

# 1    **Breast Mass Segmentation based on Multi-scale Morphological**

## 2    **Sifting and Self-grown Cascaded Random Forests**

3    Hang Min, Shekhar S. Chandra, Stuart Crozier, Andrew P. Bradley

4    School of Information Technology and Electrical Engineering, The University of Queensland,  
5    St Lucia, QLD 4072, Australia

### 6    **Abstract**

7    In this paper, we present a mammographic breast mass detection and segmentation system.  
8    Breast mass detection and segmentation are challenging due to the fact that breast masses  
9    vary in size, shape, margin and contrast. Building classifiers for this task is also difficult since  
10   the normal tissue regions often greatly outnumber the abnormal regions, causing a class  
11   imbalance in the training set. To face the first challenge, we develop a multi-scale  
12   morphological sifting (MMS) approach using linear structuring elements mapped in radial  
13   patterns to extract the mass-like objects in mammograms. The MMS paired with a multi-level  
14   Otsu thresholding method can segment masses in various shapes, margins and sizes from the  
15   background tissue accurately. To tackle the class imbalance problem, we design a self-grown  
16   cascaded random forests (CasRFs) which adaptively assembles multiple layers of random  
17   forests (RFs) and adopts a probability-ranking based under-sampling method to balance the  
18   training set for each RF. Evaluated on two public available datasets, full-digital dataset  
19   INbreast and screen-film dataset DoD BCRP, the proposed method achieves an average  
20   sensitivity of 0.93 and 0.81 at 1.64 and 3.59 false positives per image respectively. The  
21   average Dice similarity index between the segmentation and the ground truth is calculated on  
22   INbreast and was found to be 0.85.

23    **Keywords**

24    Breast masses, Mammograms, Detection, Segmentation, Morphological sifting, Random  
25    forest

26    **1    Introduction**

27    Breast cancer is one of the leading causes of cancer death among females worldwide (Jemal  
28    et al., 2011). Early detection and diagnosis can increase the chances of survival and provides  
29    patients with more treatment options (Ganesan et al., 2013). Mammography is the primary  
30    modality in breast screening that has been proven to be effective in reducing breast cancer  
31    mortality (Andreea et al., 2011). The appearance of breast masses in mammograms is one of  
32    the most important signs of development of breast cancer (Schnabel et al., 2013). The  
33    detection of breast masses is generally regarded as more challenging compared with other  
34    breast abnormalities, such as calcifications, not only due to the large variation in size and  
35    shape of breast masses, but also because breast masses can appear in low contrast in  
36    mammograms (Oliver et al., 2010). With the development of image processing and machine  
37    learning technology, computer aided detection (CAD) has been introduced to breast image  
38    interpretation. Mammographic CAD has shown its potential in assisting radiologists to  
39    improve detection rate of breast cancer (Gromet, 2008).

40    Extensive studies have been done on mammographic CAD. Reviews of breast mass detection  
41    and segmentation can be found in articles (Oliver et al., 2010; Schnabel et al., 2013).

42    Generally, the design of breast CAD consists a number of stages. Firstly, region candidates  
43    are extracted from the mammograms, representing either abnormal or normal regions. To  
44    generate the region candidates, both unsupervised and supervised methods can be employed.  
45    Unsupervised methods normally involve the process of partitioning mammograms into  
46    multiple regions/segments by analysing the characteristics of pixels such as intensity, contrast

47 and topographic representations (Schnabel et al., 2013). To generate the region candidates,  
48 techniques such as region-grow (Görgel et al., 2013), thresholding (Kozegar et al., 2013;  
49 Varela et al., 2007), frequency domain filtering (Zhang et al., 2016), clustering (Martins et al.,  
50 2009) have been explored. These region candidate generation methods mostly aim at locating  
51 the suspicious regions. The regions generated normally include both the lesions and also  
52 many normal dense regions. The unsupervised methods tend to have difficulties in capturing  
53 characteristic lesions, such as spiculated lesions with highly irregular margins and  
54 architectural distortions (Schnabel et al., 2013). An alternative way of generating candidates  
55 is to use supervised methods. Instead of unsupervised segmentation, supervised methods use  
56 machine learning techniques such as convolutional neural networks (CNN) (Dhungel et al.,  
57 2015a) to study the characteristic features that describe mammographic masses, and generate  
58 a signal at the location of breast masses. The supervised methods may have the potential to  
59 detect the characteristic lesions. However, these types of methods also face relatively high  
60 computational complexity caused by the fact that the search should be performed  
61 exhaustively through the pixels in the image and even on different scales to avoid missing  
62 small lesions (Schnabel et al., 2013). Moreover, it may only provide the location of the  
63 masses (e.g. using bounding boxes) without accurately contouring the masses in the detection  
64 process (Dhungel et al., 2015a; Liu et al., 2015).

65 After the region candidates are generated, machine learning methods are normally used to  
66 determine whether the candidates are masses or normal tissue. Various machine learning  
67 algorithms have been applied to solve this problem, such as neural networks (Varela et al.,  
68 2007), support vector machines (Görgel et al., 2013; Martins et al., 2009), deep learning and  
69 random forest (Dhungel et al., 2015a). However, data imbalance is a common issue in this  
70 step, and high class imbalance normally leads to poor classification performance for the  
71 minority class (He and Ma, 2013). Since the minority class in this case are the positive

samples (masses) and the true positive rate (TPR or sensitivity) is critical to the system, a learning algorithm that can deal with the severe class imbalance and achieve a satisfactory sensitivity is required. Various methods have been developed to address the data imbalance problem. Sampling, either under-sampling the majority class, or over-sampling the minority class (Chawla et al., 2002), is a commonly used approach to balance the data. However, the sampled data may not be able to represent the distribution of the original data (Bria et al., 2014; Kang and Cho, 2006). Another solution is to introduce evaluation metrics to the data by assigning different weights to each class in proportion to its frequency in the training set (Viola and Jones, 2002). However, it can be difficult to design a robust evaluation metric to properly value the minority class (He and Ma, 2013). This type of method may also not perform as well as the sampling methods when facing an extremely imbalanced data set (Kang and Cho, 2006). To improve the sampling method, ensemble learning has been proposed to cope with severe class imbalance. Ensemble methods use a set of learning machines to make decisions and have been proven to be successful in solving class imbalance problem (He and Ma, 2013). Ensemble learning can be paired with random under-sampling methods such as RUSBoost (Seiffert et al., 2010), or oversampling methods such as SMOTEBoost (Chawla et al., 2003). To date, there are only a few publications that have addressed the data imbalance problem in breast mass detection (Chu et al., 2015; Kozegar et al., 2013; Murthy et al., 2013).

In this work, we propose a novel algorithm that can both detect and segment breast masses with a simple and straightforward structure. Figure 1 shows the diagram of the mammographic breast mass detection and segmentation system. The system mainly consists of a multi-scale morphological sifting approach and a self-grown cascaded random forests. Firstly, region candidates are generated by using multi-scale morphological sifting (MMS) and multi-level Otsu thresholding (MLT). To extract mass-like patterns, the MMS analyses

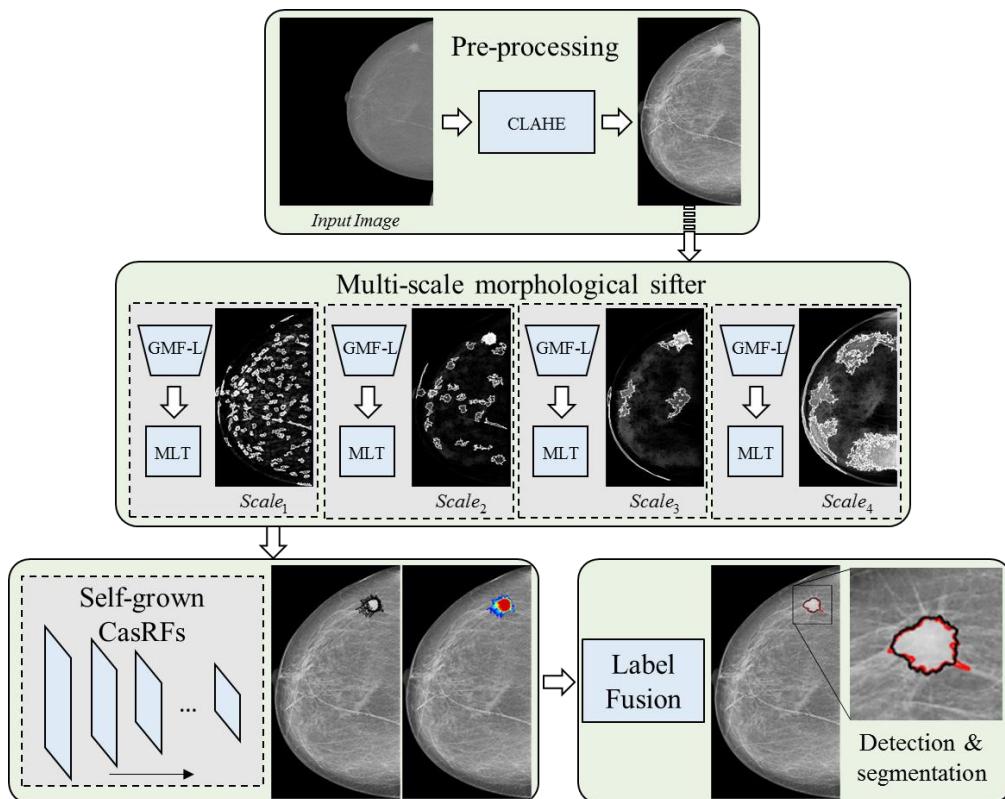
97 images by using two sets of linear structuring elements mapped in multiple orientations at  
98 each scale. The novelty of the region candidate segmentation method is that it creates a rather  
99 general morphological model of breast densities that captures the masses accurately  
100 regardless of their size and shape variations. Although morphological based methods have  
101 been used in several studies (Kharel et al., 2017; Kimori, 2011; Tan et al., 2015), however,  
102 they are mainly morphological operations used with other filters as pre-processing, while in  
103 our work the MMS plays the key role in region candidate segmentation.

104 The region candidate generation method normally produces many more normal tissue regions  
105 than mass regions, leading to a highly skewed training set. To cope with the severe data  
106 imbalance problem, a probability ranking based, adaptively self-grown cascaded random  
107 forests (CasRFs) is designed to classify the region candidates as abnormal or normal regions.

108 The CasRFs assemble a series of random forest (RF) classifiers in a cascade manner, and  
109 gradually discard the least suspicious region samples throughout the layers by using an  
110 under-sampling approach guided by probability ranking. Once the remaining negative  
111 samples can no longer be separated from the positive samples, the cascade ceases to grow.

112 Unlike many ensemble learning methods that adopted random under-sampling to attain a  
113 balanced training subset for each individual base classifier (Seiffert et al., 2010; Wei et al.,  
114 2015), the CasRFs under-sample the majority class by choosing the less lesion-like samples  
115 based on their probability to be a lesion, which ensures that the training subset at the current  
116 RF layer is highly separable. This enables the cascade to grow stably. The CasRFs are  
117 constructed in a simple and straightforward way without the need of involving complex  
118 weight matrixes and cost functions (Bria et al., 2016; Viola and Jones, 2002). The CasRFs are  
119 capable of dealing with highly skewed training sets and adaptive to the imbalance problem  
120 without the need of setting the number of layers, as in previous studies (Baumann et al., 2013;  
121 Tang et al., 2012). This work is an extension of our preliminary work (Min et al., 2017). We

122 have replaced the region candidate generation method with a novel multi-scale morphological  
 123 sifting algorithm, and adopted a different initialization approach for the self-grown CasRFs.  
 124 The new MMS method performs significantly better than the previous region candidate  
 125 generation method based on regular morphological structuring elements (Min et al., 2017).  
 126 The redeveloped self-grown CasRFs perform the training process on all training samples  
 127 from all scales rather than training at each individual scale (Min et al., 2017). This work also  
 128 presents extended analysis on the mechanism and performance of the CasRFs. Overall, the  
 129 proposed method combines mass detecting and segmenting functions in one relatively simple  
 130 structure that can be established on a regular desktop. Evaluated on two publicly available  
 131 mammographic datasets, the proposed method achieves competitive performance compared  
 132 with state-of-the-art methods that have much higher complexity (Dhungel et al., 2015a;  
 133 Dhungel et al., 2017).



134  
 135 Figure 1. The diagram of the mammographic breast mass detection and segmentation system using multi-scale  
 136 morphological sifting and self-grown cascaded random forests. CLAHE stands for contrast limited adaptive

137 histogram equalization. GMF-L stands for grayscale morphological filter using linear structuring elements.  
138 MLT stands for multi-level Otsu thresholding. The outlines of the detected patches are shown in black, and the  
139 ground truth is marked in red.

140 **2 Material and methods**

141 **2.1 Mammographic datasets**

142 The system is evaluated on two public mammographic datasets, INbreast (Moreira et al.,  
143 2012) and DoD BCRP (Bowyer et al., 1996). INbreast is a full-digital mammographic  
144 database containing 115 cases (410 images) with a  $70 \mu\text{m}$  pixel size and a 14-bit contrast  
145 resolution. There are 116 masses within a size range between  $15 \text{ mm}^2$  and  $3689 \text{ mm}^2$  in this  
146 dataset. INbreast is currently the largest publicly available, annotated full-field digital  
147 mammographic dataset (Dhungel et al., 2017). The lesion locations and boundaries are  
148 outlined by an image specialist, which enables evaluation of segmentation performance of the  
149 CAD system. DoD BCRP is a screen-film mammographic dataset whereby a training set of  
150 39 cases (80 annotated masses) and a testing set of 40 cases (81 annotated masses) are  
151 provided. The pixel size is  $43.5 \mu\text{m}$  and the contrast resolution is 12-bit. DoD BCRP contains  
152 many characteristic lesions such as spiculated, ill-defined masses and architectural distortions,  
153 which is good for testing the performance of the proposed method on highly irregular masses.  
154 INbreast is more clinical relevant compared with DoD BCRP since most screening  
155 programmes have adopted full-digital mammography. Therefore, INbreast is used as the  
156 primary evaluation dataset for our algorithm, and the performance on DoD BCRP is only  
157 presented as secondary results. We believe it is important to use public available datasets for  
158 evaluation to attain unbiased performance comparison between studies.

159 2.2 Pre-processing

160 To speed up the process, the mammograms are often resized to a lower resolution in pre-  
161 processing (Chu et al., 2015; Dhungel et al., 2015a; Eltonsy et al., 2007; Min et al., 2017).  
162 Here, we reduce the size of the mammograms by a factor of 4 using bi-cubic interpolation.  
163 The breast profile is pre-segmented using a simple threshold and then the redundant  
164 background is cropped away in INbreast. The redundant background regions in DoD BCRP  
165 are manually cropped away due to the fact that the mammograms from DoD BCRP contain  
166 various image noises and artefacts in the background. The pixel values in the image are  
167 linearly rescaled to 16-bit and the contrast limited adaptive histogram equalization (CLAHE)  
168 (Pizer et al., 1987) is then applied. The number of tiles in CLAHE is set as {4, 4}, which is  
169 arbitrarily chosen, and the contrast enhancement limit is set as 0.01 (default).

170 2.3 Region candidate generation

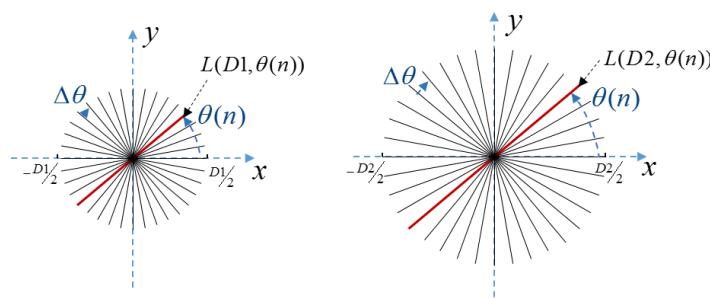
171 In order to generate the region candidates, a multi-scale morphological sifting (MMS)  
172 approach is applied to extract elements that are likely to be masses. Then multi-level Otsu  
173 thresholding (MLT) is employed to partition the image into meaningful patches that either  
174 represent mass regions or normal tissue regions.

175 2.3.1 Multi-scale morphological sifters

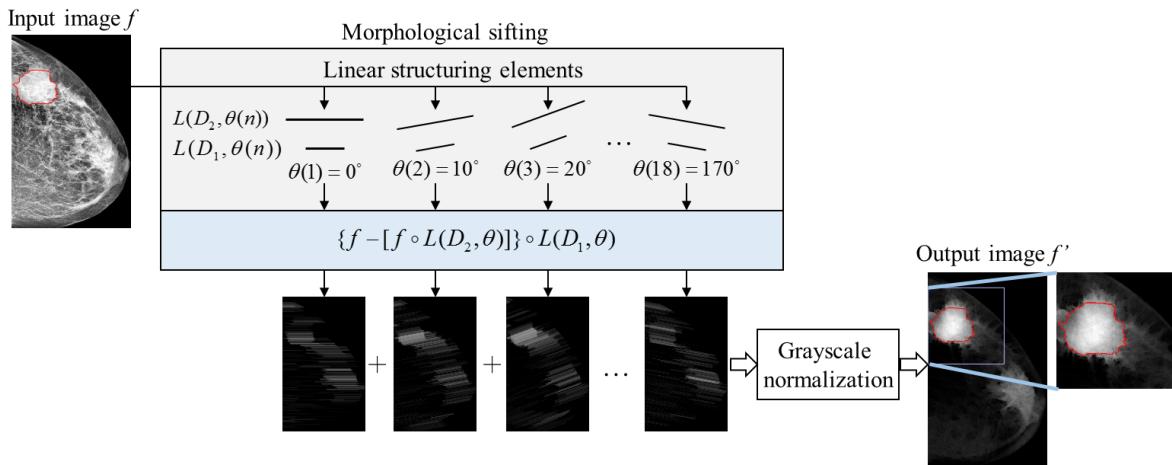
176 Morphological operations with structuring elements can be used to probe and analyse the  
177 images under the study for properties of interest (Gonzalez and Woods, 2008). Applying a  
178 dual-stage morphological filtering approach using two morphological operations (two top-hat  
179 operations or one top-hat followed by an opening) with different structuring elements can  
180 suppress the background tissue and extract the patterns of interest (Li et al., 2001; Wang,  
181 2006). By altering the size and shape of structuring elements, this dual-stage morphological  
182 approach can extract elements that fit the size and shape priors introduced by the structuring

183 elements. Here, we propose a new set of multi-scale grayscale morphological filters by using  
 184 linear structuring elements (MGMF-L) that cover a range of orientations, as malignant  
 185 mammographic densities are often surrounded by a radial pattern of linear spicules  
 186 (Karssemeijer and te Brake, 1996). Since the dual-stage morphological approach has the  
 187 ability to sieve elements of interest from the background, we name the approach  
 188 ‘morphological sifting (MS)’. The algorithm is described as follows.

