# Large-scale mammography CAD with Deformable Conv-Nets

Stephen Morrell[1] stephen.morrell@gmail.com, Zbigniew Wojna[1], Can Son Khoo[1], Sebastien Ourselin[2], and Juan Eugenio Iglesias[1]

[1] Medical Physics and Biomedical Engineering, University College London, UK
[2] Sch. Biomed. Eng. Im. Sci., King's College London

**Abstract.** State-of-the-art deep learning methods for image processing are evolving into increasingly complex meta-architectures with a growing number of modules. Among them, region-based fully convolutional networks (R-FCN) and deformable convolutional nets (DCN) can improve CAD for mammography: R-FCN optimizes for speed and low consumption of memory, which is crucial for processing the high resolutions of to 50 μm used by radiologists. Deformable convolution and pooling can model a wide range of mammographic findings of different morphology and scales, thanks to their versatility. In this study, we present a neural net architecture based on R-FCN / DCN, that we have adapted from the natural image domain to suit mammograms – particularly their larger image size – without compromising resolution. We trained the network on a large, recently released dataset (Optimam) including 6,500 cancerous mammograms. By combining our modern architecture with such a rich dataset, we achieved an area under the ROC curve of 0.879 for breast-wise detection in the DREAMS challenge (130,000 withheld images), which surpassed all other submissions in the competitive phase.

## 1 Introduction

Breast cancer is the most commonly diagnosed cancer and the second leading cause of cancer death in U.S. women [1]. Timely and accurate diagnosis is of paramount importance since prognosis is improved by early detection and treatment, notably before metastasis has occurred. Screening asymptomatic women with mammography reduces disease specific mortality by between 20% and 40% [1] but incorrect diagnosis remains problematic. Radiologists achieve an area under the ROC curve (AUC) between 0.84 and 0.88 [2], depending on expertise and use of computer aided detection (CAD).

CAD for mammography was first approved 20 years ago but some studies showed it to be ineffective [3] or counterproductive [2] because of over-reliance. Early CAD methods used simple handcrafted features and produced many false positive detections [2]. The best of these "classical", feature-engineered methods, represented by e.g., [4,5], plateaued at 90% sensitivity for masses at one false positive per image [5], and at 84% area of overlap in segmentation [4].

Deep learning (DL) has enhanced image recognition tasks, building on GPUs, larger data sets and new algorithms. Convolutional neural nets (CNNs) have been applied to mammography, outperforming classical methods. For example, Dhungel et al. [6] used CNNs to achieve state of the art results in mass classification. In a recent study, Kooi et al. [7] proposed a two-stage system in which a random forest classifier first generated proposals for suspicious image patches, and a CNN then classified such patches into malignant or normal groups. Kooi's system was trained on a large private dataset of 40,506 images (6,729 cases) of which 634 were cancerous, and achieved an AUC of 0.941 in patch classification – representing significant improvement beyond prior work.

Despite the adoption of CNNs and increasing size of datasets, mammographic analysis still lags work on natural images in dataset size and algorithm comparability. Databases like ImageNet [8] and MS COCO [9] include millions of instances. Moreover, these public datasets enable independent verification of algorithmic performance with private test sets. Very recently, the Optimam [10] and Group Health datasets have begun to approach ImageNet and MS Coco sizes. Group Health was made available under the DREAMS Digital Mammography Challenge [11] (henceforth "the Challenge"), which used a verified hidden test set to benchmark comparisons between methods. The Challenge had 1,300 participants, and was supported by the FDA and IBM among others.

Moreover, CNN architectures now include a plethora of new techniques for detection, classification and segmentation [12,13,14]. It has also been shown that integration of these tasks in unified architectures – rather than pipelining networks – not only enables more efficient end-to-end training, but also achieves higher performance than when tasks are performed independently (e.g., [14]). This is due to sharing of features that use richer locality information in the labels – segmentations or bounding boxes.

Here we present our submission to the second phase ("collaborative") of the Challenge. Our contribution includes the selection of architectures from the natural image domain, adaptation to mammography to balance the trade-off between high resolution and network size (computational tractability) for fine feature detection, e.g., microcalcifications, data augmentation and score aggregation. In particular, the presented system is – to the best of our knowledge – the highest resolution DL mammography object detection system ever trained.

