

1 Supplementary Material

2 *Pre-processing phase*

3 Prior to processing each mammogram we initially checked one of its tag fields

4 (*APs.PhotometricInterpretation*) to determine whether the image's greyscale needs to be inverted:

```

5         if isfield(APs, 'PixelIntensityRelationshipSign')
6             if (APs.PixelIntensityRelationshipSign==1)
7                 Invert color;
8             end
9         end

```

10
11 The first proper step of our algorithm is to identify the “breast-to-air” region. Image segmentation has to
12 fulfil certain predefined criteria and is one of the most difficult tasks in image processing[29]. Similar
13 procedures that have been used for removing labels, markers and patient's particulars in digitised film
14 mammogram images are used here for reducing the size of FFDM images, by removing unnecessary black
15 background. This step is carried out by identifying the vertical critical column, j' , in the mammogram
16 image $\{I \in \mathbb{R}^2\}$ of size $m \times n$. We use a vertical projection technique, of the intensity values, which is
17 computationally undemanding yet very efficient in both FFDM and digitised film mammogram images. A
18 vertical projection can be seen as a plot in which each index, in b_j , corresponds to the sum of the pixels
19 along column j .

$$20 \quad b_j = \int_{-\infty}^{\infty} I(i, j) di \quad (S1)$$

21 This integral can be approximated by summation ($b_j = \sum_{i=1}^m I(i, j)$). Then, locating the global minimum in
22 the vector b yields the vertical critical column: $j' = \{j' / b_{j'} = \min(b)\}$, $j \in [1, 2, \dots, n]$. Similarly, the projection
23 method is adopted to locate the horizontal bottom cut in the image (i.e., where the breast ends) as shown
24 in Figure S1 (compare (a) to (b)). The image is then down-sampled to 512 pixels width using bilinear
25 interpolation while preserving the image aspect ratio for the height. Next, the signal-to-noise ratio is
26 improved by means of applying a global spatial contrast adjustment to the mammograms. Limits to

contrast stretch are found by considering the non-black region after performing image intensity normalization.

Detection of breast area

To choose a threshold algorithm for detecting the breast boundary (object blob) we first examined the following 13 different threshold algorithms, in a small subset of images: *concavity method; triangle algorithm; entropy method; inter-means method; iterative inter-means; maximum likelihood (EM); mean method; median method; minimum error; iterative minimum error; moments method; p-tile method and Otsu's method*. We choose to use the *triangle algorithm* due to its invariance to intra-breast intensity variations, making it our choice for detecting breast boundaries. Our algorithm then estimates the area of the objects in the resulting binary image and retains the largest object. Whilst the latter step, removal of non-breast objects, exhibits a mild effect in FFDM images, it is of crucial importance in digitised film mammogram images where there are labels, calibration wedges and various other paraphernalia on a mammogram. The position and orientation of these labels are not standard. Belkhodja et al.[30] use the Otsu method and a Gaussian filter in their algorithm to mask out labels in digitised film images.

Segmentation and pectoral muscle detection

The next step is to invoke an image segmentation routine. Image segmentation can be defined as the partitioning of a given image into non-overlapping, constituent regions that are homogeneous with respect to some characteristic such as intensity or texture [29] Ch.10. A survey of the medical image segmentation literature is provided by [31]. In this step we use Delaunay-based image segmentation,[32] , to estimate the number of pixels with homogenous intensity. This segmentation is fully automated, fast and does not require the user to provide an initial estimate of the number of clusters. The technique resembles the dynamic thresholding method used for segmentation, but it differs in terms of divide and merge-decision making. A Delaunay triangulation (*DT*), described in detail in [33], is constructed in our case from a set of points which correspond to the probability density function (PDF) of the image intensities. The outer

boundary of a DT is simply the convex hull (CV) of the set of the feature points where vertices provide a direct access to intensity values to be used in image segmentation. Time complexity ($O(n \log n), n_{\text{histogram_bins}} \leq 2^8$) is considerably lower than other approaches such as region growing and the deformable contours/snakes algorithms.

