (P = .135). For random class sampling, the predicted distribution of labels on the test set was 8.1% fatty, 47.5% scattered, 40.1% heterogeneously dense, and 4.3% extremely dense. This differed from the predicted distribution of labels on the test set with equal class sampling, which was 13.5% fatty, 37.5% scattered, 36.8% heterogeneously dense, and 12.2% extremely dense (P <.001, Fig. 2B). The predicted binary distribution for random (44.4% dense) and equal sampling (49.0% dense) also differed (P < .001). For random class sampling, the mean sensitivity and specificity for classifying dense breast was 0.833 (95% CI 0.803-0.856) and 0.888 (95% CI 0.872-0.905), respectively. Comparatively, for equal class sampling, there was an increase in sensitivity (0.880, 95% CI 0.869-0.890, P < .05) with a decrease in specificity (0.842, 95% CI 0.828-0.857, P < .001). A display of the range of classifications for models trained with different model and training parameters for 50 patients in the testing set is shown in Figure 2A.

## Effect of Digital Mammography Data Format on Model Generalizability

A plot of projections of intensity distributions of preprocessed images showed clustering within image format,



Fig. 2. (A) A visual display of the range of classifications for models trained with different model and training parameters for 50 patients in the testing set. The radiologist interpretation is displayed in the first row. The average breast density rating across all models and radiologist interpretation is displayed in the last row and was used to order the patients from least dense (left) to most dense (right). (B) The distribution of predicted breast density labels in the testing set differed for experiments with random class sampling (left) compared with equal class sampling (right) at each minibatch. \*\*\*\*P < .001. E. dense = extremely dense; H. dense = heterogeneously dense.

delineating differences between image formats (Fig. 3B). Clustering by intensity distribution was preserved even after passing the images through a trained neural network, as shown by projections of the output of the penultimate layer, with the grouping by breast density occurring within the respective image format cluster (Fig. 3D, E). For all image format-specific models, testing set performance was decreased for other image formats compared with the image format the model was trained on (P < .001). In

contrast, a model trained on all images showed no differences in performance across image formats (P > .05, Fig. 3C).

## Crowdsourcing Assessment

The four-class  $\kappa$  between the crowd and algorithm (0.505, 95% CI 0.503-0.506) was greater than agreement between crowd and original interpreting radiologist (0.463, 95% CI 0.461-0.464, *P* < .001). Agreement with the algorithm was



**Fig. 3.** (A) Intensity distribution histogram (frequency versus intensity value) of 100 randomly selected images of each pixel format. (B) Visualization of the histogram of intensities of 3,000 preprocessed images from the testing set demonstrating clustering of images by image format. (C) Performance of models trained on specific image formats as well as all images, showing that for image format-specific models, testing set performance was decreased for other image formats compared with the image format the model was trained on. (D, E) Visualization of an intermediate layer of the trained neural network for 3,000 images in the testing set, color-coded by image format and radiologist interpretation of breast density. E. dense = extremely dense; H. dense = heterogeneously dense.



Fig. 4. Confusion matrices showing the agreement between original interpreting radiologist, algorithm, and crowd. The agreement between the algorithm and crowd (B) was greater than the agreement between crowd and original interpreting radiologist (A). The agreement between algorithm and original interpreting radiologist for the same patient studies (C) shown for reference. (D) There was higher agreement, in terms of four-class  $\kappa$ , with the algorithm than with the original interpreting radiologist from the DMIST trial for both crowdsourcing participants who read mammograms and those who do not. \**P* < .001. E. dense = extremely dense; H. dense = heterogeneously dense.

greater than agreement with the original interpreting radiologist for both crowd participants who regularly read mammograms and those who do not (Fig. 4D). As a reference, the four-class  $\kappa$  between algorithm and radiologist was 0.636 (95% CI 0.635-0.637) for the same patient studies (Fig. 4A-C).

### DISCUSSION

In this study, we investigated the performance of deep learning models in a large multi-institution and multimammography system patient cohort. Our best-performing model achieved a  $\kappa$  of 0.667, equivalent to the agreement observed by Lehman et al, which only used mammograms from a single institution and mammography system [6].

One challenge of training robust deep learning models is the availability of large annotated imaging data sets [27]. In this study, we provide empirical evidence that the size of the training set is a key determinant in the performance of neural networks, consistent with another study on abnormality classification in chest radiographs [28]. In accordance with deep learning studies in other domains, tens of thousands of annotated images are needed before model performance begins to plateau in diverse imaging cohorts, supporting the need for collaborative efforts among medical institutions [28-30].