189 Firstly, two sets of linear structuring elements are defined as  $\{L(D1, \theta(n)) | n = 0, 1, \dots, N-1\}$  and  
 190  $\{L(D2, \theta(n)) | n = 0, 1, \dots, N-1\}$ , where  $D1, D2$  stands for the magnitudes, and  $\theta(n)$  stands for the  
 191 orientation of each individual line.  $\theta(n)$  is equally spaced in  $[0^\circ, 180^\circ]$  as shown in Figure 2 (a).  
 192 There are  $N$  elements in each set and the orientation difference between two adjacent lines is  
 193  $\Delta\theta$  ( $\Delta\theta = 180^\circ / N$ ).



(a) Two sets of linear structuring elements. The lines marked in red represent a linear element pair  $[L(D1, \theta(n)), L(D2, \theta(n))]$ .



(b) The process of morphological sifting on one scale.

194  
 195 Figure 2. The process of morphological sifting using linear structuring element pairs in different orientations. (a)  
 196 shows how the linear elements are mapped in space. (b) shows how the structuring element pairs

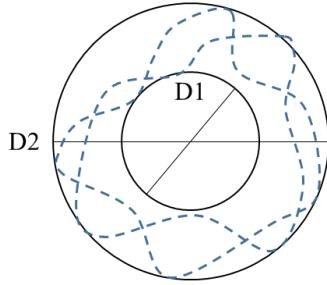
197  $\{[L(D1, \theta(n)), L(D2, \theta(n))] | n = 0, 1 \dots N-1\}$  are applied on the input image on each scale. This is an example on  
198 scale 4 and the lesion lands in scale 4 due to its size.

199 The grayscale morphological filter is described in the equation below.

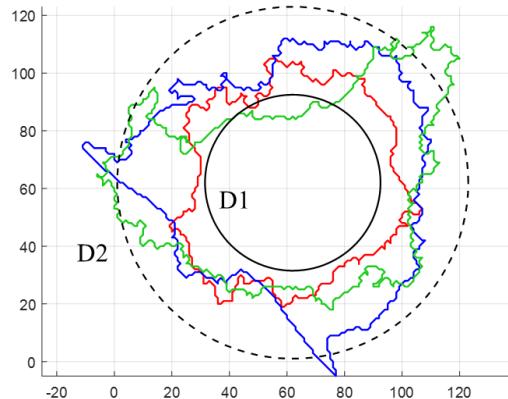
200

$$f' = \sum_{n=0}^{N-1} \{f - [f \circ L(D2, \theta(n))] \} \circ L(D1, \theta(n)) , \quad 1$$

201 where  $f$  stands for the input image. ‘ $\circ$ ’ stands for morphological opening.  $f'$  is the  
202 processed image, defined as the summation of the outputs of morphological sifters over all  
203 the orientations. Figure 2 (b) shows the process of morphological sifting. Two linear  
204 structuring elements  $[L(D1, \theta(n)), L(D2, \theta(n))]$ , with the same orientation, are paired up. These  
205 structuring element pairs are then applied onto the input image individually as shown in  
206 Figure 2 (b). The output image is generated by applying a grayscale normalization on the  
207 summation of all the result images generated by the morphological filtering. Here, we  
208 imagine that the target breast densities are composed of numerous straight ‘pen-strokes’ in  
209 different directions and sizes, just like these linear structuring elements. By extracting these  
210 ‘lines’ in different orientations, we are also extracting the target objects. At a certain scale,  
211 the MMS extracts objects whose outlines lie approximately in the space between the circles  
212 with a diameter of  $D1$  and  $D2$  as shown in Figure 3 (a). This is because that the ‘strokes’  
213 inside the target objects tend to survive the morphological operations done by structuring  
214 elements  $L(D1)$  and  $L(D2)$ .



215 (a) Theoretical model of objects that MMS filters extract.



215 (b) Real masses extracted by MMS filters on one scale (scale3).

216 Figure 3. The morphological model of the objects that MMS extracts.  $D1$  and  $D2$  here stand for the diameters  
217 of the two circles, as well as the magnitudes used in the linear structuring elements.

218 The morphological sifters are applied at multiple scales targeting lesions within different size  
219 ranges to cope with the size variation of masses. If the size range of masses is  $[Area_{\min}, Area_{\max}]$ ,  
220 the pixel size of the image is  $P$ , the number of scales used is  $M$ , and the resizing factor used  
221 in the pre-processing stage is  $R$ , we can roughly estimate the magnitude range of linear  
222 structuring elements as,

$$223 [DI_{\min}, DI_{\max}] = \left[ 2 \times \sqrt{\frac{Area_{\min}}{\pi}}, 2 \times \sqrt{\frac{Area_{\max}}{\pi}} \right]. \quad 2$$

224 Then, an ‘exponential’ morphological structuring element bank  
225  $\{L(D1(i)), L(D2(i)) | (i=1,2,\dots,M)\}$  is generated.  $M$  is the number of scales used, and the scale  
226 interval  $SI$  is

$$227 SI = \left( \frac{DI_{\max}}{DI_{\min}} \right)^{\frac{1}{M}}. \quad 3$$

228 The magnitudes of the linear structuring element pair at the  $i^{\text{th}}$  scale ( $Scale_i$ ) are

$$229 \{D1(i) = DI_{\min} \times SI^{i-1}, D2(i) = DI_{\min} \times SI^i | (i = 1, 2, \dots, M)\}. \quad 4$$

230 Therefore, at  $Scale_i$ , the MS filtered image  $f'(i)$  is

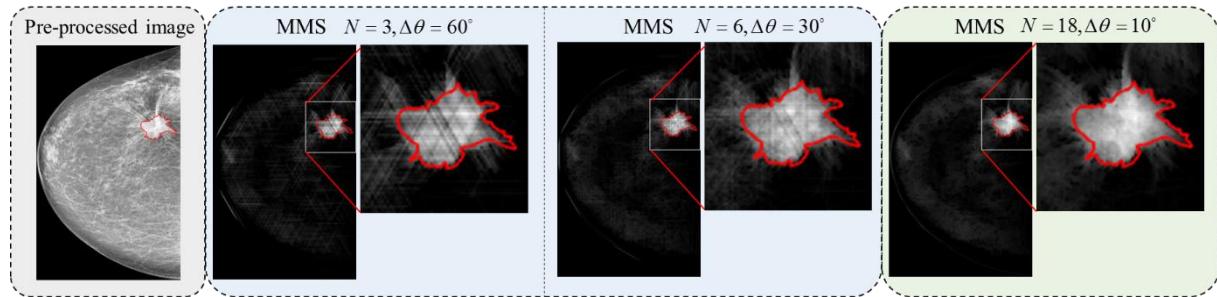
231

$$f'(i) = \sum_{n=0}^{N-1} \{f - [f \circ L(D2(i), \theta(n))] \} \circ L(D1(i), \theta(n)). \quad 5$$

232 Thus, the morphological sifter at  $Scale_i$  targets elements within the size range of

233  $[Area_{D1(i)}, Area_{D2(i)}] = [\pi/4 \times D1^2(i), \pi/4 \times D2^2(i)] . \quad 6$

234 Generally, as the number of linear structuring elements  $N$  increases, the smoother the  
 235 processed image becomes. However, this also comes with an increased computational  
 236 complexity. Figure 4 shows the examples of applying the morphological sifters for different  
 237  $N$ . Here,  $N$  is set as 18 ( $\Delta\theta = 10^\circ$ ) and  $f'$  is normalized into 16-bit in our case. The number  
 238 of scales is set as 4.

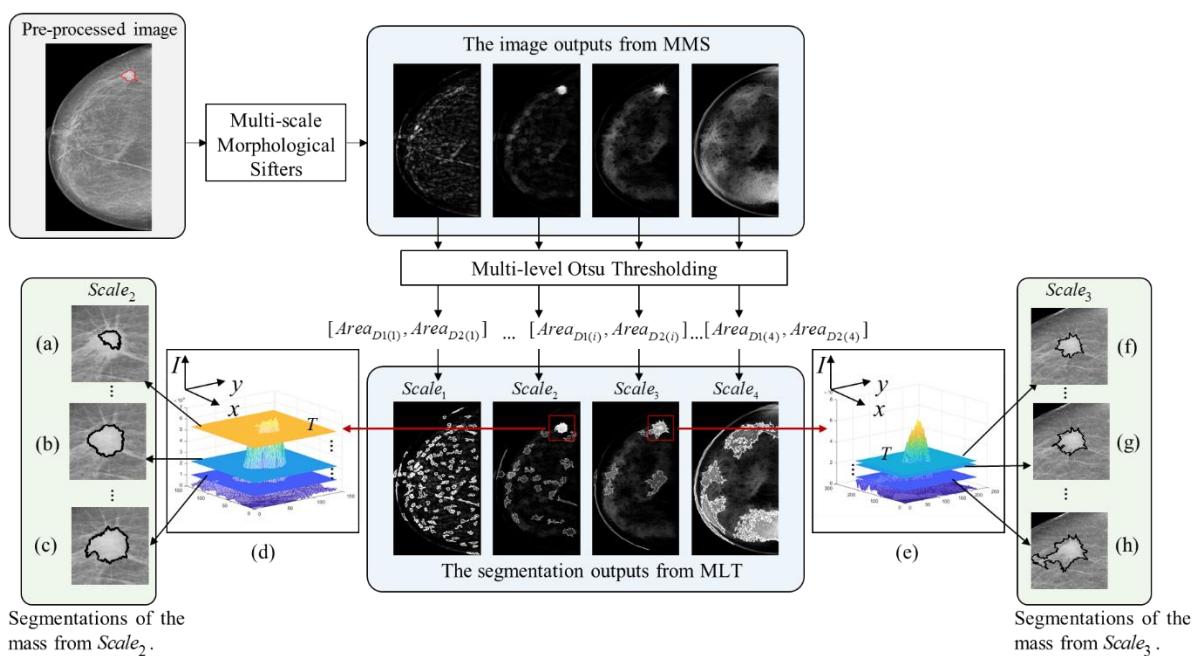


240 **Figure 4.** An example of the multi-scale morphological sifters (MMS) at one of the scales using  $N = 3$ ,  $N = 6$ ,  
 241 and  $N = 18$ . The red outlines denote the ground truth of the lesions.

### 242 2.3.2 Multi-level Otsu thresholding

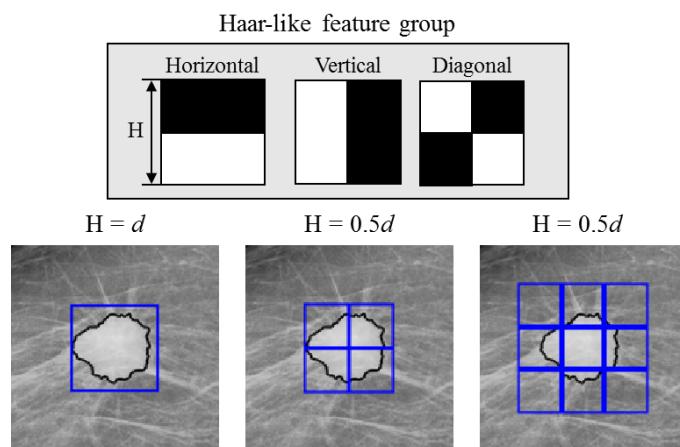
243 After multi-scale morphological sifting, the multi-level Otsu thresholding (MLT) method is  
 244 then applied on the processed images at each scale. Otsu's method is a classic thresholding  
 245 method (Otsu, 1979), which maximizes the between-class variance based on the histogram of  
 246 an image (Gonzalez and Woods, 2008). Multi-level Otsu is an extension of the original Otsu  
 247 method and is a reliable way to generate the optimal thresholds to segment an image (Backes  
 248 and Bruno, 2008; Liao et al., 2001). **Figure 5** shows the mechanism of region candidate  
 249 segmentation using MMS and MLT. The MLT generates a series of thresholds  $[T_k | k = 1, \dots, K]$ .  
 250 For each threshold  $T_k$ , if a pixel value in the filtered image  $f'(x, y) > T_k$ , then the pixel value in

251 the output binary image  $B(x, y) = 1$ , if  $f'(x, y) \leq T_k$ , then  $B(x, y) = 0$ .  $K$  thresholds should  
 252 generate  $K$  binary images representing the segmentations at each level of threshold. However,  
 253 in order to save space, at each scale, we plot the contours of the segmentations generated  
 254 from different thresholds together on the same image as shown in ‘The segmentation outputs  
 255 from MLT’ in Figure 5. Generally, the more scales and the more levels in multilevel Otsu,  
 256 the more likely for the algorithm to capture the accurate outlines of the lesions, but with a  
 257 higher computational complexity. Here, we use four scales and 16-level multilevel Otsu to  
 258 segment the images. Only the patches within the targeted size range (equation 6) at each scale  
 259 are selected as the region candidates as shown in Figure 5.



260  
 261 Figure 5. Region candidate segmentation based on multi-scale morphological sifting and multi-level Otsu  
 262 thresholding on 4 scales. (d) and (e) are the x-y-Intensity plots of the lesion area in  $Scale_2$  and  $Scale_3$ , where  
 263 the flat surfaces show how different thresholds slice through the intensity value and generate the segmentations.  
 264 (a) ~ (c) show several segmentation examples at the lesion area on  $Scale_2$ , while (f) ~ (h) show some examples  
 265 at the lesion area on  $Scale_3$ . Both  $Scale_2$  and  $Scale_3$  generate relatively satisfactory segmentations of the  
 266 lesion.

267 2.4 Region candidate classification using self-grow cascaded random forests  
 268 After the region candidates are generated, intensity features such as average intensity of  
 269 pixels within the candidate region, kurtosis, skewness, contrast (te Brake et al., 2000), shape  
 270 features such as eccentricity, extent, solidity, circularity (Cascio et al., 2006), radius, and  
 271 texture features such as entropy, grey-level co-occurrence matrix (GLCM), inertial  
 272 momentum (Cascio et al., 2006), are extracted. Haar-like features (Viola and Jones, 2001) are  
 273 also used, and they are defined as in Figure 6 below. Three types of Haar features, horizontal,  
 274 vertical and diagonal are operated at two sizes  $d$  and  $0.5d$ , where  $d$  is the longest distance  
 275 between the pixels on the boundary and the centre of the region.



276  
 277 Figure 6. Haar-like feature group.  $d$  is the longest distance between the pixels on the boundary of the region and  
 278 the centre of the region.

279 The previous stage generates many more non-lesion patches than mass patches. The ratio  
 280 between positive to negative samples is approximately 1:160 in the training set, which  
 281 indicates that the training set is highly imbalanced. Here, we design an ensemble learning  
 282 method, self-grown cascaded random forests (self-grown CasRFs), to deal with the data  
 283 imbalance problem. The random forest (RF) is chosen as the base classifier due to its  
 284 advantages that it has relatively high accuracy, efficiency and is less prone to overfitting  
 285 (Breiman, 2001).

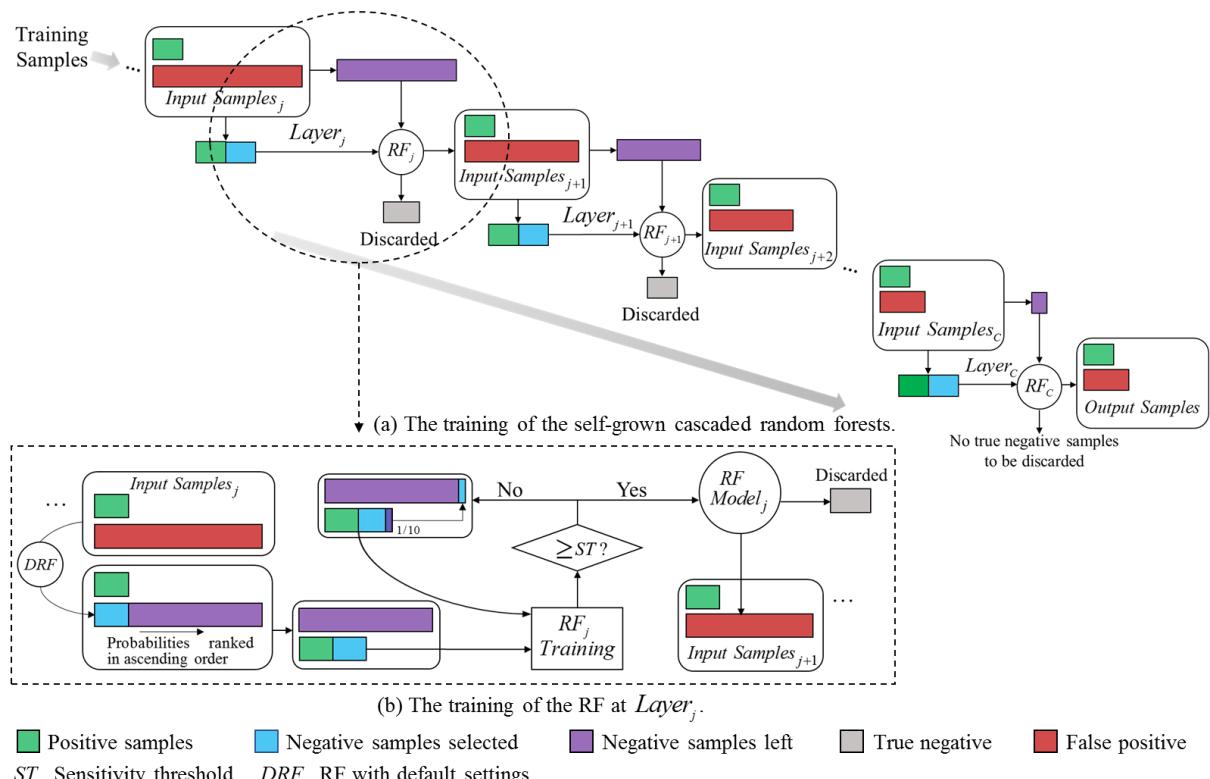
286 The growth of the multi-layered self-grown cascaded random forests is shown in Figure 7 (a).  
287 There are two main parameters for RF, number of trees used in the forest (ntree), and number  
288 of features used in each tree (mtry). At  $Layer_j$ , all the positive and current negative samples  
289 (all negative samples for  $Layer_i$ ) are put through a RF with default settings (DRF, ntree = 500,  
290  $mtry = \sqrt{\text{number of features}}$ ) as shown in Figure 7 (b). The negative samples are then sorted in  
291 ascending order of (posterior) probabilities generated from the DRF. To balance the two  
292 classes, if the number of positive samples is  $NP$ , the first  $NP$  negative samples are chosen to  
293 establish a balanced training set. Since the negative samples have been ranked in ascending  
294 order, the chosen negative samples are the least mass-like candidates and therefore relatively  
295 easy to distinguish from the positive candidates. Then a RF is trained on the balanced training  
296 set using a grid search for parameters ntree between [100,1000], and mtry between