Let I be the entire mammogram image. Image segmentation can be seen as a process that partitions I into n sub-regions (I_1, I_2, \dots, I_n) based on a clear measure such as homogeneity in such a way that regions in the segmented image (S) comply to:

$$\text{i) } \bigcup_{i=1}^n I_i = I$$

$$\text{ii) } I_i \cap I_j = \emptyset, \forall i \neq j.$$

$$\text{iii) } I_i \text{ is a connected set, } i=1, 2 \dots n.$$

After this segmentation step, a pseudo-colour pattern (Red (R), Green (G), Blue (B)) is created using the composite of the original mammogram and its segmented version (S). We choose (arbitrarily) to represent the original image with the R component, in the RGB primary colour space. While the segmented image is treated as the G component in the native colour space, the blue B (black here) component is zeroed out to eliminate third component interference. Different perturbations can obviously be obtained. However, our choice is to have dense area (pectoral muscle and fibroglandular tissue) marked with a green gamut and the radiolucent tissue marked as red-orange; see the pseudo-colour pattern in Figure S1 (b). We examined two colour spaces namely, HSV and YIQ. It is interesting to note that I and Q components in the YIQ colour space have a very appealing property, namely that they provide an excellent separation of the red-orange colour from the green gamut. The HSV, on the other hand, is non-linear, therefore it is more computationally demanding than the YIQ transformation, especially in analyses of large volumes of medical data. The breast area segmentation step can be written in a simplified form based on the following signum function.

$$\text{sign}(Q(x, y)) = \begin{cases} -: Q(x, y) \in \text{fatty_region_foreground} \\ 0: Q(x, y) \in \text{image_background} \\ +: Q(x, y) \in \text{Dense_tissue_foreground} \end{cases} \quad (\text{S2})$$

See Figure S1 (c, e) for an example of the application of the signum function to real mammographic data. In order to remove all or most of the pectoral muscle, a mask needs to be created. Unlike existing methods where the pectoral muscle mask is constructed under the assumption that its boundary can be approximated with a single straight line on the top left corner for LMLOs, we propose the use of the convex hull allowing for fragmented straight lines. To create the mask, the smallest convex polygon (convex hull) is constructed from the binary image, $Q(x, y) < 0$ (see, Figure S1 (c-d)).