In our investigation of model parameters, pretraining and ensembling led to improvements in performance. Pretraining neural networks followed by fine-tuning in the target domain (also known as transfer learning) has become a well-established paradigm for medical imaging applications to achieve high performance [12,29]. In our study, we noted that pretraining on ImageNet improved performance for the breast density classification task. Further improvement in performance was seen with ensembling of independently trained models, which is analogous to how a consensus of experts is more likely to be correct than any single expert [31]. Interestingly, neural network architecture did not have a significant effect on performance despite differences in model complexity and design.

One important consideration when training a model is the objective function used to optimize the algorithm, also known as a cost function. Our experiments have shown that the choice of cost function had a significant effect on model performance, mainly because each cost function makes different assumptions about the nature of the labels. Specifically, mean absolute error, mean squared error, and ordinal regression assume that the categories are ordered, but categorical cross-entropy does not. Furthermore, mean absolute error and mean squared error assume the distance between adjacent classes is equal, whereas ordinal regression does not. In our application, breast density is classified on an ordered scale with undefined distances between adjacent classes (ie, the distances between fatty and scattered compared with heterogeneously dense and extremely dense cannot be quantified), making ordinal regression the most appropriate cost function. This is validated in our experiments, in which we find that ordinal regression exhibited the highest performance, although this was significantly different to only mean absolute error.

We also did not notice significant difference between random and equal class sampling on model performance in terms of  $\kappa$  coefficient. Class sampling is an important consideration in cases in which there are differences in the number of patient samples from each class (ie, when the majority class significantly outnumbers the minority class). In our study, we have more patients with scattered and heterogeneously dense breasts (44.2% and 37.5%, respectively) than with fatty and extremely dense breasts (11.8% and 6.6%, respectively), which is the expected distribution as breast density has a normal distribution. Under random class sampling, the neural network would be exposed to more training examples of scattered and heterogeneously dense breasts than of fatty and extremely dense breasts. Equal class sampling can be used to mitigate this inherent class imbalance by ensuring that the neural network is adequately exposed to all classes [32]. However, it is also important to note that with equal class sampling, the distribution of predicted labels changes-specifically, minority classes are predicted with higher frequency and majority classes are predicted with lower frequency, as shown by our experimental results. The net result of this is that the sensitivity of predicting dense breast improves with equal class sampling. Moreover, equal class sampling leads to lower specificity for classification of dense breast. From a policy perspective, this can lead to more patients being notified that they have dense breast. Additional imaging performed on these patients may lead to increases in the number of false-positives. This is a key example of how the manner in which deep learning models are trained can have implications for clinical care.

One critical hurdle that prevents the deployment of deep learning models in the clinical work environment is their relatively poor generalizability across institutional differences, such as patient demographics, disease prevalence, scanners, and acquisition settings. In fact, other deep learning studies that have shown poor generalizability of deep learning models when applied to data from different institutions than the one they were trained on [33,34]. In our study, we found that models trained on specific digital mammography data formats do not generalize to other data formats, and it was only after training on images from all digital mammography data formats did our model showed high performance on all data formats. Indeed, several deep learning studies for mammographic breast density assessment were only validated on patient cohorts from a single institution or digital mammography system [6,35,36]. Some possible differences between different digital mammography systems or versions of systems include the x-ray tube target, filter, digital detector technology, and control of automatic exposure [37]. Our results add to the growing body of literature that states that deep learning models do not necessarily generalize when applied to data that differs from that which the model was trained with.

Various studies have shown the utility of crowdsourcing and citizen science in biological and medical image annotation [38-41]. Crowdsourcing for annotation and evaluation is advantageous because it is scalable, high throughput, cost-efficient, and accurate [42-44]. As such, we performed a crowdsourcing assessment of our algorithm. Notably, there was a diversity of experience of the participants in crowdsourcing, with its inclusion of students, residents, and radiologists who do not routinely read mammograms. As such, it is unsurprising that the agreement between the crowd and algorithm was lower than the agreement between the original interpreting radiologist and algorithm. Interestingly, the crowd (both participants who routinely read mammograms and those who did not) had higher agreement with the algorithm than with the original interpreting radiologist. This may reflect the consistency of the algorithm in its assessment compared with the various interpreting radiologists from different sites in the Digital Mammographic Imaging Screening Trial study. In other words, a single algorithm may allow for greater consistency than having different human radiologists rate each imaging study.