297  $[\frac{1}{2}\sqrt{\text{number of features}}, 2\sqrt{\text{number of features}}]$ , to find a model that reaches a sensitivity threshold

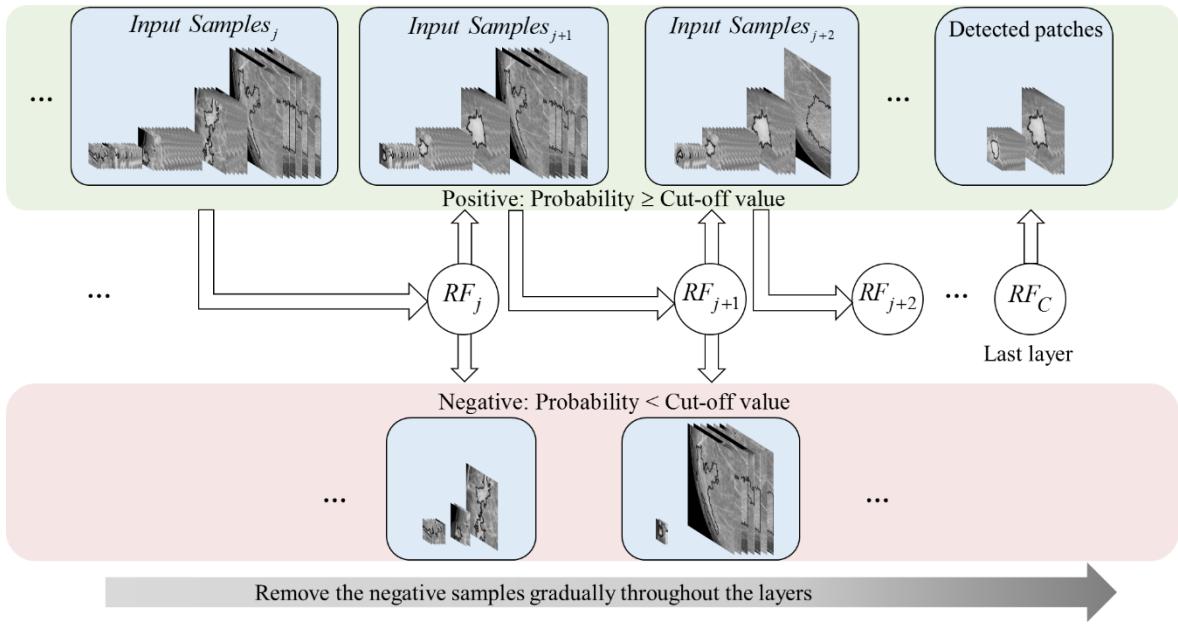
298  $ST$ . If no RF can reach  $ST$  in the current round of grid search, top 10% of the negative  
299 samples (with the highest probabilities, more likely to be lesions) are removed from the  
300 selected training set as shown in Figure 7 (b). After a RF model that meets  $ST$  is attained, the  
301 currently unused negative samples are classified by the RF model at  $Layer_j$  and the true  
302 negative (TN) samples are discarded from the training. The remaining FPs and all positive  
303 samples then become  $Input Samples_{j+1}$  and the learning moves on to  $Layer_{j+1}$ . As the cascade  
304 grows out more layers, the relatively easy non-lesion samples are gradually removed and the  
305 classification task becomes increasingly difficult. When the number of FPs stops to decrease,  
306 it indicates that the positive and negative samples can no longer be separated apart by the  
307 CasRFs, and the cascade stops to grow, as shown in Figure 7 (a). The formation of CasRFs is  
308 adaptive to the imbalance problem, requiring no need to manually choose the number of  
309 layers. The CasRFs are trained on the samples gathered from all the scales in the previous

310 stage. Here, a MATLAB RF package adapted by Abhishek Jaiantilal (Jaiantilal, 2013) based  
 311 on (Liaw and Wiener, 2002) is used.

312 The testing process is similar to the training process. If testing samples are classified as  
 313 negative (probability < cut-off value) by a RF classifier at the current layer, then they are  
 314 removed and identified as normal regions. If they are classified as positive (probability  $\geq$  cut-  
 315 off value), then they are put into the next RF classifier in the cascade. The testing phase of  
 316 CasRFs is illustrated in Figure 8. Since the design of the CasRFs is to filter out the easier  
 317 (less lesion-like) negative samples first and then the difficult (more lesion-like) ones later, the  
 318 negative samples in the testing data that are less likely to be masses are gradually discarded  
 319 from the layers of RF models and only the samples left at the last layer are regarded as the  
 320 suspicious regions.



321  
 322 Figure 7. The training of the self-grown cascaded random forests. (a) shows an overview of the training of the  
 323 whole cascade, and (b) shows the training of a layer. DRF stands for RF with default settings.



324

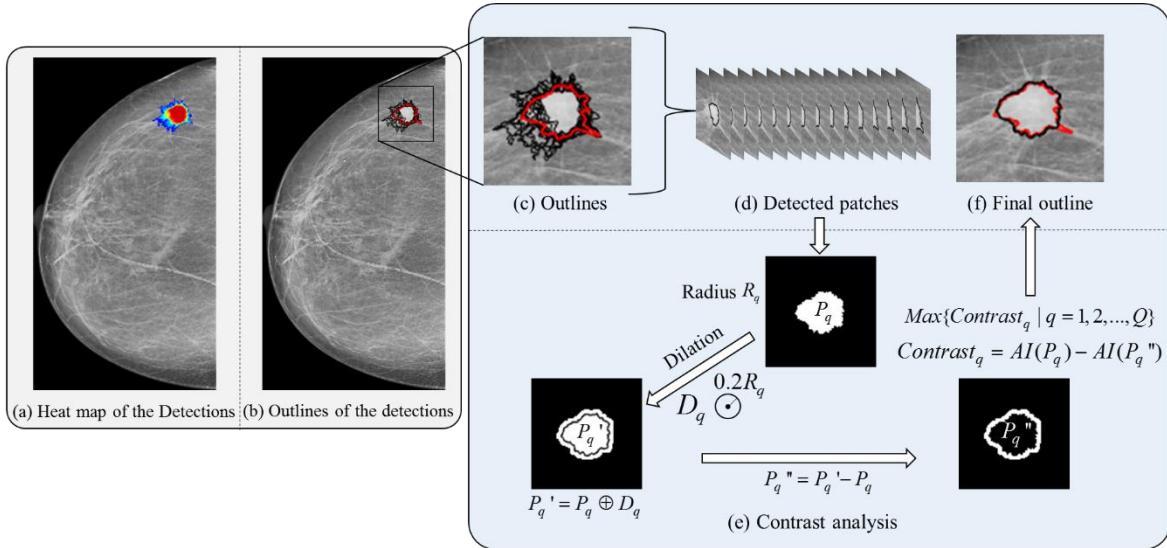
Remove the negative samples gradually throughout the layers

325

Figure 8. A schematic of the testing process using the CasRFs.

326 **2.5 Post-processing and visualizing the results**

327 After the classification, the outlines of the detected patches are plotted on the mammograms  
 328 and a heat map of the probability is also generated as shown in Figure 9 (a), (b). However,  
 329 there are often multiple identified suspicious patches overlapping with each other as shown in  
 330 Figure 9 (c). To improve the visualization of the detection output and evaluate the  
 331 segmentation accuracy, a label fusion method is designed to merge the overlap regions into  
 332 one region by analysing the contrast between the inner area and the outer area of the region.  
 333 Figure 9 (e) shows the technique of calculating the local contrast. Here, the contrast is  
 334 defined as the subtraction between the average intensity of pixels within patch  $P_q$  and the  
 335 average intensity of pixels within the background area  $P_q''$ . The patches are firstly merged  
 336 with in their scales and then between the scales. The outline of the region with the highest  
 337 contrast is kept as the final contour of the detected region as shown in Figure 9 (f).



338

339 Figure 9. A schematic of the label fusion through local contrast analysis. (a) is the heat map of the detections  
 340 from the CasRFs and all the detected patches are plotted on (b). (d) shows the detected patches and they go  
 341 through a contrast analysis approach as shown in (e).  $P_q$  represents the inner region and  $P_q''$  represents the  
 342 outer region of a detected patch.  $AI$  stands for average intensity of pixels within a certain region. The outline of  
 343 the patch that has the highest contrast is selected as the final outline of the detection. (f) shows the final fused  
 344 outline.

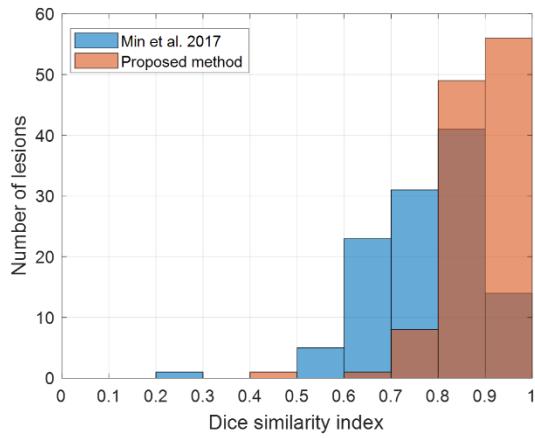
### 3 Experiments

345 Our system is evaluated on two public available datasets, INbreast and DoD BCRP. For the  
 346 evaluation on INbreast, we use the same evaluation method, repeated random sub-sampling  
 347 validation (Dubitzky et al., 2007), as studies (Dhungel et al., 2015a, 2016; Dhungel et al.,  
 348 2017). In these studies, the INbreast dataset is randomly split into training, validation and  
 349 testing sets five times. In our work, we used the same testing set while we combine the  
 350 training and validation sets together as a larger training set in each validation since random  
 351 forest cross-validates itself internally when building the trees. For the evaluation on DoD  
 352 BCRP, we simply use the pre-separated training and testing sets specified in this dataset.  
 353 During the classification of region candidates, the training samples are labelled as positive or  
 355 negative for training according to their dice similarity index (DSI) (Dice, 1945) to the ground

truth. The DSI between a region  $A$  and ground truth  $B$  is defined as  $2 \times |A \cap B| / (|A| + |B|)$  (Dice, 1945). The patches with the highest DSI are labelled as positive and those with a DSI lower than 0.1 are labelled as negative in the training sets. Overall, major inputs required by the system include a size range of the lesions (which has been specified in the INbreast documentation (Moreira et al., 2012)), a number of linear structuring elements used in MMS (set as 18 in this study, recommend values larger than 10), a number of scales used in MMS (set as 4 in this study), a number of levels in the multi-level thresholding (set as 16), and a sensitivity threshold  $ST$  for CasRFs (set as 0.99 for INbreast and 0.97 for DoD BCRP). The cut-off value inside random forests during training remains as the default value 0.5. All experiments are carried on a Dell desktop with Intel Core i7-4790 CPU @ 3.60GHz, 16GB RAM.

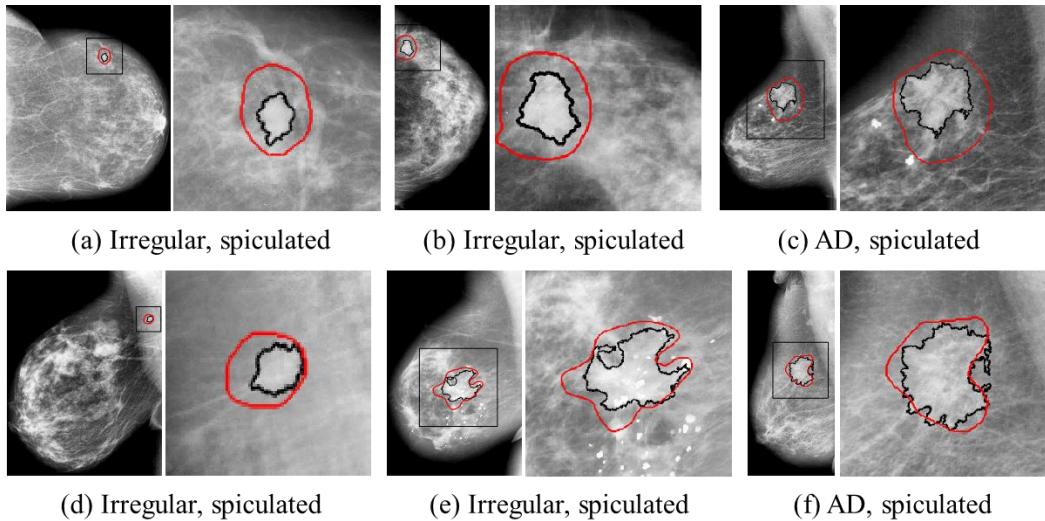
## 4 Results

The MMS filter working with MLT can capture 100% of the masses in INbreast at a minimum DSI with the ground truth of 0.4, and 98% of the segmentations have a DSI with the ground truth higher than 0.7. Figure 10 shows a performance comparison of the region candidate segmentation between our previous method (using circular structuring elements) and the proposed method. It can be seen the proposed method produces significantly better segmentations compared to our previous work (Min et al., 2017). Since DoD BCRP provides no accurate segmentation of lesions in the annotations, we could not evaluate the ROI segmentation on this dataset quantitatively, instead, we present some examples of the region candidate segmentation in Figure 11, where the masses are spiculated with an irregular shape or architectural distortion. It can be seen that the ROI segmentation method based on MMS and MLT still performs well on irregular masses.



379

380 Figure 10. Performance comparison of region candidate segmentation between the proposed method and our  
381 previous study (Min et al., 2017).



382

383 Figure 11. Examples of the region candidate segmentation on DoD BCRP. The red lines stand for the ground  
384 truth and the black lines stand for the segmentation of the lesion using MMS and MLT. Note that only the  
385 segmentations of the lesions are outlined in these images. AD stands for architectural distortion.

386 Before the label fusion, the system outputs heat maps based on the probabilities generated  
387 from CasRFs. Figure 12 shows examples of the heat maps showing the confidence of the  
388 detections. The boxplots of FP rate and DSI on INbreast dataset are shown in Figure 13. After  
389 the label fusion, if a detected region has a  $DSI \geq DT$  ( $DT \in \{0,1\}$ ), it is regarded as a true  
390 positive, if not, it is regarded as a true negative. For validation, the repeated random sub-  
391 sampling method is used.(Dhungel et al., 2017; Dubitzky et al., 2007) When the threshold

392  $DT$  is set as 0.2, the average sensitivity across all the validations is 0.93 at 4.73 FPs per  
 393 image (FPI) before the label fusion and at an average FPI of 1.64 FPs/image after the label  
 394 fusion by contrast analysis. The average sensitivity  $S_A$  and average FPI  $FPI_A$  here are defined  
 395 as,

$$396 S_A = \frac{\sum_{v=1}^V \text{number of lesions detected } (v)}{\sum_{v=1}^V \text{number of lesions}(v)}, \quad 7$$

$$397 FPI_A = \frac{\sum_{v=1}^V \text{number of FPs } (v)}{\sum_{v=1}^V \text{number of images}(v)}, \quad 8$$

398 where  $v$  stands for the  $v^{th}$  validation and  $V$  stands for the number of validations which is 5 in  
 399 this work. A normalized partial area under the FROC curves (PAUC) is also calculated to  
 400 provide an integral measurement of the performance. The PAUC is defined as

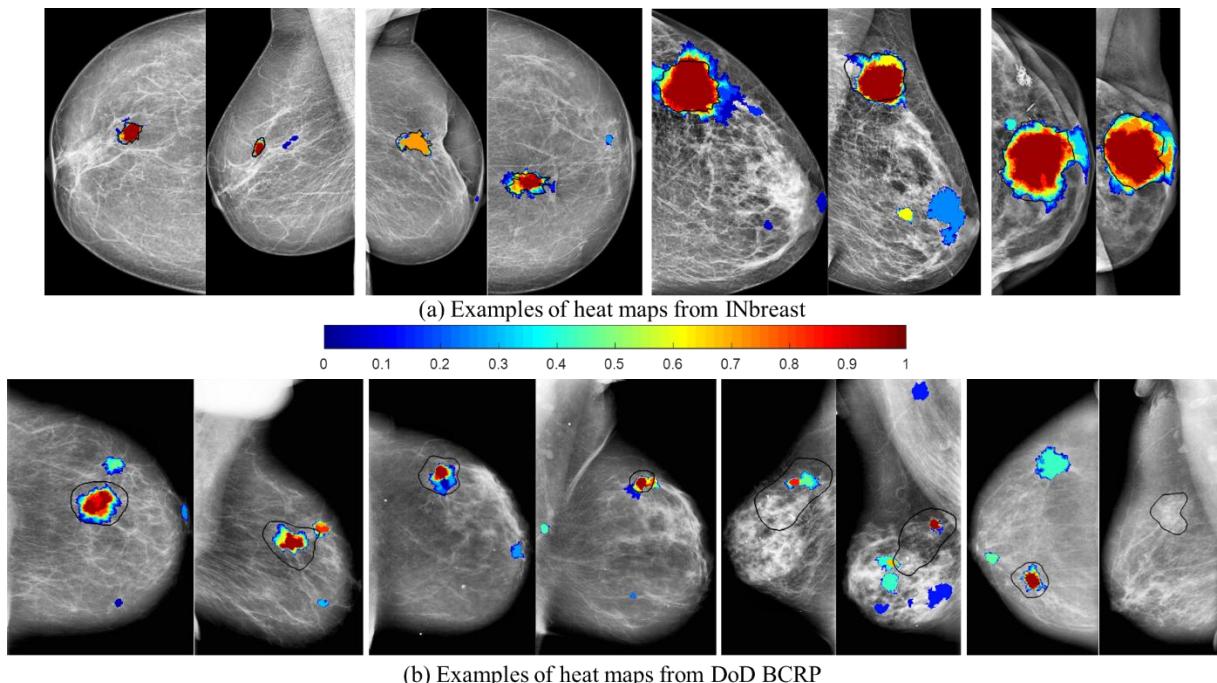
$$401 PAUC = \frac{1}{|\beta - \alpha|} \int_{\alpha}^{\beta} TPR dFPI, \text{ where } [\alpha, \beta] \text{ stands for the range of FPI and set as [0, 5] in our case.}$$

402 The PAUC is 0.91, and 0.89 for the proposed method and our previous method respectively  
 403 when  $DT = 0.2$ . The PAUC is 0.90 and 0.88 for the proposed method and our previous  
 404 method when  $DT = 0.5$ .