## 2   Methods

### 2.1   Network architecture

Even though our objective is classifying whole images, we chose a detection architecture to exploit the rich bounding box information in our training dataset (Optimam), and also to increase the interpreteability of the results, which is useful for clinicians. Our choice of meta-architecture is Region-based Fully Convolutional Networks (R-FCN) [15], which are more memory-efficient than the popular Faster Region-based Convolutional Neural Nets (F-RCNN) [16]. R-FCNs were
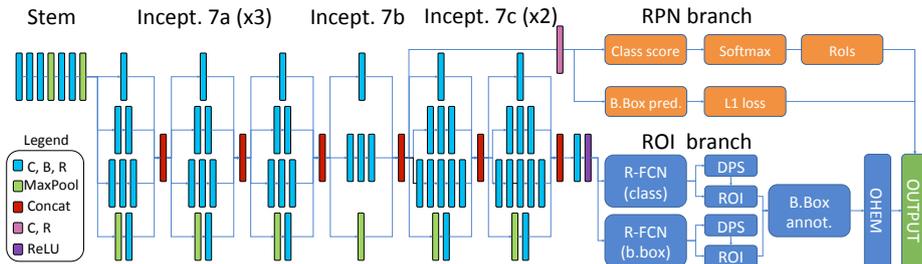
**Fig. 1.** Network architecture used in this study. C, B and R stand for convolution, batch norm and ReLU layers; RPN for region proposal network; DPS (ROI) for deformable position sensitive ROI; and OHEM [17] for online hard example mining.

enhanced with Deformable Convolutional Networks (DCN) [12], which dynamically model spatial transformations for convolutions and Regions of Interest (ROI) Pooling, depending on the data's current features:

$$y(\boldsymbol{p}_0) = \sum_{\boldsymbol{p}_n \in \mathcal{R}} \boldsymbol{w}(\boldsymbol{p}_n) \cdot \boldsymbol{x}(\boldsymbol{p}_0 + \boldsymbol{p}_n + \Delta\boldsymbol{p}_n),$$

where $y$ is the filter response at a location $\boldsymbol{p}_0$; $\mathcal{R}$ is a neighborhood around $\boldsymbol{p}_0$; and $\boldsymbol{w}$ and $\Delta\boldsymbol{p}_n$ are learnable sets of weights and offsets, respectively. The versatility of adaptive convolution and pooling enables DCNs to model a wider spectrum of shapes and scales, which is appropriate in mammography – where features of interest can be of very different sizes (from barely perceptible microcalcifications to large masses) and forms (foci, asymmetries, architectural distortions).

A diagram of our architecture is shown in Fig.1. It starts with a detection backbone, followed by two parallel branches: a region proposal network (RPN) branch, and a region of interest (ROI) branch – as per the R-FCN meta-architecture. The RPN branch [16] proposes candidate ROIs, which are applied on the score maps from the Inception 7b module. The ROI branch uses deformable position sensitive (DPS) score maps to generate class probabilities. Analogously to deformable convolutions, deformable ROI pooling modules include a similar parallel branch (Fig. 4 in [12]), in order to compute the offsets. Deformable pooling can directly replace its plain equivalent and can be trained with back propagation.

## 2.2 Adaptations

**Backbone:** Our backbone is descended from Inception v3 [18] from which we selected the first 7 layers (the "stem") and modules 7A, 7B and 7C. Choosing Inception, for which pre-trained weights from natural images were published, allowed transfer learning which showed beneficial for mammographic image analysis [19]. We included early layers on the assumption these are more consistent between domains. We chose consecutive layers to preserve co-adaption of weighs

where possible. We compared to other recent architectures [20,21,22] in a pilot dataset but results were weaker; we did not pursue them.

**Resolution-related trade-offs:** Current GPU memory constraints preclude full size mammographic images in deep CNNs – yet radiologists regularly zoom in to the highest level of detail. In particular, malignant microcalcifications may only be discerned at ∼50 μm resolution (approximately $4,000 \times 5,000$ pixels). Leading CNNs are designed for maximal GPU memory usage when fed natural images, which have two orders of magnitude fewer pixels, so the CNNs must be trimmed judiciously for mammography. This results in a trade-off between backbone choice, module selection, image downsampling, batch size, and number of channels in the different layers. We selected Inception v3 for its superior trade-off between parameter parsimoniousness and accuracy in natural images [23]. Its successor, Inception ResNet v2, was used in [12] but would have restricted images to $1,300 \times 1,300$ pixels. We included fewer repeats of modules 7A, 7B and 7C: three, one and two repeats, respectively; pilot experiments showed fewer layers to be sufficient for mammograms, whose content is less heterogeneous than natural images. We reduced batch size to one image per GPU. The channels are co-adapted in the pretrained weights and we ultimately retained them all. These choices, combined with meta-architecture and framework choices, enabled input size of $2,545 \times 2,545$ pixels (i.e. minimum downsample factor of 0.42), the highest resolution used for mammography classification to the best of our knowledge.