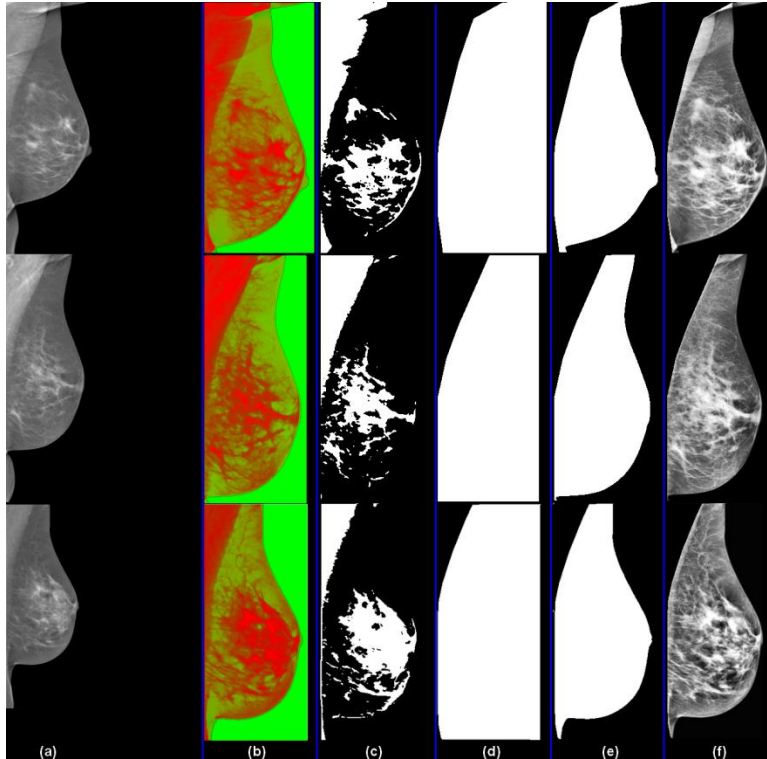


Figure.S1. Pre-processing of mammograms, (a) original mammograms, (b) pseudo-colour generation after applying the horizontal and vertical cropping, (c) the negative signal in the $Q(x, y)$ colour space, detecting the non-reddish area (d) convex hull of (c), (e) the final extracted breast mask, and (f) breast region after applying the contrast limited adaptive histogram equalization (CLAHE).

The accurate detection of breast and pectoral muscle boundaries is a challenging task since in some cases the latter's intensity blends together with the breast area where it becomes hard even for an observer to

identify the border of the pectoral muscle. We therefore wanted to compare our approach in this step to alternative pre-existing methods. A common approach is to use active contours (ACs), also known as Snakes or deformable contours [21] p. 299. Among those studies which used ACs in conjunction with mammography are [34,35], and the algorithm was recently evaluated in [36]. The ACs approach is based on an ingenious process that forms a set of points which aim to enclose a target shape, the shape to be extracted. The ACs approach works perfectly well in retrieving the true boundary of a fuzzy shape, thanks to its internal efficient energy minimization algorithm, making it an ideal option for semi-automatic processes. However, it fails to converge when the contour is tangent to the velocity vector and it can deform around a non-targeted object. For this reason the initialization phase can impact highly on the final result. Therefore, to enable a fair comparison with our developed algorithm, we assisted ACs to overcome these obstacles. We compared the result from our algorithm and the ACs approach in terms of the estimates of area PD to which they lead and against a third estimate (ground truth) of PD (obtained by the semi-automated approach Cumulus in a subset of 39 FFDM images. The Pearson correlation coefficients (correlation with Cumulus PD) were estimated as $r=0.90$ and $r=0.79$ for our PD proposed approach and the active contour algorithm, respectively. We measured time complexity on the larger dataset of 1011 processed FFDM images. The average time complexity was recorded as mean=0.2712 sec and mean=8.1614 sec for our PD proposed approach and the active contour algorithm, respectively. Both procedures are fast, although our approach can lead to remarkable time gain when considering a large volume of images as will be necessary for the KARMA cohort study. For instance, to process 190,000 images (the number of images which have currently been collected within KARMA) would take in excess of two weeks using the active contour algorithm as compared to merely 15 hours using our proposed approach.

Once the image of the breast has been segmented, we enhance the local contrast to reveal features present in the breast by applying contrast limited adaptive histogram equalisation (CLAHE) with the number of tiles being equal to $[8 \times 8]$. We have shown in previous studies pertaining to diabetes that CLAHE is a vital input for many imaging based medical oriented sciences; for further insights into its micro-biology

114 applications see[37,38]. CLAHE is described mathematically as follows. Let I be the image with the
 115 segmented breast, the CLAHE transform \bar{h} is given by:

$$116 \quad \bar{h}(I) = [I_{\max} - I_{\min}] \rho(f) + I_{\min} \quad (S3)$$

117 where $\rho(f)$ represents the CDF (cumulative distribution function) denoted by $\rho(f) = \rho(X \leq x) = \int_0^x f(t)dt$.