There are several limitations to our study. The first is that we only had one radiologist, from a cohort of radiologists, perform interpretation for each patient study. Future studies will incorporate multiple readers for each patient study. In addition, for models initialized with random weights, we did not optimize training hyperparameters such as the learning rate schedule or the duration of training [45]. It is possible that optimization would improve the performance of the randomly initialized model, but in this study, we show the performance advantage of pretrained neural networks with minimal hyperparameter tuning. Furthermore, in our investigation of augmentation, we only explored random flips and rotations, though future work will explore other augmentation techniques such as intensity scaling and elastic deformations [46]. Lastly, our algorithm was only developed to assess mammographic breast density. Future work can extend our algorithm and crowdsourcing evaluation for more complex tasks such as assigning BI-RADS categories.

This study was developed in conjunction with the ACR AI-LAB, a framework for democratization of artificial intelligence. The goal of the ACR AI-LAB is to provide an interface for clinicians and scientists to work together to develop deep learning models. We highlight several fundamental features needed for artificial intelligence democratization: First, we demonstrate the possible data, model, and training parameters that can influence the performance of the model. These parameters will be available as options in AI-LAB. We also show importance of diverse training data for model generalizability, supporting collaborative development of algorithms across institutions which the AI-LAB will facilitate. Lastly, we show how crowdsourced annotations can be used to evaluate algorithm performance, which users will be able to do on the AI-LAB platform.

## Implications

We showcase the various data, training, and model parameters that can influence model performance, highlighting pretraining, cost function, and sampling approach as important parameters. Furthermore, we found that model performance deteriorates when training and testing on different imaging data formats. In performing this study in tandem with the development of the ACR AI-LAB, we demonstrate important principles and pitfalls that radiologists and data scientists have to consider when training neural network models. Our hope is that radiologists who use the AI-LAB can refer to this study as an educational tool when utilizing the AI-LAB to train their own deep learning models.

## TAKE-HOME POINTS

- The choice of data, model, and training parameters can impact deep learning model performance for evaluation of mammographic breast density. Notably, when training was performed with randomly sampled images from the data set versus sampling equal number of images from each density category, the model predictions were biased away from the low-prevalence categories such as extremely dense breasts.
- The performance of the model degrades when evaluated on data formats that differ from the one that we trained on, emphasizing the importance of multiinstitutional training sets.
- Crowdsourcing can be an effective means of evaluating model performance.
- These options for model training and evaluation will be made available in the ACR AI-LAB, a platform for democratizing artificial intelligence that was developed in tandem with this study.

## ACKNOWLEDGMENTS

Research reported in this publication was supported by a training grant from the National Institute of Biomedical Imaging and Bioengineering of the National Institutes of Health under award number 5T32EB1680 to K. Chang and J. B. Patel and by the National Cancer Institute of the National Institutes of Health under Award Number F30CA239407 to K. Chang. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. This publication was supported from the Martinos Scholars fund to K. Hoebel. Its contents are solely the responsibility of the authors and do not necessarily represent the official views of the Martinos Scholars fund to K. Hoebel. This publication from the Martinos Scholars fund to K. Hoebel. Its contents are solely the responsibility of the authors and do not necessarily represent the official views of the Martinos Scholars fund. This study was supported by

National Institutes of Health grants U01 CA154601, U24 CA180927, and U24 CA180918 to J. Kalpathy-Cramer. This research was carried out in whole or in part at the Athinoula A. Martinos Center for Biomedical Imaging at the Massachusetts General Hospital, using resources provided by the Center for Functional Neuroimaging Technologies, P41EB015896, a P41 Biotechnology Resource Grant supported by the National Institute of Biomedical Imaging and Bioengineering, National Institutes of Health. We acknowledge the computing resources provided by the MGH & BWH Center for Clinical Data Science and Amazon Web Services Machine Learning Research Awards.

### ADDITIONAL RESOURCES

Additional resources can be found online at: https://doi. org/10.1016/j.jacr.2020.05.015.