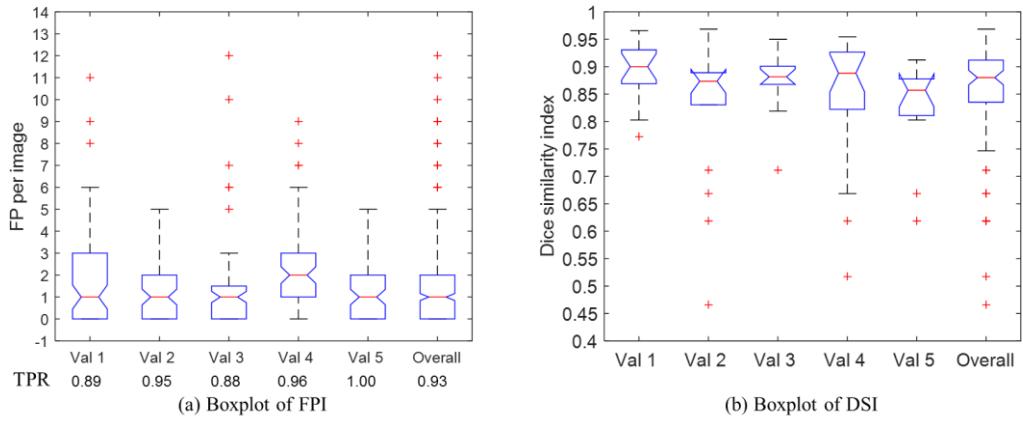
405 The average DSI of the segmentation on INbreast is  $0.85 \pm 0.10$  (minimum DSI  $> 0.4$ ). The  
 406 receiver operating characteristic (ROC) and free response operating characteristic (FROC)  
 407 curves before the label fusion are shown in Figure 14 (a) and (b). Figure 15 shows the FROC  
 408 curves on the testing set after the label fusion under the condition of  $DT = 0.2$  and  $DT = 0.5$ ,  
 409 comparing with previous studies (Dhungel et al., 2017; Min et al., 2017). Figure 16 shows the  
 410 relations between TPR, FPI and cut-off value in CasRFs before and after the label fusion. As  
 411 shown in the FROC curve in Figure 15(c), the detection system can achieve an average  
 412 sensitivity of  $0.91 \pm 0.06$  at approximately 1.2 FPs per image when the average DSI is

413  $0.86 \pm 0.08$  (minimum DSI > 0.6). The segmentation examples of lesions in various sizes and  
414 shapes from INbreast are shown in Figure 17.

415 Additional evaluation of the proposed method is also conducted on DoD BCRP. Since DoD  
416 BCRP only provides an approximate manual annotation, it can be difficult to evaluate the  
417 segmentation performance. Here, we decide if the overlap ratio (Ben-Ari et al., 2017;  
418 Dhungel et al., 2015a; Reichel and Cole, 2016) between a detected region and the ground  
419 truth is higher than 0.7, the lesion is regarded as detected. Here, the overlap ratio between two  
420 regions  $A$  and  $B$  is defined as  $A \cap B / \min(A, B)$  (Reichel and Cole, 2016; Szymkiewicz, 1934).  
421 Thus the sensitivity of detection on DoD BCRP is 0.81 when the sensitivity threshold  $ST$  in  
422 the CasRFs is set as 0.97 and the cut-off value in the individual RFs is 0.5 by default, at a FPI  
423 of 3.59 after label fusion. The comparison between this work and other state-of-art methods  
424 in terms of TPR, FPI and segmentation accuracy (DSI) where applicable is shown in Table 1.

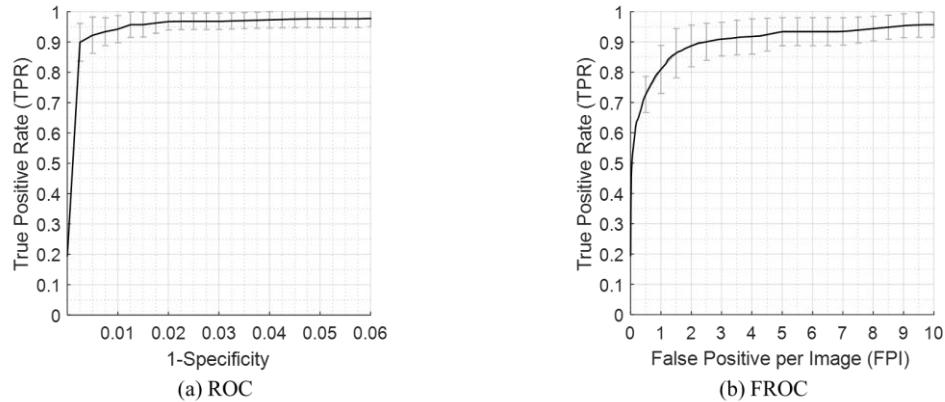


425  
426 Figure 12. The examples of detection heat maps and the colour bar showing the mapping of probability values  
427 into the colour map. The ground truth is labelled with black lines.



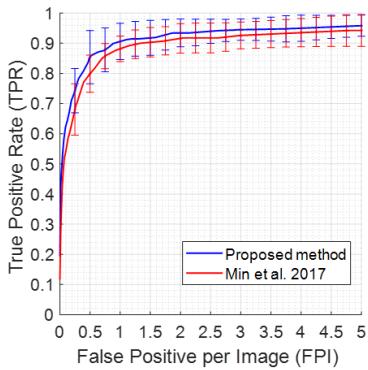
428

429 Figure 13. Boxplots of sensitivity & FP rate and DSI. ‘Val’ is short for ‘validation’ and TPR stands for true  
 430 positive rate (sensitivity). (a) shows the boxplots of FPI for all the validations and the TPR at each validation. A  
 431 boxplot of the overall FPI together with the average TPR across all the validations is also shown in (a). (b)  
 432 shows the boxplots of DSI (segmentation accuracy) at each validation and also an overall boxplot of all the  
 433 detected masses among all the validations. The median of the overall average DSI is 0.88.

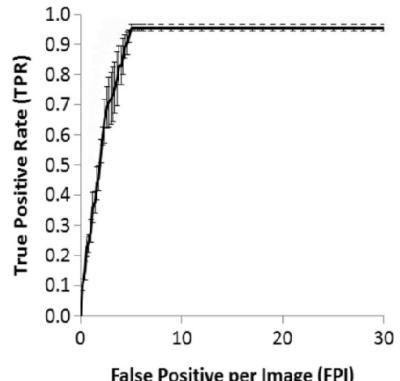


434

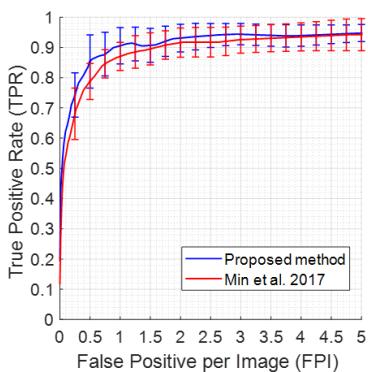
435 Figure 14. ROC and FROC of the classification on the testing samples from INbreast before label fusion.



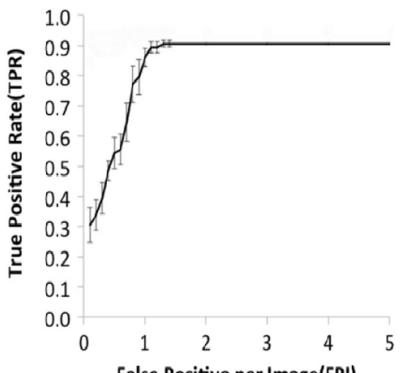
(a) FROC of the proposed method and Min et al. 2017 ( $DT=0.2$ )



(b) FROC of Dhungel et al. 2017 ( $DT=0.2$ )



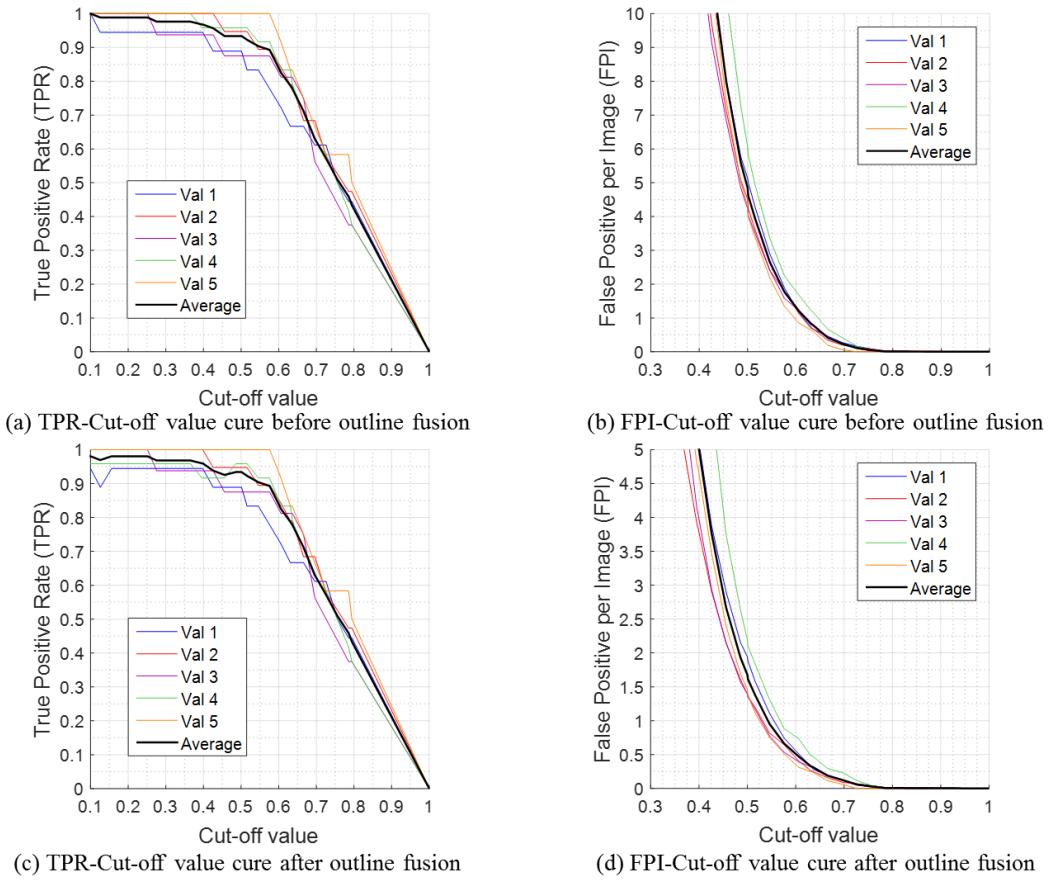
(c) FROC of the proposed method and Min et al. 2017 ( $DT=0.5$ )



(d) FROC of Dhungel et al. 2017 ( $DT=0.5$ )

436

437 Figure 15. The FROC curves on the testing dataset after label fusion. (a) and (b) show the FROC curves of the  
 438 proposed method, our previous work (Min et al., 2017) and study (Dhungel et al., 2017) with  $DT = 0.2$ . (c) and  
 439 (d) show the FROC curves of the proposed method, our previous work (Min et al., 2017) and study (Dhungel et  
 440 al., 2017) with  $DT = 0.5$ .



441

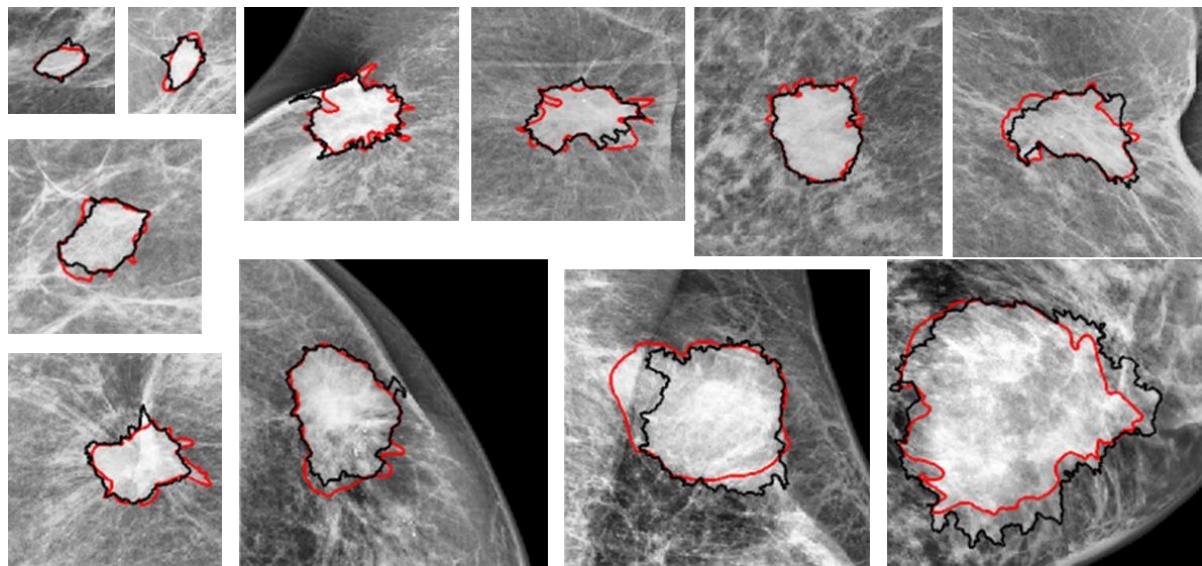
442 Figure 16. The relationship between TPR, FPI and the cut-off value (default 0.5) before and after the label  
 443 fusion on INbreast. ‘Val’ is short for ‘validation’. The black line represents the average TPR and FPI when  
 444 given different cut-off values ( $DT = 0.2$ ).

445 Table 1. Performance comparison between the proposed method and previous publications. The third column  
 446 indicates whether the method generates both detection and segmentation of the lesions or only detection. ✓  
 447 means the method detects as well as segments the lesions. TPR stands for true positive rate (sensitivity), FPPI  
 448 stands for FPs per image, and DSI stands for dice similarity index.

Methods	Dataset	Segmentation	Avg. TPR @ FPPI	Avg. DSI
(Kozegar et al., 2013)	INbreast	✓	0.87@3.67	-
(Dhungel et al., 2015a)	INbreast		0.96@1.2, 0.87@0.8	-
(Dhungel et al., 2016) & (Dhungel et al., 2017)	INbreast	✓	0.90@1.3	0.85
(Min et al., 2017)	INbreast	✓	0.94@1.99	0.80
The proposed method	INbreast	✓	0.93@1.64, 0.91@1.2	0.85, 0.86
(Beller et al., 2005)	DoD BCRP	✓	0.70@8	-

(Dhungel et al., 2015a)	DoD BCRP		0.75@4.8,0.7@4	-
(Min et al., 2017)	DoD BCRP	✓	0.77@3.93	-
The proposed method	DoD BCRP	✓	0.81@3.59	-

449



450

451 Figure 17. Segmentation examples of masses in various shapes and sizes from INbreast dataset. The black lines  
452 represent the segmentation from our method, and the red lines are the ground truth.

453 **5 Discussion**

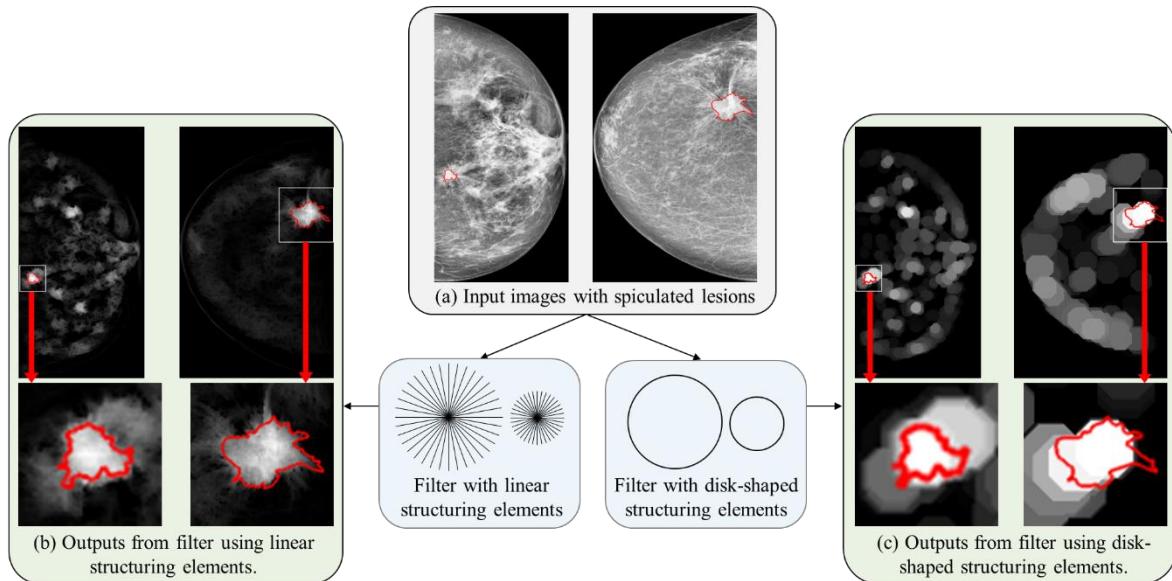
454 **5.1 Unsupervised segmentation using morphological sifting and multi-level thresholding**

455 In this work, a ROI segmentation method based on MMS and MLT is designed to separate  
456 the lesions from the background tissue. As mentioned in section 1, the main challenges for  
457 breast mass segmentation are the large variations of lesions in size, shape and contrast (Oliver  
458 et al., 2010). Generally, lesions appear in irregular shape and spiculated margin are highly  
459 suspicious in malignancy. In order to cope with these challenges, the multi-scale  
460 morphological sifters are developed to extract the regions of interest (ROI) within the mass  
461 size range. Applying the morphological sifters at multiple scales can limit the size variation at  
462 each individual scale. The new morphological sifting method adopts linear structuring  
463 elements that are mapped in a multiple directions and length to imitate the pattern of masses

as shown in Figure 2. After MMS and MLT, the average DSI between the best segmented region candidates and the ground truth is  $0.89 \pm 0.07$  among all the cases in INbreast, which indicates the segmentation method performs well regardless of the variations in the mammographic appearance of the lesions. In previous studies, disk-shaped structuring elements have been used to extract mass-like patterns (Chu et al., 2015; Min et al., 2017). They are effective in extracting densities of a semi-circular shape, however they may not be as effective on masses with highly irregular shapes and spiculated margins that are more likely to be malignant (Karssemeijer and te Brake, 1996). Figure 18 shows the comparison between applying linear (this work) and disk-shape structuring elements (Min et al., 2017) to extract spiculated masses with relatively irregular shapes. Figure 10 has also provided a statistic comparison of the segmentation accuracy for lesions between these two methods.