118
 119 Table S1: The derived statistical and textural features
 120

Feature (* denotes features derived from a single region, the rest are derived from the 12 regions)	Intensity based	Morphology/Shape based
PD: Percent density (CASAM-Area)		•
$f1$: The DC value of the 2D discrete cosine transform	•	
$f2$: The coefficient of the last decomposition of single-level discrete 2D wavelet transform	•	
$f3$: The max coefficient of the 2D discrete Fourier transform of the local range filter	•	
$f4$: Kurtosis of the region of interest (ROI) intensity	•	
$f5$: Skewness of the ROI intensity distribution	•	
$f6$: Entropy of the ROI intensity	•	
$f7$: Mean of the entropy filtered ROI	•	
$f8$: Entropy of the complex imaginary part of the convolved region with log-Gabor filters	•	
$f9$: Entropy of the complex real part of the convolved region with log-Gabor filters	•	
$f10$: Entropy of the magnitude part of the convolved region with log-Gabor filters	•	
$f11$: The max coefficient of the 2D discrete Fourier transform applied to the ROI	•	
$f12$: Entropy of the normalised co-occurrence matrix of the ROI	•	
$f13$: The magnitude of the difference mean of the region and the mean of the complement ROI	•	
$f14$: The max coefficient of the 2D discrete Fourier transform of the Hessian filter applied to the ROI	•	
$f15$: Entropy of the entire breast area*	•	
$f16$: The max coefficient of the log of the magnitude part of the discrete cosine transform applied to the ROI	•	
$f17$: The max coefficient of the log of the magnitude part of the discrete Fourier transform applied to the ROI	•	
$f18$: The 4th central moment of the ROI		•
$f19$: Number of particles within the ROI		•
$f20$: Solidity of the ROI		•
$f21$: Eccentricity of the ROI		•
$f22$: Euler Number of the ROI		•
$f23$: Number of particles within the entire breast*		•
$f24$: (Skewness of the normalised singular value decomposition of the ROI) / (standard deviation of the ROI)	•	
$f25$: Singular value decomposition of the ROI	•	
$f26$: Euler number within the breast (binary image)	•	
$f27$: Interquartile range- the difference between the 75 th and the 25 th percentiles of the intensity values of the ROI	•	
$f28$: The 1 st Fractal descriptor of the entire breast*	•	
$f29$: Mean intensity of the breast*	•	
$f30$: Skewness of the eroded ROI intensity	•	
$f31$: Skewness of the intensity distribution of the breast*	•	
$f32$: Kurtosis of the projection along the Y axis of the ROI		•
$f33$: Kurtosis of the projection along the X axis of the ROI		•
$f34$: Mean perimeter of the ROI		•
$f35$: Mean circularity of smallest particles 1-100 pix*		•
$f36$: Mean intensity of the ROI	•	
$f37$: Homogeneity of the breast area*	•	
$f38$: Median intensity of the ROI	•	

<i>f39</i> :Skewness of the gradient of the intensity within the ROI	•	
<i>f40</i> :Energy property of the co-occurrence matrix of the selected area - horizontal shift (2 pixels)	•	
<i>f41</i> :Energy property of the co-occurrence matrix of the selected area - diagonal shift (8 pixels)	•	
<i>f42</i> :Variance of the intensity within the breast area*	•	
<i>f43</i> :Mean of the local binary pattern (LBP) of the whole breast region*		•
<i>f44</i> :Kurtosis of the local LBP histogram of the whole breast region*		•
<i>f45</i> :Mean intensity of the pectoral muscle*	•	
<i>f46</i> :RunLengthCode: Short Run Emphasis (SRE), Long Run Emphasis(LRE), Gray Level Non-Uniformity (GLN), Run Percentage (RP), Run Length Non-Uniformity (RLN), Low Gray Level Run Emphasis (LGRE), High Gray Level Run Emphasis (HGRE)*	•	
<i>f47</i> :Entropy of the sum of the average values of co-occurrence matrix of the breast along the horizontal and diagonal directions *	•	
<i>f48</i> :Standard deviation of FD (fractal descriptor) for the selected area *	•	
<i>F49</i> : Lacunarity of FD for the selected area *	•	
<i>f50</i> : Laws Texture Energy Measures: Entropy Level, Entropy Edge, Entropy Spot, Entropy Ripple *	•	
<i>f51</i> : Microcalcifications *		
<i>f52</i> : Intensity ratio(dense/breast) *	•	
<i>f53</i> : Dense area size *		•
<i>f54</i> : Breast area *		•
<i>f55</i> : Percentage Fatty *		•

Table S2. The table depicts the twelve different regions used in our approach from which features in Table S1 are derived.