#### REFERENCES

- Siegel RL, Miller KD, Jemal A. Cancer statistics, 2019. CA Cancer J Clin 2019;69:7-34.
- Duffy SW, Tabár L, Chen H-H, et al. The impact of organized mammography service screening on breast carcinoma mortality in seven Swedish counties. Cancer 2002;95:458-69.
- Tabár L, Vitak B, Chen HH, Yen MF, Duffy SW, Smith RA. Beyond randomized controlled trials: organized mammographic screening substantially reduces breast carcinoma mortality. Cancer 2001;91: 1724-31.
- Boyd NF, Byng JW, Jong RA, et al. Quantitative classification of mammographic densities and breast cancer risk: results from the Canadian National Breast Screening Study. JNCI J Natl Cancer Inst 1995;87:670-5.
- Razzaghi H, Troester MA, Gierach GL, Olshan AF, Yankaskas BC, Millikan RC. Mammographic density and breast cancer risk in white and African American women. Breast Cancer Res Treat 2012;135: 571-80.
- Lehman CD, Yala A, Schuster T, et al. Mammographic breast density assessment using deep learning: clinical implementation. Radiology 2019;290:52-8.
- Liberman L, Menell JH. Breast imaging reporting and data system (BI-RADS). Radiol Clin North Am 2002;40:409-30.
- Sprague BL, Conant EF, Onega T, et al. Variation in mammographic breast density assessments among radiologists in clinical practice. Ann Intern Med 2016;165:457.
- Spayne MC, Gard CC, Skelly J, Miglioretti DL, Vacek PM, Geller BM. Reproducibility of BI-RADS breast density measures among community radiologists: a prospective cohort study. Breast J 2012;18:326-33.
- **10.** Brandt KR, Scott CG, Ma L, et al. Comparison of clinical and automated breast density measurements: implications for risk prediction and supplemental screening. Radiology 2016;279: 710-9.
- Youk JH, Gweon HM, Son EJ, Kim J-A. Automated volumetric breast density measurements in the era of the BI-RADS Fifth Edition: a comparison with visual assessment. AJR Am J Roentgenol 2016;206: 1056-62.
- **12.** Esteva A, Kuprel B, Novoa RA, et al. Dermatologist-level classification of skin cancer with deep neural networks. Nature 2017;542: 115-8.
- **13.** Chang K, Beers AL, Bai HX, et al. Automatic assessment of glioma burden: a deep learning algorithm for fully automated volumetric and bi-dimensional measurement. Neuro Oncol 2019;21:1412-22.

- Li MD, Chang K, Bearce B, et al. Siamese neural networks for continuous disease severity evaluation and change detection in medical imaging. NPJ Digit Med 2020;3:48.
- Pisano ED, Gatsonis C, Hendrick E, et al. Diagnostic performance of digital versus film mammography for breast-cancer screening. N Engl J Med 2005;353:1773-83.
- Beers A, Brown J, Chang K, et al. DeepNeuro: an open-source deep learning toolbox for neuroimaging Neuroinformatics, 2020, https:// doi.org/10.1007/s12021-020-09477-5
- Glorot X, Glorot X, Bengio Y. Understanding the difficulty of training deep feedforward neural networks. Proc Int Conf Artif Intell Stat (AISTATS'10) Soc Artif Intell Stat. May 13-15, 2010; Sardinia, Italy.
- He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. 2016. IEEE Conf. Comput. Vis. Pattern Recognit, June 27-30 2016; Las Vegas, Nevada.
- Huang G, Liu Z, van der Maaten L, et al. Densely connected convolutional networks. Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR), Honolulu, HI, 2017, pp. 2261-2269, https://doi.org/10.1109/CVPR.2017.243.
- 20. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z. Rethinking the inception architecture for computer vision. Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit., June 27-30 2016, Las Vegas, Nevada
- Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. Int Conf Learn Represent, May 7-9 2015; San Diego, California.
- 22. Russakovsky O, Deng J, Su H, et al. ImageNet large scale visual recognition challenge. Int J Comput Vis 2015;115:211-52.
- 23. Cheng J, Wang Z, Pollastri G. A neural network approach to ordinal regression. Proc. Int. Jt. Conf. Neural Networks 2008. https://doi.org/ 10.1109/IJCNN.2008.4633963. Sept 3-6, 2008: Prague, Czech Republic.
- 24. Dietterich TG. Ensemble methods in machine learning. Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), 2000;1857:1-15. https://doi.org/10.1007/3-540-45014-9\_1. June 21-23, 2000; Cagliari, Italy.
- McInnes L, Healy J. UMAP: Uniform manifold approximation and projection for dimension reduction arXiv 2018.
- Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics 1977. https://doi.org/10.2307/2529310.
- 27. Jia Deng, Wei Dong, Socher R, Li-Jia Li, Kai Li, Li Fei-Fei. ImageNet: A large-scale hierarchical image database. 2009. IEEE Conf. Comput. Vis. Pattern Recognit, 2009. https://doi.org/10.1001/jamaophthalmol. 2018.1934. June 20-25, 2009; Miami, Florida.
- Dunnmon JA, Yi D, Langlotz CP, Ré C, Rubin DL, Lungren MP. Assessment of convolutional neural networks for automated classification of chest radiographs. Radiology 2019;290:537-44.
- 29. Gulshan V, Peng L, Coram M, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. JAMA 2016;304:649-56.
- Chang K, Balachandar N, Lam C, et al. Distributed deep learning networks among institutions for medical imaging. J Am Med Informatics Assoc 2018;25:945-54.