**Data augmentation:** Training is more effective if additional augmented data are included. Each training image was rotated through $360°$ in $90°$ increments and flipped horizontally, thus included eight times per epoch. We opted for the benefit of using rotated images despite induced "anatomical" noise, i.e., implausible anatomy. We also used the same four rotations and flipping at inference. We did not use random noise or random crops (which are standard in the natural image domain), to avoid omitting lesions at the image edge.

**Aggregation of multiple views:** Screening exams usually consist of two views (cranio-caudal, CC, and medio-lateral oblique, MLO) of each breast, giving four images per exam. For each of these images we generated 8 predictions, when including augmentations. Each subjects's probability of malignancy lesion was calculated by computing the mean over views and augmentations for each laterality, and then taking the maximum over the two sides; other combination rules were explored but yielded inferior results (see results in Section 3.3).

## 3    Experiments and results

### 3.1    Data

Two datasets were released in 2016/17, which were substantially larger than prior digital mammography datasets. These are the Optimam [10] dataset, which we used in training, and Group Health (GH), used in testing. The ground truths were determined by biopsy. Using standardised data makes comparisons objective and reproducible, e.g., as for ImageNet in the natural image domain. We believe these new benchmarks will allow attribution to future architectures.

Group Health images are a representative sample of 640,000 screening mammogram images, approximately 0.5% cancerous, provided in the Challenge. All machines were Hologic and maximum image size was $3,300 \times 4,100$ pixels. The GH data were not downloadable, excepting a small pilot set of 500 images for prototyping. The data are kept on IBM's cloud and are not accessible directly. Challenge participants could only upload models, run inference on the cloud, and receive a score. While this hampers testing experiments, it preserves patient confidentiality and ensures veracity of results. We used two subsets of GH in our study: 1. GH-13K, a subset of approximately 13,000 images with cancer prevalence inflated artificially by a factor of four; and 2. GH-Validation, a representative subset with 130,000 images, which was used for final testing and ranking and which the organisers intend to keep open for future testing, providing a hitherto absent way to benchmark performance.

Optimam consists of 78,000 selected digital screening and symptomatic mammograms including approximately 7,500 findings with bounding boxes of which 6,500 were cancerous. It included preprocessed and magnification images. Mammography machines were mainly Hologic and GE, and maximum image size was $4,000 \times 5,000$ pixels. Most teams in phase 2 of the Challenge used Optimam.

## 3.2 Experimental setup

We trained our network on all Optimam images with findings. We trained our architecture on three different classes: negative, benign and malignant findings. For subsequent analyses, we used the score of the malignant class as prediction score. Pixel intensities were normalised across different manufacturers and devices using the corresponding lookup tables. We chose MXNet `http://mxnet.incubator.apache.org/`, for its memory-efficiency which surpasses most frameworks including TensorFlow. In terms of parameters, we used the same values as in the original publications describing the different DL modules, with two main differences: 1. We changed the scale of the input images, as explained in Section 2.2; and 2. In order to reflect the much lower risk of overlap and occlusion observed in mammography compared with natural images, we reduced the RPN positive overlap parameter from 0.7 to 0.5, and the proposal NMS threshold from 0.7 to 0.1. Training took 48h on two NVIDIA TitanX GPUs.

We chose AUC on breast or patient classification as measure of diagnostic accuracy. AUC is used frequently to estimate the diagnostic performance of both CAD and radiologists. Compared with metrics like specificity at sensitivity or partial AUC, which are usually applied at high sensitivity levels and may be used to evaluate screening radiologists, AUC measures diagnostic accuracy across all probability thresholds, so is a comprehensive metric. Use cases for DL algorithms may range form automated flagging of some positive cases for immediate follow up – where high precision at high confidence thresholds is key – to safely excluding normal mammograms – where high negative predictive value at low confidence thresholds prevails.