The new multi-scale morphological sifting method clearly show advantages over our previous method. As to the region candidate generation methods proposed by recent studies (Casti et al., 2016; Görgel et al., 2013; Kozegar et al., 2013; Zhang et al., 2016), these approaches mainly aim at locating suspicious regions (e.g. to find the focal area of the lesions) and did not report the segmentation accuracy of ROI generation, while the proposed MMS and MLT based method generates relatively accurate segmentations. Moreover, this approach is not affected by the existence of the pectoral muscle. There are special cases where in the MLO view, lesions can fully or partially overlap with the pectoral muscle as shown in Figure 11 (d), and removing pectoral muscle means at least part of the lesion will also be removed. Many studies remove the pectoral muscle in pre-processing (Casti et al., 2016; Soulami et al., 2017; Tan et al., 2015; Zhang et al., 2016), while the proposed method does not need to, therefore it can extract lesions that overlap with the pectoral muscle. This ROI generation method is the foundation of the detection and segmentation of the whole system. Compared with studies based on convolutional neural network (CNN) (Dhungel et al. 2015, 2016 and

489 2017), the MMS sifter is much simpler and computationally efficient. It removes the  
 490 irrelevant parts of the breast, extracts the suspicious regions accurately at the early stage of  
 491 the algorithm and provides exact regions of interest for the later learning step. This helps to  
 492 avoid the need of iterative searching as in CNN based algorithms.



493  
 494 Figure 18. Comparisons between using linear and disk-shape structuring elements. The two examples contain  
 495 lesions with spiculated margin. The red lines represent the ground truth.

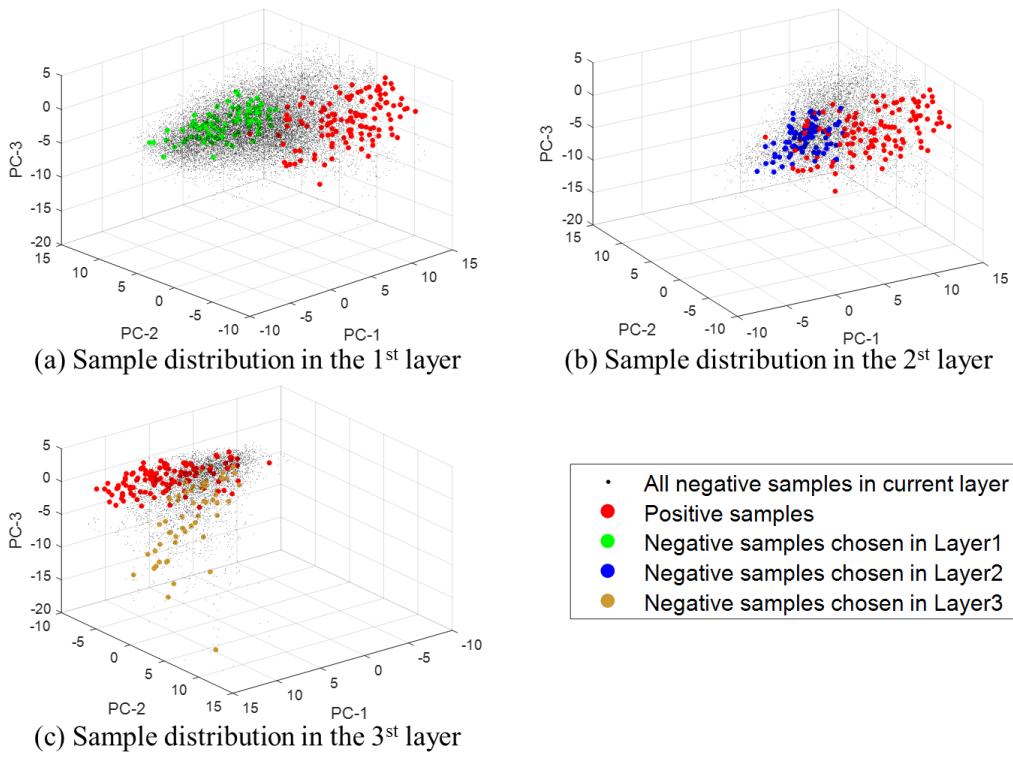
496 Although the morphological sifting method performs well on most types of lesions, it is still  
 497 not quite as effective in coping with some architectural distortions that have vague or  
 498 scattered spicules, and lack significant central densities. An architectural distortion (AD) is  
 499 defined as an interruption of normal breast pattern with lines, often appearing as a star-shaped  
 500 distortion, with no definite mass visible (Moreira et al., 2012). Figure 21 shows two cases of  
 501 AD from INbreast and DoD BCRP. Figure 21 (a) shows an AD with visible central densities  
 502 that is well segmented and detected by the proposed algorithm. Figure 21 (b) shows an AD  
 503 that has no obvious (very small) focal area, and scattered spicules in various angles. The  
 504 segmentation of the AD in Figure 21 (b) only captures a small density near the centre of the  
 505 AD, instead of identifying the centre of the lesion precisely or contouring the whole lesion.

506 This is mainly due to the fact that there are not as many ‘strokes’ passing the morphological  
507 sifters in an area where there are no obvious central densities and the spicules are scattered.  
508 However, the frequency of this type of architectural distortion is relatively low. There are  
509 around 5 ADs appearing in this pattern only from DoD BCRP, which our method has  
510 difficulties to generate relatively satisfactory segmentations. The proposed algorithm is  
511 supposed to be a general model to extract majority types of masses, while these characteristic  
512 ADs are relatively rare and may need methods specially designed to identify (Rangayyan et  
513 al., 2010). There are methods designed to detect line structures that can be used to detect  
514 architectural distortions. Study (Zwiggelaar et al., 2004) compared four linear structure  
515 detecting approaches in mammographic images (line operator, orientated bins, Gaussian  
516 derivatives and ridge detector). Study (Matsubara et al., 2015) adopted direction and  
517 background filters to extract the line structures to detect ADs. Rami Ben-Ari *et al* extracted  
518 ROIs by sampling the parenchymal tissues (segmented by an unsharp mask filter) where ADs  
519 are likely to be found (Ben-Ari et al., 2017). However, it is still unclear whether these  
520 methods would be effective on capturing ADs that have relatively faint, scattered spicules  
521 and no significant central densities.

## 522 5.2 Tackling class imbalance using self-grown CasRFs

523 After the generation of region candidates, there is a huge imbalance between the positive and  
524 negative samples. To cope with the class imbalance problem, we adopted an ensemble of  
525 random forests structured as a cascade. The novelty of the CasRFs lies in the under-sampling  
526 guided by probability-ranking and its ability to grow adaptively according to the skewness of  
527 the training data. As we have stated in Section1, many previous studies adopt naïve random  
528 under-sampling when training each individual base classifier. However, this type of random  
529 under-sampling could lead to the possibility that the positive and negative samples in the

530 subset may not be separable. If the subset is too challenging for the classifier, it could affect  
 531 the performance of the whole cascade at an early stage. By using the probability-ranking  
 532 based under-sampling, only the more discriminative samples (according to the probability  
 533 ranking) are used to train the RF at a current layer. By doing this, the training process is not  
 534 only stabilized, but also giving the RF a relatively less challenging training subset. Figure 19  
 535 shows an example of the negative samples chosen in the first three layers during the training  
 536 process. It can be seen that the negative samples chosen in the three layers are the samples  
 537 that are clearly separable from the positive samples. The selected negative samples start from  
 538 being rather far away from the positive samples in layer 1 to being closer to the positive  
 539 samples in layer 2 and 3.



540  
 541 Figure 19. Example of the negative samples selected to form the training subset with all the positive samples in  
 542 layer 1, layer 2 and layer 3 during the training in one of the validations. The dimension of the sample data has  
 543 been reduced by principal component analysis (PCA). The first three principal components are used in this  
 544 figure.

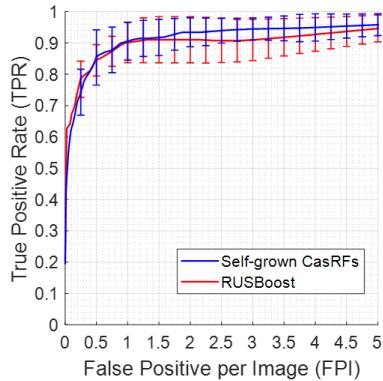
545 The CasRFs also give the user control over the detection rate during the training since  
546 attaining high sensitivity is crucial in mass detection. To reach the ideal sensitivity set by the  
547 user during the training of each individual RF, the CasRFs gradually reduce the number of  
548 negative samples (starting from the more difficult ones to remove) in the initially balanced  
549 training set. By doing this, we in essence reverse the imbalance by letting the positive  
550 samples become the majority class and the RF will naturally focus more on the detection of  
551 positive samples, without any need to adjust weights to increase the cost of misclassifying the  
552 minority samples (Bria et al., 2016; Viola and Jones, 2002) or manipulate the threshold in the  
553 RF (Wei et al., 2015).

554 The CasRFs only use part of the negative samples during the whole training process, since  
555 after building the RF model at a current layer, the unused negative samples are put through  
556 the model and some of them that do not contain any useful information will be eliminated  
557 without the need to be used for training in the next layer. This can be observed in Figure 19,  
558 where there are fewer negative samples in Figure 19 (b) than in Figure 19 (a) since some of  
559 the non-selected negative samples have been discarded by the RF model built in layer 1. This  
560 will also reduce the computational cost especially when the training set is highly skewed.

561 Unlike studies (Baumann et al., 2013; Tang et al., 2012) which used fixed number of layers,  
562 the CasRFs can adapt to highly skewed training data without the need to set the number of  
563 layers in the cascade (i.e. the number of RFs), by setting an explicit stopping criteria. The  
564 CasRFs perform relatively well when facing a class imbalance ratio of 1:160 in this work,  
565 and also achieved satisfactory results in the previous study where the imbalance ratio can be  
566 as high as 1:2000 (Min et al., 2017), while in other imbalance learning related studies, the  
567 imbalance ratio is around 1:6 (Wei et al., 2015), 3:10 (Viola and Jones, 2002) and 1:20  
568 (Kozegar et al., 2013) respectively. The feature extraction and training can be carried out on a

569 regular desktop without the need of using graphic processing units like deep learning based  
570 studies (Dhungel et al., 2017; Ribli et al., 2017).

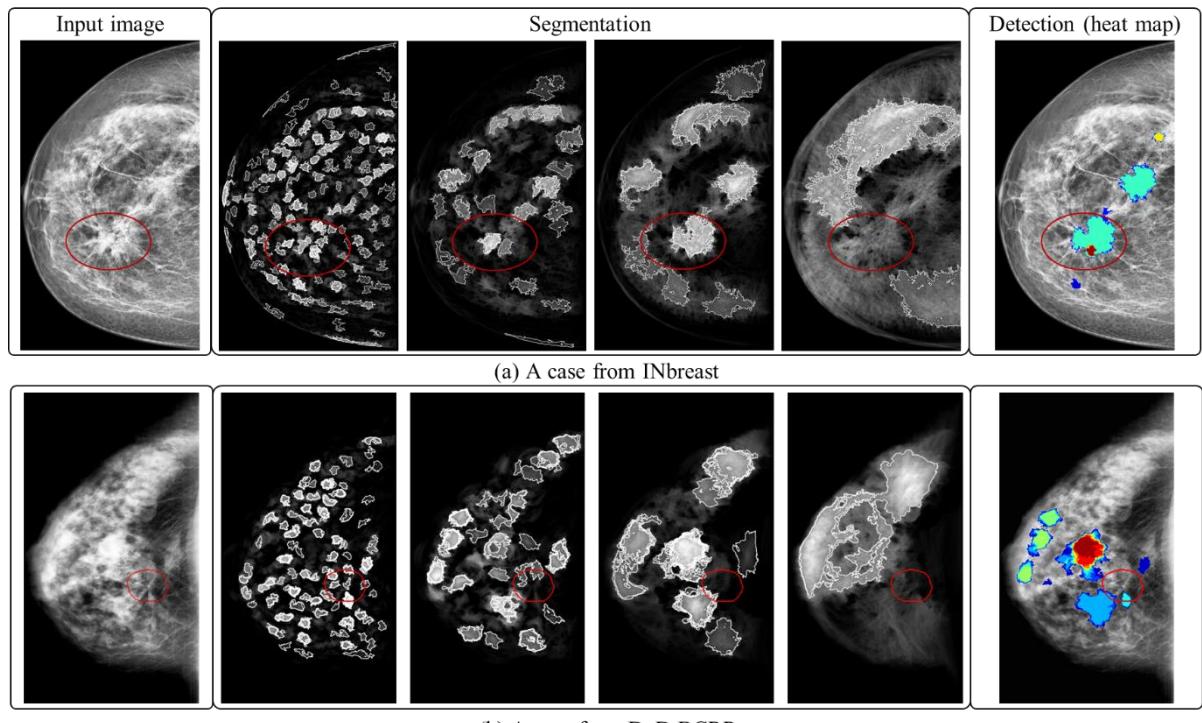
571 The design of the CasRFs was firstly introduced in our preliminary work (Min et al., 2017).  
572 However, during the training in this work, all the negative samples are used to initialize the  
573 first layer instead of random selecting a subset of negative samples as in (Min et al., 2017).  
574 And the samples from all scales are put together as the training set instead of training the  
575 CasRFs on each individual scale as in (Min et al., 2017), thus we could get a more sufficient  
576 number of positive samples (masses). If trained on each individual scale, the system achieves  
577 an average sensitivity of 0.90 at a FP rate of 1.45FPs/image at the default cut-off value of 0.5.  
578 This shows that in our case, having more positive samples in the training set does benefit the  
579 sensitivity, however also slightly increase the FP rate. We also compare the classification  
580 performance of the self-grown CasRFs with RUSBoost (Seiffert et al., 2010). The RUSBoost  
581 uses decision tree as base learner, with the maximal number of decision splits per tree set as  
582 the number of training instances. For the ensemble structure, the number of learning circles is  
583 set as 1000, and the learning rate for shrinkage is set as 0.1. Figure 20 shows the FROC  
584 curves of the proposed self-grown CasRFs and RUSBoost. The PAUC is 0.91 and 0.89 for  
585 CasRFs and RUSBoost respectively. The CasRFs also take less time to classify the samples  
586 than RUSBoost. For instance, given a certain testing image represented by its feature matrix,  
587 it takes approximately 4s for CasRFs and 14s for RUSBoost. The better overall performance  
588 and less execution time make the CasRFs more favourable than RUSBoost.



589

590

Figure 20. The FROC curves of the self-grown CasRFs and the RUSBoost.



591

592 Figure 21. Two examples of the segmentation affected by architectural distortions. (a) is from INbreast and it  
593 has visible central densities within the distortion. (b) is from DoD BCRP and it contains a spiculated distortion  
594 with scattered spicules and no central density. The red lines indicate the ground truth.

### 595 5.3 Visualizing the detection and segmentation results

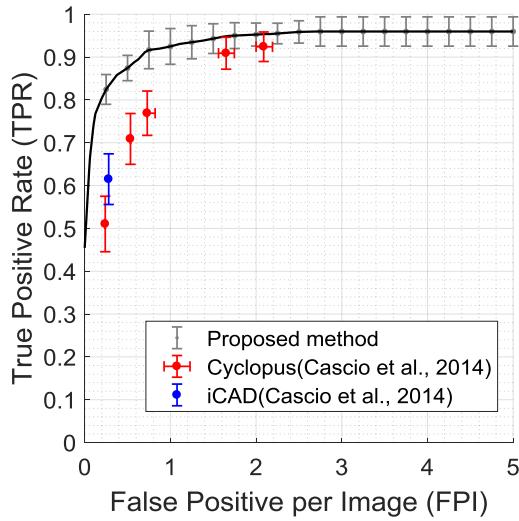
596 After classifying the samples using CasRFs, each patch sample was assigned a probability  
597 score by summing up the scores generated from all the RFs in the cascade. These scores are  
598 used in generating heat maps as shown in Figure 12 where only the positive patches are

599 marked. The heat map is a very useful tool to visualize the confidence of the detection and  
600 can be used as the final detection output. Figure 16 (a), (b) analyse the relations between TPR,  
601 FPI and cut-off value of CasRFs (before label fusion). To further clean up the outlines of the  
602 detections, we adopted a simple label fusion method to merge the overlapped regions by  
603 choosing the patch with the highest contrast against the background. Comparing Figure 16 (a)  
604 and (c), it can be noticed that the label fusion method does have the possibility of sacrificing  
605 the sensitivity sometimes as TPI drops within cut-off value 0.4~0.45 in validation 4. This  
606 happens due to the incorrect merging of overlap regions. An alternative fusion method is to  
607 choose the patch with the highest predicted probability. However, the average DSI calculated  
608 under this method is  $0.84 \pm 0.10$ , which is slightly lower than the one under contrast analysis.  
609 We believe there are more sophisticated ways to merging the overlap raw detections such as  
610 expectation-maximization algorithm (Warfield et al., 2004) and atlas-based method  
611 (Langerak et al., 2010).

#### 612 5.4 Performance evaluation

613 Since full-field digital mammographic (FFDM) has replaced the screen-film due to the higher  
614 image quality (de Munck et al., 2016) and DoD BCRP lacks an accurate ground truth, we  
615 mainly evaluated the system's performance on INbreast and only provided the performance  
616 on DoD BCRP as a supplemental result. Table 1 shows the performance comparison between  
617 this method and previous publications evaluated on INbreast and DoD BCRP. Our new  
618 method yields competitive results in both mass detection and segmentation. The validation  
619 approach on INbreast by (Eltonsy et al., 2007) is not known, and therefore, it could not be  
620 determined if our method's performance is certainly more favourable. Previous study  
621 (Dhungel et al., 2015a) shows the best mass detecting performance using deep learning and  
622 fixed two layers of RFs. However, it only locates the suspicious regions using bounding

623 boxes, without a precise segmentation of the regions. (Dhungel et al., 2016) and (Dhungel et  
624 al., 2017) present both detection and segmentation results. In these studies, the suspicious  
625 regions are located using method from the previous study (Dhungel et al., 2015a), then  
626 lesions are segmented from the background using another deep learning based algorithm  
627 (Dhungel et al., 2015b). Compared with (Dhungel et al., 2016) and (Dhungel et al., 2017),  
628 our method achieves a higher average sensitivity and DSI at a slightly lower FP rate as shown  
629 in Table 1. It is worth noting that our method has a much simpler structure with only one  
630 stage of machine learning compared with the multiple learning stages adopted in previous  
631 studies, each having very high computational complexity (Dhungel et al., 2015a, 2016;  
632 Dhungel et al., 2017). We also provide a comparative guide to the mass detection  
633 performances of two commercial systems, Cyclopus CAD (CyclopusCAD Ltd, Palermo,  
634 Italy), SecondLook (iCAD Inc. OH, USA), reported in study (Cascio et al., 2014). Evaluated  
635 on a FFDM dataset, CyclopusCAD achieved a sensitivity of 76.9% at 0.73 FPs per image,  
636 and iCAD yielded a sensitivity of 61.5% at 0.28 FPs per image for breast mass detection.  
637 Figure 22 shows the FROC curves of the commercial systems and the proposed method. Note  
638 that the FPI is solely calculated on normal images since the FPI of the two commercial  
639 systems in Figure 22 was only calculated on healthy cases (Cascio et al., 2014). It can be seen  
640 that the FROC curve of the proposed method in Figure 22 is above the ones of the two  
641 commercial systems. However, we still cannot be certain that our method performs better  
642 than these commercial systems since they are not tested on the same data.