Region Label	Description
Region 1	Finds the <i>dense</i> region in mammographic images using the <i>p-tile method</i>
Region 2	Finds the <i>dense</i> region in mammographic images using the <i>entropy of the histogram</i> .
Region 3 (*)	Finds the <i>dense</i> region in mammographic images using the <i>moment preserving thresholding method</i>
Region 4	Finds the <i>dense</i> region in mammographic images using the <i>Otsu's method</i> .
Region 5	Finds the <i>dense</i> region in mammographic images using the <i>Iterative Otsu's method</i> .
Region 6	Finds the <i>dense</i> region in mammographic images using the <i>arithmetic mean of pixels</i> .
Region 7 (~)	Finds the <i>fatty</i> region in mammographic images using our customized threshold (see Fig.S1(c) & Eq.S2)
Region 8 (~)	Finds the <i>fatty</i> region in mammographic images using the <i>entropy of the histogram</i> .
Region 9 (~)	Finds the <i>fatty</i> region in mammographic images using the <i>p-tile method</i>
Region 10 (~)	Finds the <i>fatty</i> region in mammographic images using the <i>Otsu's method</i> .
Region 11 (~)	Finds the <i>fatty</i> region in mammographic images using the <i>moment preserving thresholding method</i>
Region 12	Finds the <i>dense</i> region in mammographic images using our customized threshold (see Fig.S1(c) & Eq.S2)

(~) Denotes the logical “NOT” of a binary matrix. (*) CASAM-Area is calculated based on Region 6 in Eq.1 (in the main manuscript).

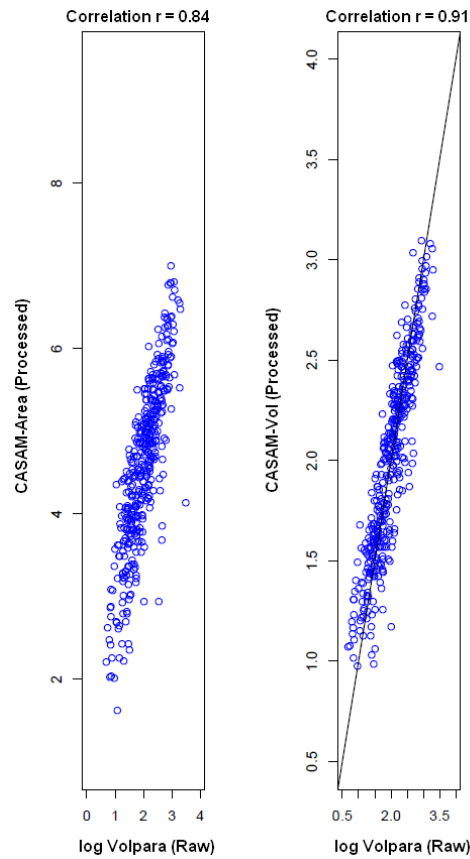


Figure.S2. Scatter plot showing the correlation between log (VolparaTM) on raw images and CASAM-Area and the predicted CASAM-Vol on the processed images for the training data set (403).

Table S3: p-values assessing the (residual) association between automated measures of mammographic density and case-control status (after adjustment for one other density measure). (a) with partial adjustment (age and BMI), (b) with full adjustment (age, BMI, menopausal status, HRT use, parity and age at first birth).

(a)		Variables adjusted for		
Covariates (n=1059)	Standard	Standard+ Volpara TM	Standard+ CASAM-Area (Processed)	Standard+ CASAM- Vol
Volpara TM	0.005	-	0.479	0.823
CASAM-Area (Processed)	0.012	0.163	-	1.000
CASAM-Vol (Processed)	0.023	0.088	0.235	-
(b)				
Volpara TM	0.010	-	0.527	0.560
CASAM-Area (Processed)	0.023	0.173	-	0.777
CASAM-Vol (Processed)	0.065	0.054	0.158	-