- Brown JM, Campbell JP, Beers A, et al. Automated diagnosis of plus disease in retinopathy of prematurity using deep convolutional neural networks. JAMA Ophthalmol 2018;136(7):803-10.
- 32. Van Hulse J, Khoshgoftaar TM, Napolitano A. Experimental perspectives on learning from imbalanced data. Proceedings of the 24th international conference on machine learning 2007;935-42.
- 33. Zech JR, Badgeley MA, Liu M, Costa AB, Titano JJ, Oermann EK. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. PLoS Med 2018;15:e1002683.
- 34. Albadawy EA, Saha A, Mazurowski MA. Deep learning for segmentation of brain tumors: impact of cross-institutional training and testing: impact. Med Phys 2018;45:1150-8.
- **35.** Wu N, Phang J, Park J, et al. Deep neural networks improve radiologists' performance in breast cancer screening IEEE Trans Med Imaging 2020;39:1184-94.
- **36.** Mohamed AA, Berg WA, Peng H, Luo Y, Jankowitz RC, Wu S. A deep learning method for classifying mammographic breast density categories. Med Phys 2018;45:314-21.
- 37. Keavey E, Phelan N, O'Connell AM, et al. Comparison of the clinical performance of three digital mammography systems in a breast cancer screening programme. Br J Radiol 2012;85:1123-7.
- 38. Shih G, Wu CC, Halabi SS, et al. Augmenting the National Institutes of Health chest radiograph dataset with expert annotations of possible pneumonia. Radiol Artif Intell 2019;1:e180041. https://doi.org/1 0.1148/ryai.2019180041.
- Halabi SS, Prevedello LM, Kalpathy-Cramer J, et al. The RSNA pediatric bone age machine learning challenge. Radiology 2019;290:498-503.
- 40. Flanders AE, Prevedello LM, Shih G, et al. Construction of a machine learning dataset through collaboration: the RSNA 2019 brain CT hemorrhage challenge. Radiol Artif Intell 2020;2: e190211.
- Filice RW, Stein A, Wu CC, et al. Crowdsourcing pneumothorax annotations using machine learning annotations on the NIH chest Xray dataset. J Digit Imaging 2019;33:1-7.
- Su H, Deng J, Fei-Fei L. Crowdsourcing annotations for visual object detection. AAAI Work. - Tech. Rep, 2012. July 22-23, Toronto, Canada.
- **43.** Irshad H, Montaser-Kouhsari L, Waltz G, et al. Crowdsourcing image annotation for nucleus detection and segmentation in computational pathology: evaluating experts, automated methods, and the crowd. Pacific Symp Biocomput 2015:294-305. January 4-8 2015, Waimea, HI, USA.
- **44.** Candido dos Reis FJ, Lynn S, Ali HR, et al. Crowdsourcing the general public for large scale molecular pathology studies in cancer. EBio-Medicine 2015;2:681-9.
- **45.** He K, Girshick R, Dollár P. Rethinking ImageNet Pre-training arXiv 2018.
- 46. Isensee F, Petersen J, Klein A, et al. nnU-Net: self-adapting framework for U-net-based medical image segmentation. Inform. aktuell 2019. https://doi.org/10.1007/978-3-658-25326-4\_7. Image processing for medicine 2019, March 17-19 2019, Lübeck, Germany.