643

644 Figure 22. FROC curves of the proposed method on normal images and two commercial systems reported in  
645 study (Cascio et al., 2014).

646 As to the performance on DoD BCRP, the MMS using linear structuring elements does show  
647 advantages in identifying spiculated lesions (except for some architectural distortions as  
648 discussed in section 5.1) compared to the disk-shape structuring elements as presented in  
649 (Min et al., 2017). In DoD BCRP, the majority of the lesions are spiculated lesions in  
650 irregular shapes, and the rest are architectural distortions and lesions in oval or round shapes.  
651 The missed lesions are either spiculated in irregular shapes or ADs. Only 50% (9/18, overlap  
652 ratio > 0.7) of the ADs are detected in the testing set from DoD BCRP. The low sensitivity on  
653 ADs is not unusual since study (Baker et al., 2003) reported only fewer than 40% of the ADs  
654 (image sensitivity) from the dataset used for evaluation were detected by two most widely  
655 available commercial CAD systems. The reason the proposed method does not perform well  
656 on ADs is that ADs are more characteristic and rarer compared with other types of lesions in  
657 the evaluation dataset. Moreover, the proposed method is still more oriented towards  
658 detecting lesions with visible central densities. In the future, we would like to explore the  
659 detection of architectural distortions, especially the ADs that have no central densities with  
660 scattered and thin radiating lines as discussed in 5.1.

661    **6 Conclusion**

662    In this paper, we presented a mammographic breast mass detection and segmentation system  
663    using multi-scale morphological sifting and a self-grown, dynamically multi-layered,  
664    cascaded random forests algorithm. The MMS using linear structuring elements is able to  
665    extract lesions in various shapes and sizes while removing the background tissue. The self-  
666    grown CasRFs provide an effective and adaptive solution to imbalanced learning by adopting  
667    a probability-ranking based under-sampling approach and distributing the class imbalance  
668    throughout the layers of the cascade. In general, the proposed system achieves promising  
669    results in both detection and segmentation on digital and screen-film mammograms. **In the**  
670    **further work, we would like to explore the possibility of combining the MMS based region**  
671    **candidate generation method with CNN to avoid hand-selecting features and exhaustive**  
672    **searching for CNN itself. We would consider adopting transfer learning (West et al., 2007)**  
673    **since the size of the training set might be limited for building CNN from scratch. Further**  
674    **study is also needed to improve the system's ability to detect architectural distortions.**

675    **7 Acknowledgement**

676    We would like to thank Dr. Neeraj Dhungel for providing the INbreast validation sets used in  
677    this work.

678    Hang Min is supported by the China Scholarship Council.

679    **8 References**

680    DoD BCRP Spiculated Mass Detection Evaluation Data.  
681    [http://marathon.csee.usf.edu/Mammography/DDSM/BCRP/bcrp\\_mass\\_01.html](http://marathon.csee.usf.edu/Mammography/DDSM/BCRP/bcrp_mass_01.html)  
682    Andreea, G.I., Pegza, R., Lascu, L., Bondari, S., Stoica, Z., Bondari, A., 2011. The role of imaging  
683    techniques in diagnosis of breast cancer. J. Curr. Health Sci 37, 241-248.

- 684 Backes, A.R., Bruno, O.M., 2008. A new approach to estimate fractal dimension of texture images,  
685 International Conference on Image and Signal Processing. Springer, pp. 136-143.
- 686 Baker, J.A., Rosen, E.L., Lo, J.Y., Gimenez, E.I., Walsh, R., Soo, M.S., 2003. Computer-aided  
687 detection (CAD) in screening mammography: sensitivity of commercial CAD systems for  
688 detecting architectural distortion. Am. J. Roentgenol. 181, 1083-1088.
- 689 Baumann, F., Ehlers, A., Vogt, K., Rosenhahn, B., 2013. Cascaded Random Forest for Fast Object  
690 Detection, Scandinavian Conference on Image Analysis. Springer, pp. 131-142.
- 691 Beller, M., Stotzka, R., Müller, T.O., Gemmeke, H., 2005. An example-based system to support the  
692 segmentation of stellate lesions, Bildverarbeitung für die Medizin 2005. Springer, pp. 475-479.
- 693 Ben-Ari, R., Akselrod-Ballin, A., Karlinsky, L., Hashoul, S., 2017. Domain specific convolutional  
694 neural nets for detection of architectural distortion in mammograms, 2017 IEEE 14th International  
695 Symposium on Biomedical Imaging (ISBI 2017), . IEEE, pp. 552-556.
- 696 Bowyer, K., Kopans, D., Kegelmeyer, W., Moore, R., Sallam, M., Chang, K., Woods, K., 1996. The  
697 digital database for screening mammography, Third international workshop on digital  
698 mammography, p. 27.
- 699 Breiman, L., 2001. Random forests. Machine learning 45, 5-32.
- 700 Bria, A., Karssemeijer, N., Tortorella, F., 2014. Learning from unbalanced data: a cascade-based  
701 approach for detecting clustered microcalcifications. Med. Image Anal. 18, 241-252.
- 702 Bria, A., Marrocco, C., Molinara, M., Tortorella, F., 2016. An effective learning strategy for cascaded  
703 object detection. Information Sciences 340-341, 17-26.
- 704 Cascio, D., Fauci, F., Iacomi, M., Raso, G., Magro, R., Castrogiovanni, D., Filosto, G., Lenzi, R.,  
705 Vasile, M.S., 2014. Computer-aided diagnosis in digital mammography: comparison of two  
706 commercial systems. Imaging Med. 6, 13.
- 707 Cascio, D., Fauci, F., Magro, R., Raso, G., Bellotti, R., De Carlo, F., Tangaro, S., De Nunzio, G.,  
708 Quarta, M., Forni, G., 2006. Mammogram segmentation by contour searching and mass lesions  
709 classification with neural network. IEEE Trans. Nucl. Sci. 53, 2827-2833.

- 710 Casti, P., Mencattini, A., Salmeri, M., Ancona, A., Mangeri, F., Pepe, M.L., Rangayyan, R.M., 2016.  
711 Contour-independent detection and classification of mammographic lesions. Biomedical Signal  
712 Processing and Control 25, 165-177.
- 713 Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. SMOTE: synthetic minority over-  
714 sampling technique. Journal of artificial intelligence research 16, 321-357.
- 715 Chawla, N.V., Lazarevic, A., Hall, L.O., Bowyer, K.W., 2003. SMOTEBoost: Improving Prediction  
716 of the Minority Class in Boosting, In: Lavrač, N., Gamberger, D., Todorovski, L., Blockeel, H.  
717 (Eds.), Knowledge Discovery in Databases: PKDD 2003: 7th European Conference on Principles  
718 and Practice of Knowledge Discovery in Databases, Cavtat-Dubrovnik, Croatia, September 22-26,  
719 2003. Proceedings. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 107-119.
- 720 Chu, J., Min, H., Liu, L., Lu, W., 2015. A novel computer aided breast mass detection scheme based  
721 on morphological enhancement and SLIC superpixel segmentation. Med. Phys. 42, 3859-3869.
- 722 de Munck, L., de Bock, G.H., Otter, R., Reiding, D., Broeders, M.J.M., Willemse, P.H.B., Siesling, S.,  
723 2016. Digital vs screen-film mammography in population-based breast cancer screening:  
724 performance indicators and tumour characteristics of screen-detected and interval cancers. Br. J.  
725 Cancer 115, 517-524.
- 726 Dhungel, N., Carneiro, G., Bradley, A.P., 2015a. Automated Mass Detection in Mammograms using  
727 Cascaded Deep Learning and Random Forests, 2015 International Conference on Digital Image  
728 Computing: Techniques and Applications (DICTA). IEEE, pp. 1-8.
- 729 Dhungel, N., Carneiro, G., Bradley, A.P., 2015b. Deep structured learning for mass segmentation  
730 from mammograms, 2015 IEEE International Conference on Image Processing (ICIP). IEEE, pp.  
731 2950-2954.
- 732 Dhungel, N., Carneiro, G., Bradley, A.P., 2016. The automated learning of deep features for breast  
733 mass classification from mammograms, International Conference on Medical Image Computing  
734 and Computer-Assisted Intervention. Springer, pp. 106-114.
- 735 Dhungel, N., Carneiro, G., Bradley, A.P., 2017. A deep learning approach for the analysis of masses  
736 in mammograms with minimal user intervention. Med. Image Anal. 37, 114-128.

- 737 Dice, L.R., 1945. Measures of the amount of ecologic association between species. Ecology 26, 297-
- 738 302.
- 739 Dubitzky, W., Granzow, M., Berrar, D.P., 2007. Fundamentals of data mining in genomics and
- 740 proteomics. Springer Science & Business Media.
- 741 Eltonsy, N.H., Tourassi, G.D., Elmaghrary, A.S., 2007. A Concentric Morphology Model for the
- 742 Detection of Masses in Mammography. IEEE Trans. Med. Imaging 26, 880-889.
- 743 Ganesan, K., Acharya, U.R., Chua, C.K., Min, L.C., Abraham, K.T., Ng, K.-H., 2013. Computer-
- 744 aided breast cancer detection using mammograms: a review. IEEE Reviews in Biomedical
- 745 Engineering 6, 77-98.
- 746 Gonzalez, R.C., Woods, R.E., 2008. Digital image processing / Rafael C. Gonzalez, Richard E.
- 747 Woods, 3rd ed.. ed. Harlow : Pearson/Prentice Hall, Harlow.
- 748 Görgel, P., Sertbas, A., Ucan, O.N., 2013. Mammographical mass detection and classification using
- 749 Local Seed Region Growing-Spherical Wavelet Transform (LSRG-SWT) hybrid scheme. Comput.
- 750 Biol. Med. 43, 765-774.
- 751 Gromet, M., 2008. Comparison of computer-aided detection to double reading of screening
- 752 mammograms: review of 231,221 mammograms. Am. J. Roentgenol. 190, 854-859.
- 753 He, H., Ma, Y., 2013. Imbalanced learning: foundations, algorithms, and applications. John Wiley &
- 754 Sons.
- 755 Jaantilal, A., 2013. *Randomforest-matlab* [Online]. Available:
- 756 <https://github.com/jrderuiter/randomforest-matlab> [Accessed].
- 757 Jemal, A., Bray, F., Center, M.M., Ferlay, J., Ward, E., Forman, D., 2011. Global cancer statistics. CA
- 758 Cancer J. Clin. 61, 69-90.
- 759 Kang, P., Cho, S., 2006. EUS SVMs: Ensemble of under-sampled SVMs for data imbalance problems,
- 760 International Conference on Neural Information Processing. Springer, pp. 837-846.
- 761 Karssemeijer, N., te Brake, G.M., 1996. Detection of stellate distortions in mammograms. IEEE Trans.
- 762 Med. Imaging 15, 611-619.

- 763 Kharel, N., Alsadoon, A., Prasad, P., Elchouemi, A., 2017. Early diagnosis of breast cancer using  
764 contrast limited adaptive histogram equalization (CLAHE) and Morphology methods, Information  
765 and Communication Systems (ICICS), 2017 8th International Conference on. IEEE, pp. 120-124.
- 766 Kimori, Y., 2011. Mathematical morphology-based approach to the enhancement of morphological  
767 features in medical images. *J. Clin. Bioinforma.* 1, 33.
- 768 Kozegar, E., Soryani, M., Minaei, B., Domingues, I., 2013. Assessment of a novel mass detection  
769 algorithm in mammograms. *J. Cancer Res. Ther.* 9, 592.
- 770 Langerak, T.R., van der Heide, U.A., Kotte, A.N., Viergever, M.A., Van Vulpen, M., Pluim, J.P.,  
771 2010. Label fusion in atlas-based segmentation using a selective and iterative method for  
772 performance level estimation (SIMPLE). *IEEE Trans. Med. Imaging* 29, 2000-2008.
- 773 Li, H., Wang, Y., Liu, K.R., Lo, S.-C., Freedman, M.T., 2001. Computerized radiographic mass  
774 detection. I. Lesion site selection by morphological enhancement and contextual segmentation.  
775 *IEEE Trans. Med. Imaging* 20, 289-301.
- 776 Liao, P.-S., Chen, T.-S., Chung, P.-C., 2001. A fast algorithm for multilevel thresholding. *J. Inf. Sci.*  
777 Eng. 17, 713-727.
- 778 Liaw, A., Wiener, M., 2002. Classification and regression by randomForest. *R news* 2, 18-22.
- 779 Liu, L., Li, J., Wang, Y., 2015. Breast mass detection with kernelized supervised hashing, 2015 8th  
780 International Conference on Biomedical Engineering and Informatics (BMEI). IEEE, pp. 79-84.
- 781 Martins, L.d.O., Cardoso de Paiva, A., Corrêa Silva, A., Braz Junior, G., Gattass, M., 2009. Detection  
782 of Masses in Digital Mammograms using K-Means and Support Vector Machine. *ELCVIA.*  
783 Electronic letters on computer vision and image analysis 8, 39-50.
- 784 Matsubara, T., Ito, A., Tsunomori, A., Hara, T., Muramatsu, C., Endo, T., Fujita, H., 2015. An  
785 automated method for detecting architectural distortions on mammograms using direction analysis  
786 of linear structures, 2015 37th Annual International Conference of the IEEE Engineering in  
787 Medicine and Biology Society (EMBC). IEEE, pp. 2661-2664.

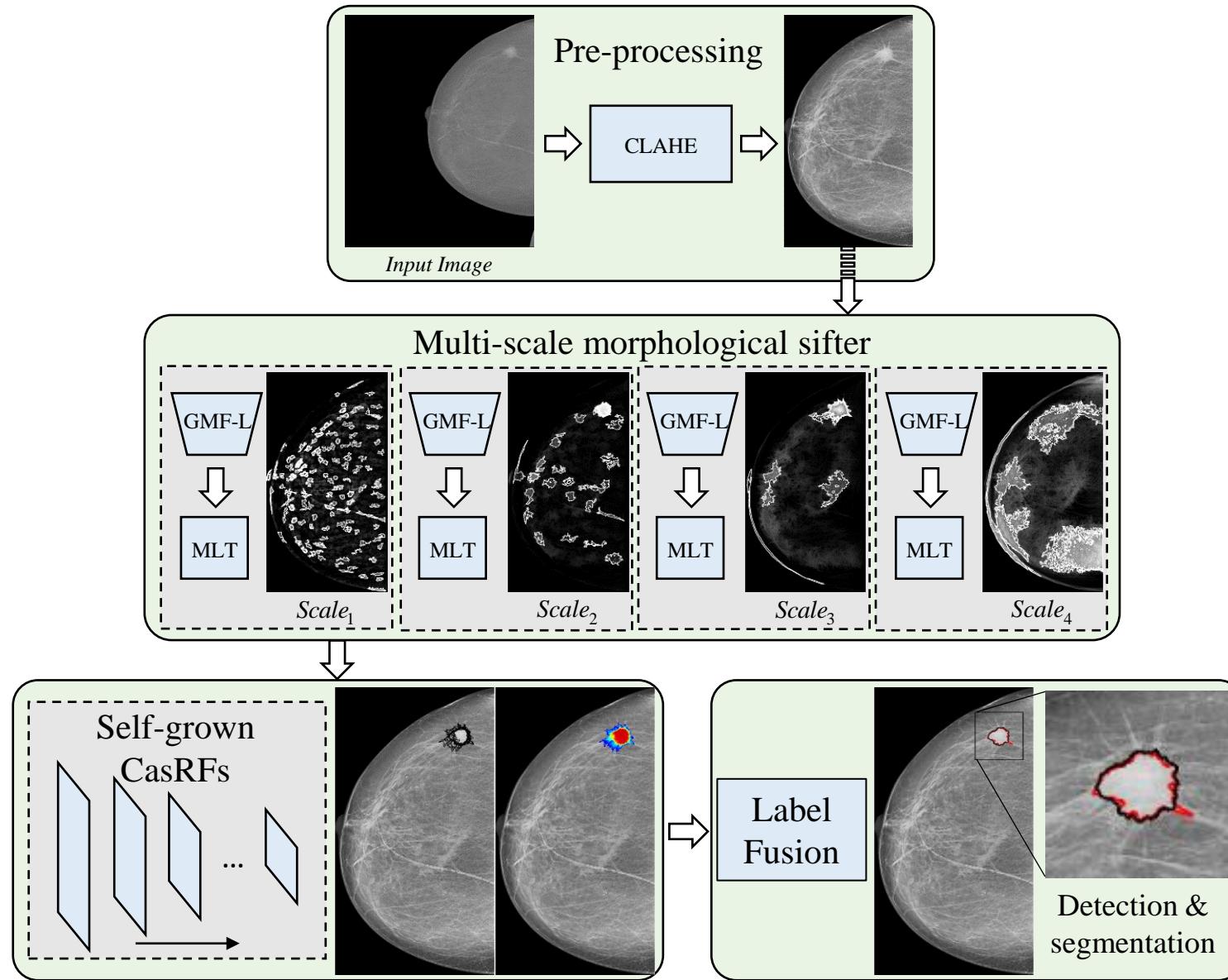
- 788 Min, H., Chandra, S.S., Dhungel, N., Crozier, S., Bradley, A.P., 2017. Multi-scale mass segmentation  
789 for mammograms via cascaded random forests, 2017 IEEE 14th International Symposium on  
790 Biomedical Imaging (ISBI 2017). IEEE, pp. 113-117.
- 791 Moreira, I.C., Amaral, I., Domingues, I., Cardoso, A., Cardoso, M.J., Cardoso, J.S., 2012. INbreast:  
792 toward a full-field digital mammographic database. Acad. Radiol. 19, 236-248.
- 793 Murthy, S.N., Kumar, A., Sheshadri, H., 2013. Mass Detection and Classification using Machine  
794 Learning Techniques in Digital Mammograms. International Journal of Computer Applications 76.
- 795 Oliver, A., Freixenet, J., Martí, J., Pérez, E., 2010. A review of automatic mass detection and  
796 segmentation in mammographic images. Med. Image Anal. 14, 87-110.
- 797 Otsu, N., 1979. A threshold selection method from gray-level histograms. IEEE Transactions on  
798 Systems, Man, and Cybernetics 9, 62-66.
- 799 Pizer, S.M., Amburn, E.P., Austin, J.D., Cromartie, R., Geselowitz, A., Greer, T., ter Haar Romeny,  
800 B., Zimmerman, J.B., Zuiderveld, K., 1987. Adaptive histogram equalization and its variations.  
801 Computer Vision, Graphics, and Image Processing 39, 355-368.
- 802 Rangayyan, R.M., Banik, S., Desautels, J.L., 2010. Computer-aided detection of architectural  
803 distortion in prior mammograms of interval cancer. J. Digit. Imaging 23, 611-631.
- 804 Reichel, U.D., Cole, J., 2016. Entrainment analysis of categorical intonation representations. Proc.  
805 P&P, Munich, Germany.
- 806 Ribli, D., Horváth, A., Unger, Z., Pollner, P., Csabai, I., 2017. Detecting and classifying lesions in  
807 mammograms with Deep Learning. arXiv preprint arXiv:1707.08401.
- 808 Schnabel, J.A., Giger, M.L., Karssemeijer, N., 2013. Breast Image Analysis for Risk Assessment,  
809 Detection, Diagnosis, and Treatment of Cancer. Annu. Rev. Biomed. Eng. 15, 327-357.
- 810 Seiffert, C., Khoshgoftaar, T.M., Van Hulse, J., Napolitano, A., 2010. RUSBoost: A hybrid approach  
811 to alleviating class imbalance. IEEE Transactions on Systems, Man, and Cybernetics-Part A:  
812 Systems and Humans 40, 185-197.
- 813 Soulami, K.B., Saidi, M.N., Tamtaoui, A., 2017. A CAD System for the Detection of Abnormalities  
814 in the Mammograms Using the Metaheuristic Algorithm Particle Swarm Optimization (PSO), In:

- 815 El-Azouzi, R., Menasche, D.S., Sabir, E., De Pellegrini, F., Benjillali, M. (Eds.), Advances in  
816 Ubiquitous Networking 2: Proceedings of the UNet'16. Springer Singapore, Singapore, pp. 505-  
817 517.
- 818 Szymkiewicz, D., 1934. Une conribution statistique à la géographie floristique. Acta Societatis  
819 Botanicorum Poloniae 11, 249-265.
- 820 Tan, M., Qian, W., Pu, J., Liu, H., Zheng, B., 2015. A new approach to develop computer-aided  
821 detection schemes of digital mammograms. Phys. Med. Biol. 60, 4413.
- 822 Tang, D., Liu, Y., Kim, T.-K., 2012. Fast Pedestrian Detection by Cascaded Random Forest with  
823 Dominant Orientation Templates, BMVC, pp. 1-11.
- 824 te Brake, G.M., Karssemeijer, N., Hendriks, J.H., 2000. An automatic method to discriminate  
825 malignant masses from normal tissue in digital mammograms1. Phys. Med. Biol. 45, 2843.
- 826 Varela, C., Tahoces, P.G., Méndez, A.J., Souto, M., Vidal, J.J., 2007. Computerized detection of  
827 breast masses in digitized mammograms. Comput. Biol. Med. 37, 214-226.
- 828 Viola, P., Jones, M., 2001. Rapid object detection using a boosted cascade of simple features,  
829 Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern  
830 Recognition, 2001. CVPR 2001. . IEEE, pp. I-511-I-518 vol. 511.
- 831 Viola, P., Jones, M., 2002. Fast and robust classification using asymmetric adaboost and a detector  
832 cascade. Adv. Neural Inf. Process. Syst. 2, 1311-1318.
- 833 Wang, Y. 2006. *Hierarchical Masses Detection Algorithms Based on SVM in Mammograms*. Master's  
834 Thesis, Xidian University, China.
- 835 Warfield, S.K., Zou, K.H., Wells, W.M., 2004. Simultaneous truth and performance level estimation  
836 (STAPLE): an algorithm for the validation of image segmentation. IEEE Trans. Med. Imaging 23,  
837 903-921.
- 838 Wei, Z.-S., Yang, J.-Y., Shen, H.-B., Yu, D.-J., 2015. A cascade random forests algorithm for  
839 predicting protein-protein interaction sites. IEEE Trans. NanoBiosci. 14, 746-760.
- 840 West, J., Ventura, D., Warnick, S., 2007. Spring research presentation: A theoretical foundation for  
841 inductive transfer. Brigham Young University, College of Physical and Mathematical Sciences.

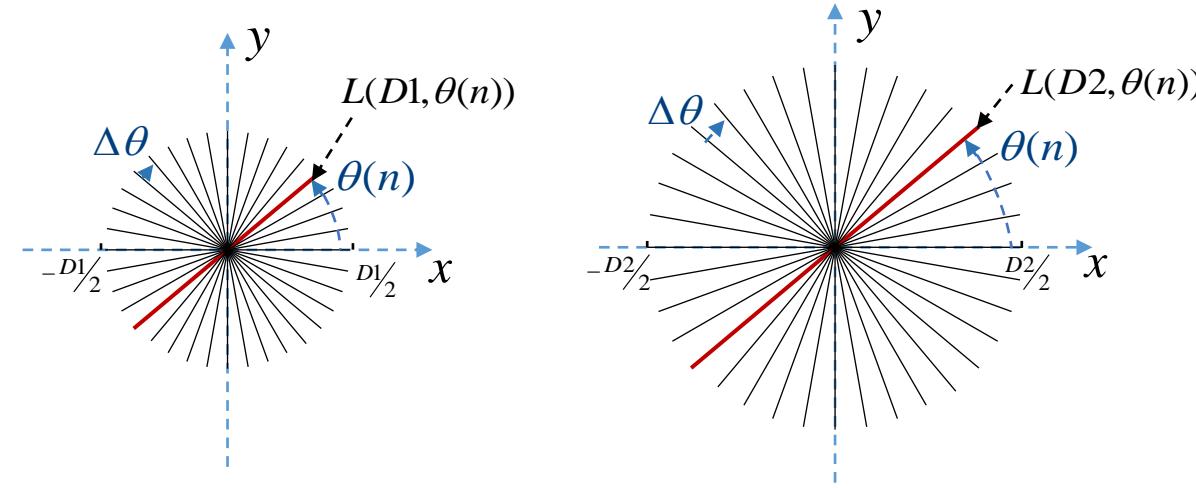
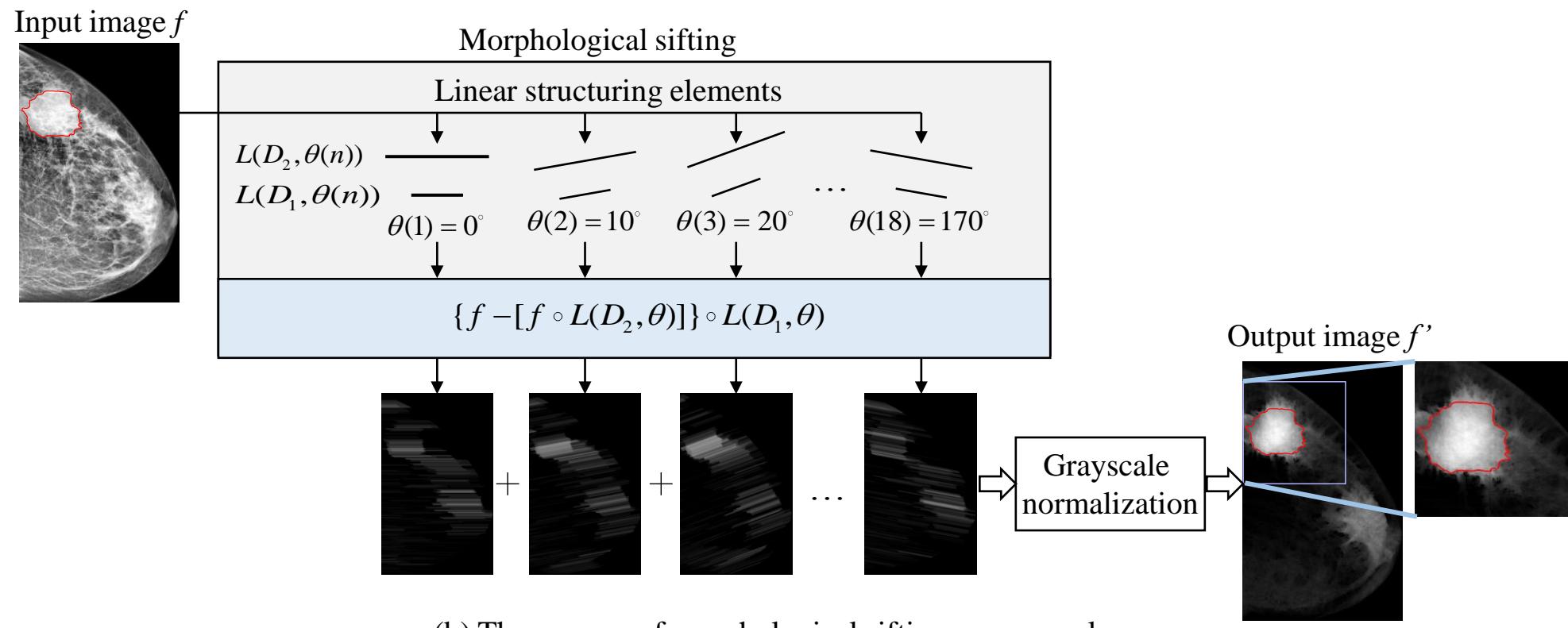
- 842 Zhang, Y.-D., Wang, S.-H., Liu, G., Yang, J., 2016. Computer-aided diagnosis of abnormal breasts in  
843 mammogram images by weighted-type fractional Fourier transform. Advances in Mechanical  
844 Engineering 8, 1687814016634243.
- 845 Zwiggelaar, R., Astley, S.M., Boggis, C.R., Taylor, C.J., 2004. Linear structures in mammographic  
846 images: detection and classification. IEEE Trans. Med. Imaging 23, 1077-1086.

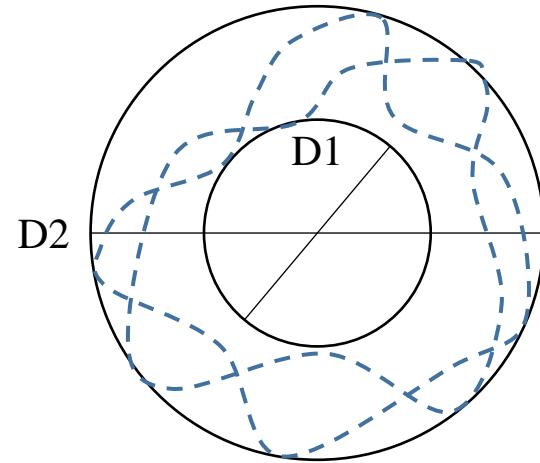
847

Figure

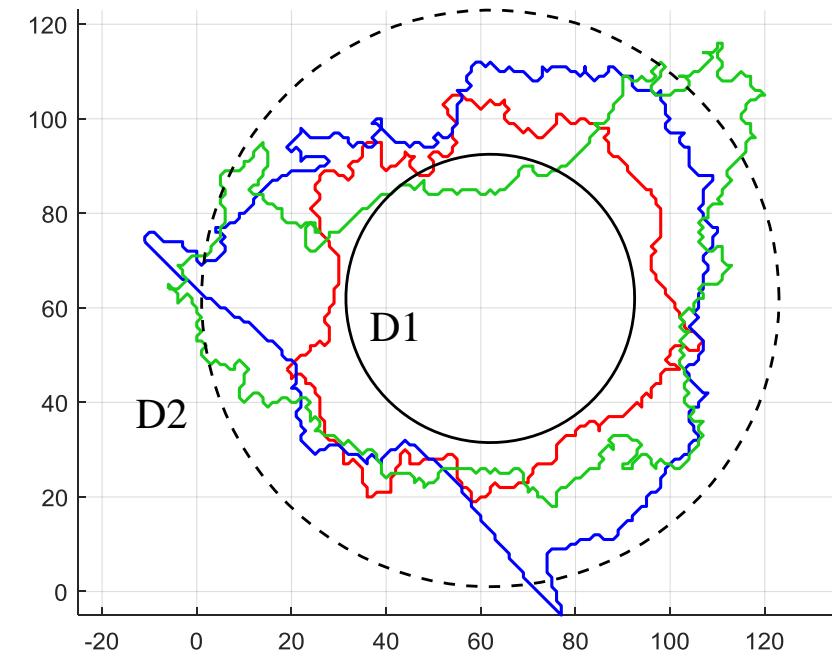


Figure

(a) Two sets of linear structuring elements. The lines marked in red represent a linear element pair  $[L(D_1, \theta(n)), L(D_2, \theta(n))]$ .

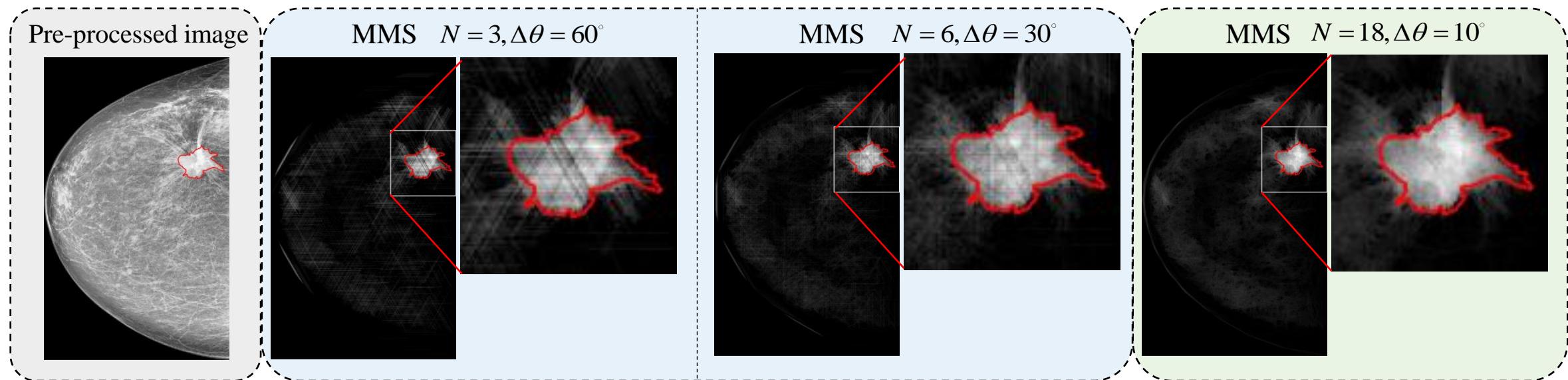


(a) Theoretical model of objects that MMS filters extract.

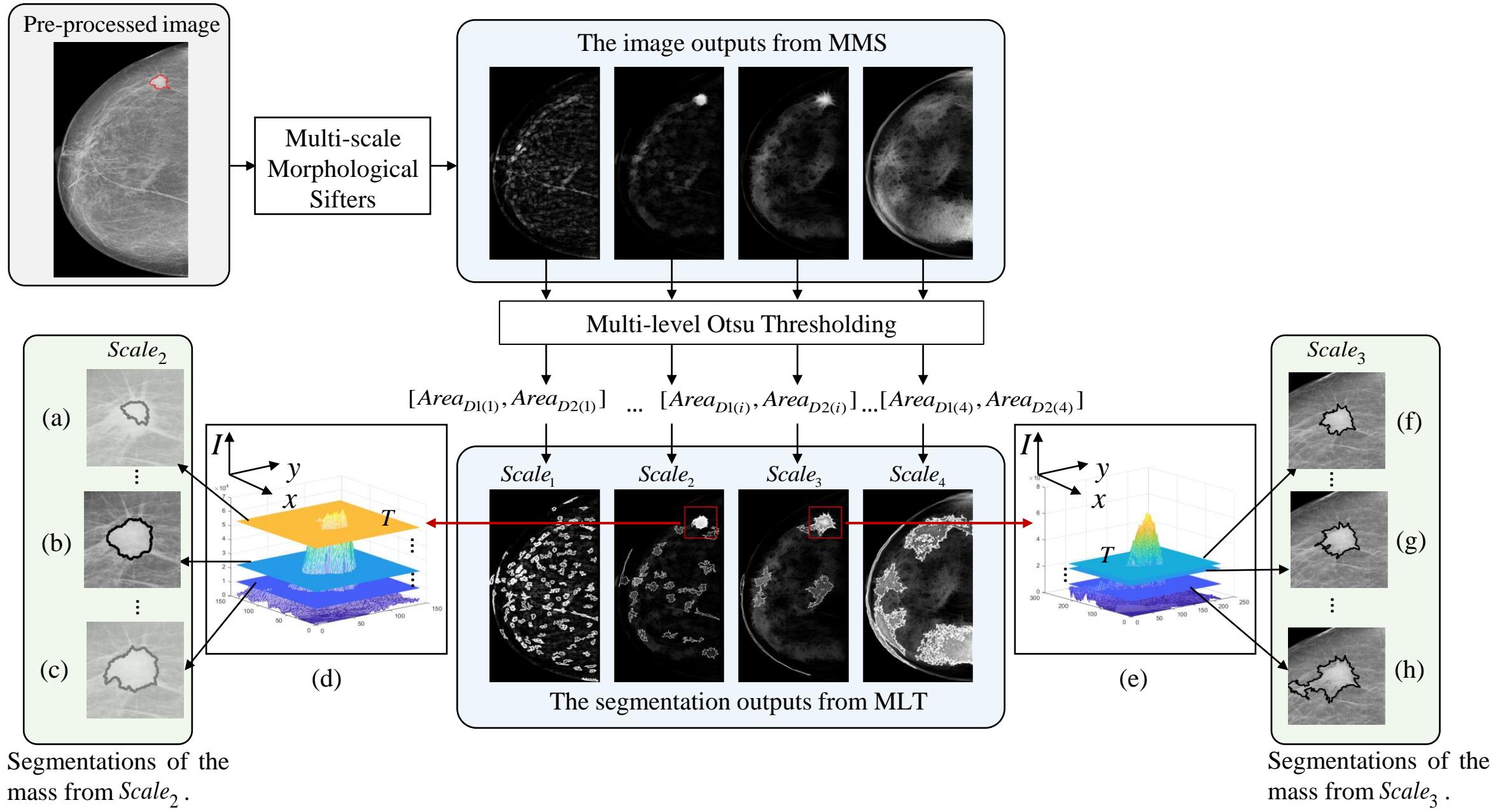


(b) Real masses extracted by MMS filters on one scale (scale3).