We conducted three sets of experiments. First, we tested a number of design choices on the pilot set, then second on the GH-13K dataset, changing one or

**Table 1.** Summary of GH-13K results. Legend for consistent variables: R - ResNet backbone; M - mean probability over augmentations and views, maximum over laterality; I - Inception backbone; N - No augmentation at inference. The scale is in pixels.

| Exp. # | Changed variable | AUC before | AUC after | Consistent variables |
|--------|------------------|------------|-----------|----------------------|
| 1 | Train and Test Scale: 2145 to 2545 | 0.8352 | 0.8227 | R, M, N |
| 2 | Train and Test Scale: 2145 to 2545 | 0.8595 | 0.8667 | M, I |
| 3 | Increased test image size: 2500 to 2900 | 0.8173 | 0.8143 | I |
| 4 | Increased test image size: 2900 to 3300 | 0.8143 | 0.8039 | I |
| 5 | Changed backbone: Resnet to Inception | 0.8352 | 0.8584 | M, N |
| 6 | Added augmentation at inference | 0.8584 | 0.8667 | M, I |
| 7 | Max over breast's images changed to mean | 0.8591 | 0.8667 | I |
| 8 | As for 8 above and inference scale: 2,454 to 2,545 | 0.8511 | 0.8667 | I |
| 9 | Adapted Inception ResNet v2. Image size of $1,600 \times 1,600$ | N/A | 0.7366 | M, I, N |

two key variables at a time while holding others constant. We tested backbone choices, scales, train and test augmentation and aggregation. Each evaluation on GH-13K took a day on the IBM cloud. Finally, we submitted our final model for testing on GH-Validation, which took approximately 8 days.

### 3.3  Results

Table 1 summarises results from our GH-13K experiments. Experiments 1-5 explored combinations of backbones and scales. In terms of backbone, empirical results showed that Inception was superior to ResNet [24] (1 vs. 2 and 5). In terms of scales, the main conclusion was that performance peaked when train and test inference was run at the higher scale (2, 8); however, accuracy dropped when testing size exceeded training (3, 4). Experiments 6-8 assessed the impact of augmentation (6), as well as answering the question of how to aggregate augmentation scores (7). In general, these experiments showed that: a) augmentation at inference does help; and b) within a single breast, taking the mean over views and augmented images outperformed the alternative of taking the maximum (we still use the maximum across lateralities). Row 9 shows a test AUC of 0.74 on GH-Validation from the competitive phase of the Challenge using a classification net (an adapted Inception ResNet v2), TensorFlow with the largest possible resolution under that setup (1,600×1,600, much smaller than our current model, thus leading to a 0.15 decrease in AUC).

Based on these results, we submitted our final architecture described in Section 2 for testing on the large GH-Validation dataset. In the first sub-challenge, which records the AUC by breast purely on imaging and blinded to demographic information, we achieved AUC=0.879 (standard deviation: 0.00914), see Fig. 2(e), which is 0.005 above the top AUC in the competitive phase of the Challenge. It was also the highest single-model AUC in the collaborative phase, 0.014 below an ensemble of detection models, and higher than all patch-based models. The second sub-challenge is on subject-wise AUC, with access to both im-

ages and demographics. Despite ignoring demographics, our architecture gave AUC=0.868, behind only the top score in the competitive phase (a patch-based curriculum-trained model) by 0.006. Twenty-five method descriptions from this phase are available at `synapse.org`, but details of the collaborative phase, including performance of patch-based models trained on Optimam, is embargoed pending publication by the Challenge. Fig. 2 shows sample outputs from GH.
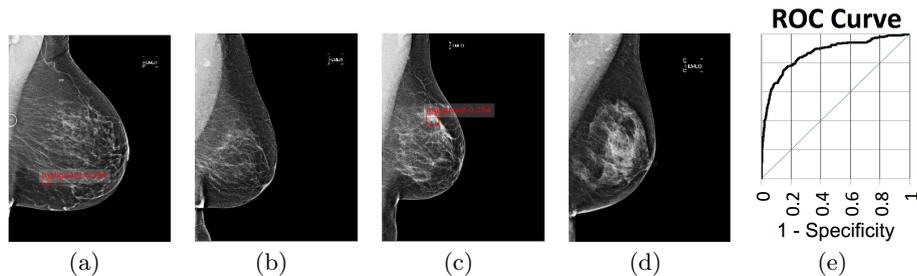


|    (a)    |    (b)    |    (c)    |    (d)    |    (e)    |

**Fig. 2.** (a) True positive prediction ($p = 0.90$ probability of malignancy) of an inconspicuous lesion on a left MLO of a 73 year old woman. (b) True negative (malignancy: $p = 0.06$) for left MLO view of a 66 year old woman. (c) False positive ($p = 0.78$) on left MLO of a 43 year old woman, due to hyper-intense region. (d) False negative ($p = 0.03$) for left CC view of a 61 year old woman. (e) ROC by breast, AUC = 0.879.

## 4 Discussion and conclusion

We have presented a two-stage detection network trained on strongly labelled data which achieved 0.879 AUC by breast on a large unseen representative screening test set, operating at high resolution. Important questions remain, e.g., the impact of deformable modules, ameliorating batch size = 1 in batch normalisation, image rotation, the use of architectures that handle multiple scales (FPNs), single pass models, comparison to the different setup in [7] and clinical applicability. Exploring these directions, along with integrating demographic features into the architecture, will be in a future journal extension. As mammographic training databases grow and the rapid progress in DL for machine vision continues, we hope this will provide a first benchmark in mammogram classification on an public yet hidden test set.

## References

1. American Cancer Society: What are the key statistics about breast cancer?
2. Lehman, C., Wellman, R., Buist, D., Kerlikowske, K., Tosteson, A., Miglioretti, D.: Diagnostic Accuracy of Digital Screening Mammography With and Without Computer-Aided Detection. JAMA Internal Medicine **175** (2015) 1828

3. Gross, C.P., Long, J.B., Ross, J.S., Abu-Khalaf, M.M., Wang, R., Killelea, B.K., Gold, H.T., Chagpar, A.B., Ma, X.: The Cost of Breast Cancer Screening in the Medicare Population. JAMA Internal Medicine **173**(3) (feb 2013) 220

4. Jiang, M., B, S.Z., Zheng, Y.: Mammographic Mass Segmentation with Online Learned Shape and Appearance Priors. MICCAI 2016 **9902** (2016) 35–43

5. Karssemeijer, N., te Brake, G.M.: Detection of stellate distortions in mammograms. IEEE Transactions on Medical Imaging **15**(5) (1996) 611–619

6. Dhungel, N., Carneiro, G., Bradley, A.P.: The Automated Learning of Deep Features for Breast Mass Classification from Mammograms. MICCAI 2016 **9902** (2016) 106–114

7. Kooi, T., Litjens, G., Ginneken, B.V., Gubern-mérida, A., Sánchez, C.I., Mann, R., Heeten, A.D., Karssemeijer, N.: Large scale deep learning for computer aided detection of mammographic lesions. Medical Image Analysis **35** (2017) 303–312

8. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. IJCV **115**(3) (2015) 211–252

9. Lin, T., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.: Microsoft COCO: Common objects in context. ECCV (2014) 740–755

10. Royal Surrey County Hospital: The Optimam Mammography Image Database

11. Bionetworks, S.: The Digital Mammography DREAM Challenge (2016)

12. Dai, J., Qi, H., Xiong, Y., Li, Y., Zhang, G., Hu, H., Wei, Y.: Deformable convolutional networks. In: CVPR. (2017) 764–773

13. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR. Volume 1. (2017) 4

14. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Computer Vision (ICCV), 2017 IEEE International Conference on, IEEE (2017) 2980–2988

15. Dai, J., Li, Y., He, K., Sun, J.: R-FCN: Object Detection via Region-based Fully Convolutional Networks. (may 2016)

16. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Advances in neural information processing systems. (2015) 91–99

17. Shrivastava, A., Gupta, A., Girshick, R.: Training region-based object detectors with online hard example mining. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 761–769

18. Szegedy, C., Vanhoucke, V., Shlens, J., Wojna, Z.: Rethinking the Inception Architecture for Computer Vision. (2015)

19. Carneiro, G., Nascimento, J., Bradley, A.P.: Unregistered multiview mammogram analysis with pre-trained deep learning models. MICCAI (2015)

20. Hu, J., Shen, L., Sun, G.: Squeeze-and-Excitation Networks. (2017)

21. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely Connected Convolutional Networks. arXiv preprint (2016) 1–12

22. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated Residual Transformations for Deep Neural Networks. (nov 2016)

23. Canziani, A., Paszke, A., Culurciello, E.: An Analysis of Deep Neural Network Models for Practical Applications. Arxiv (2016) 7

24. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. (2016) 770–778