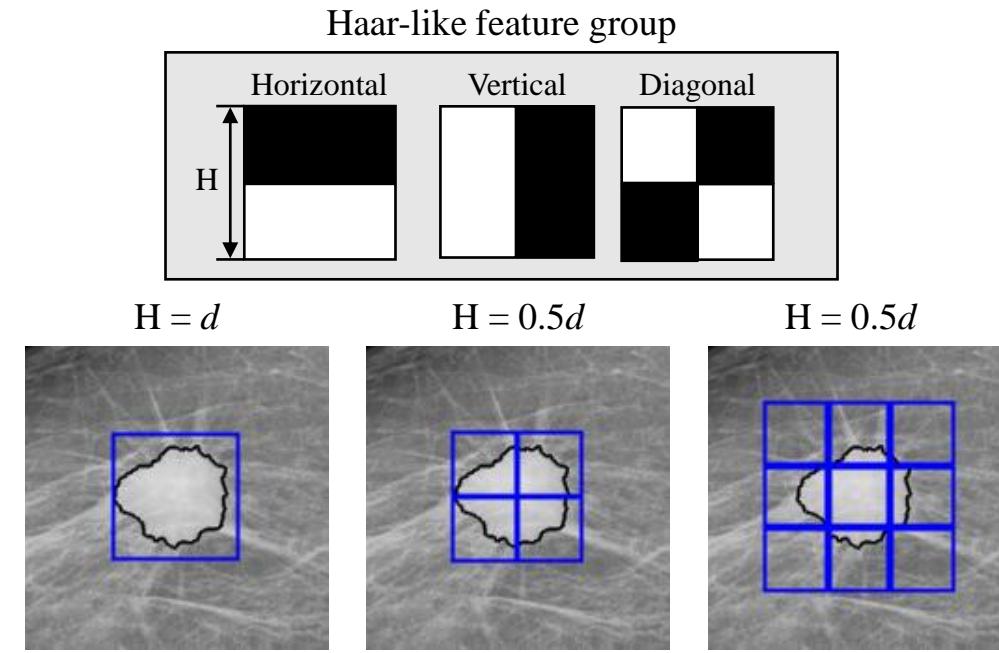
Figure



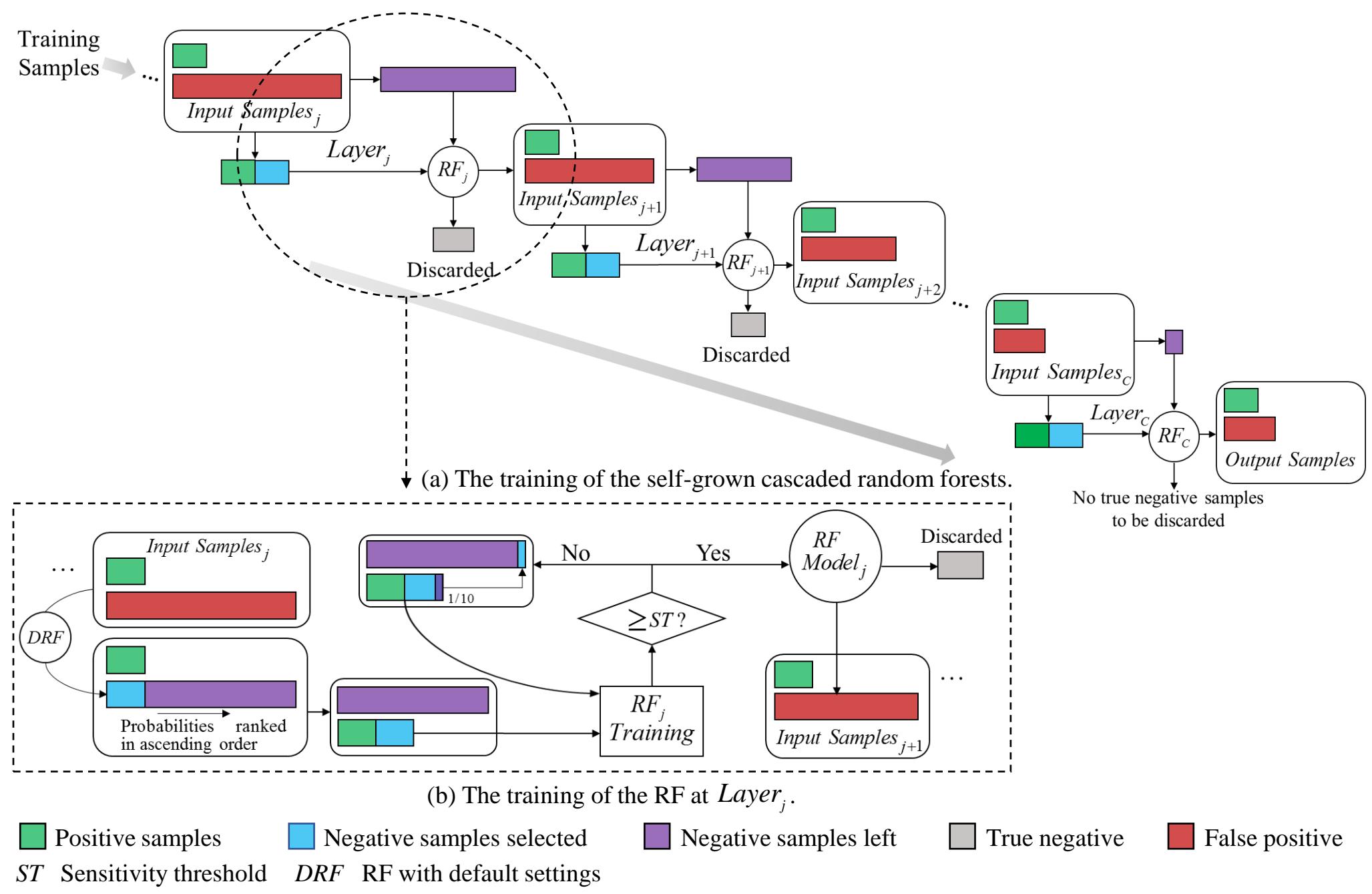
Figure



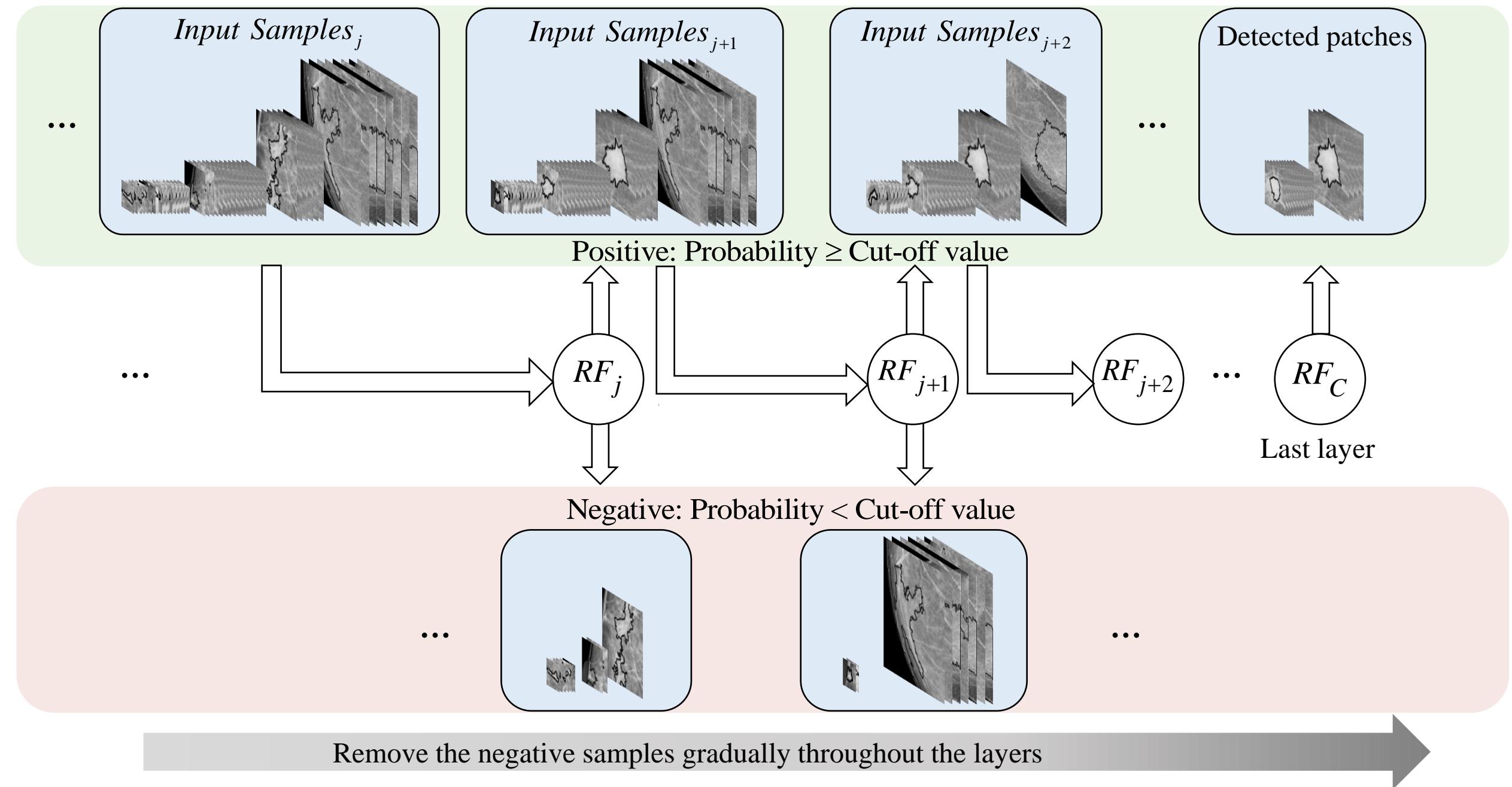
Figure



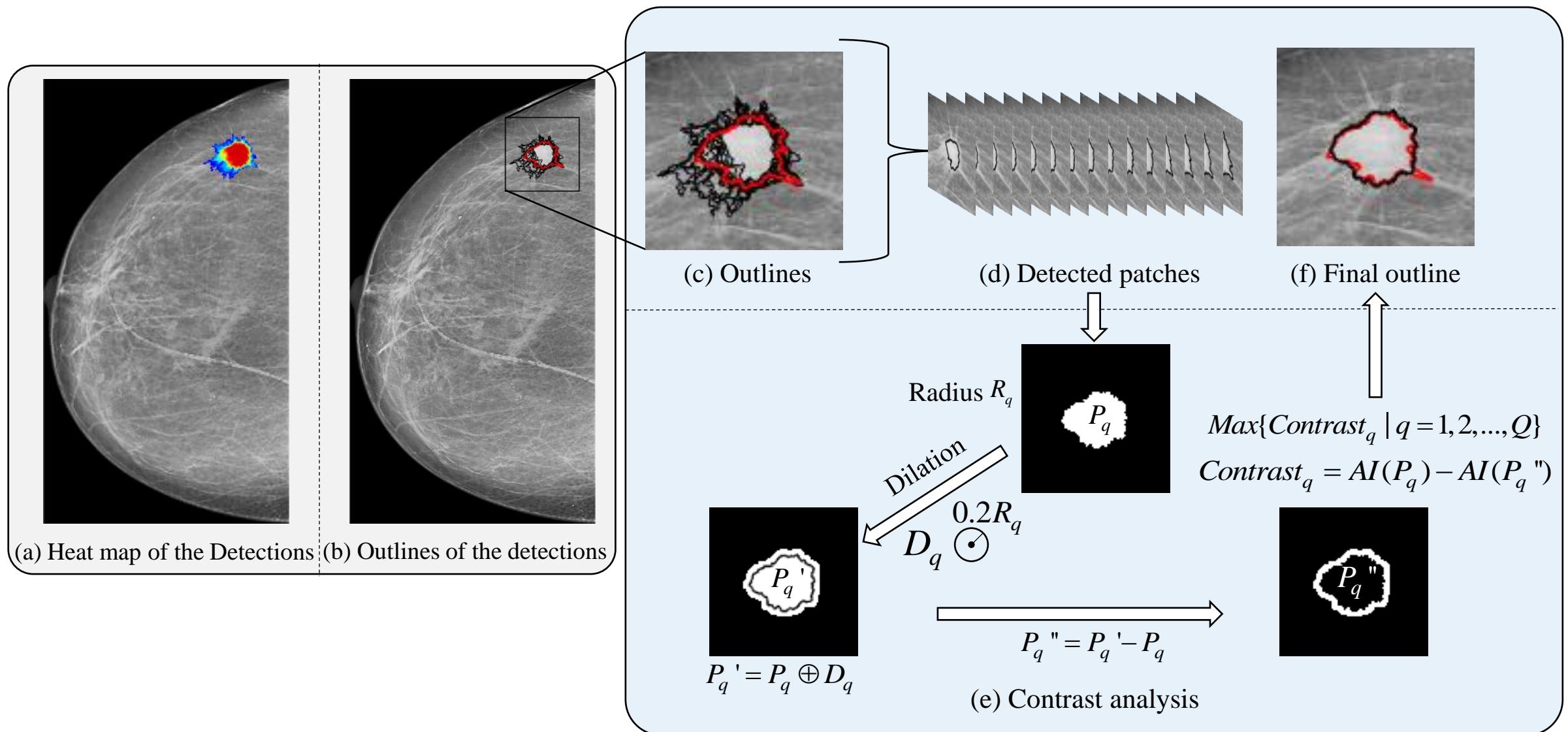
Figure



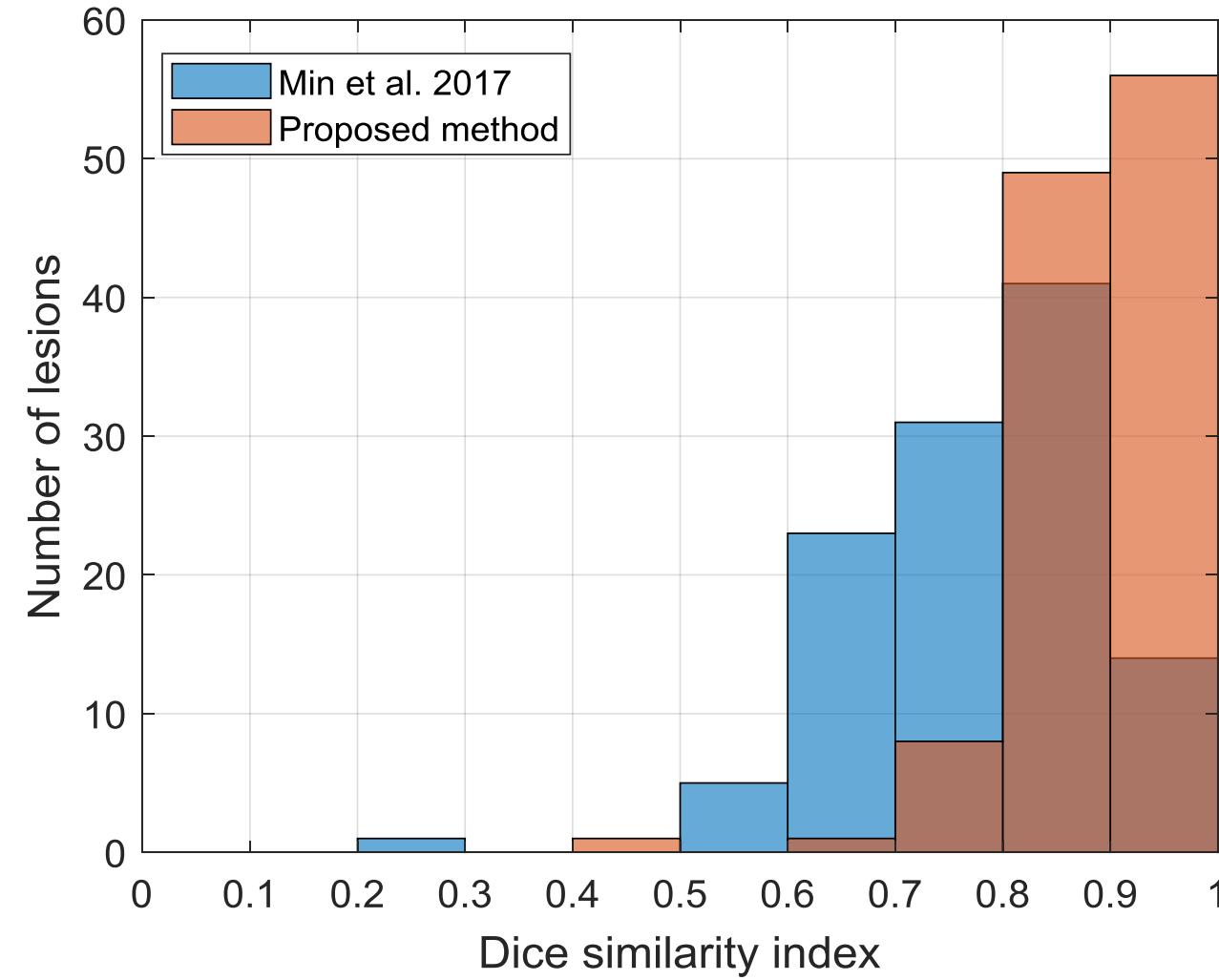
Figure

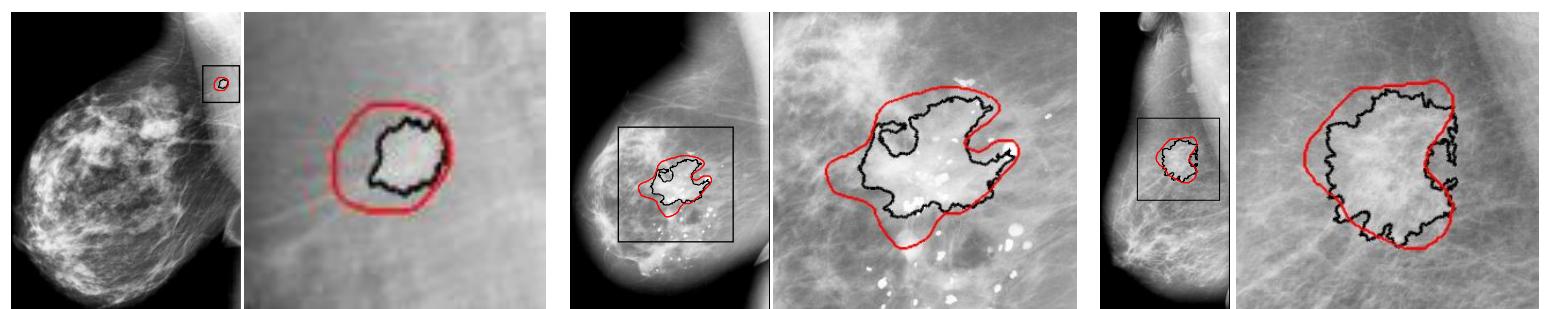
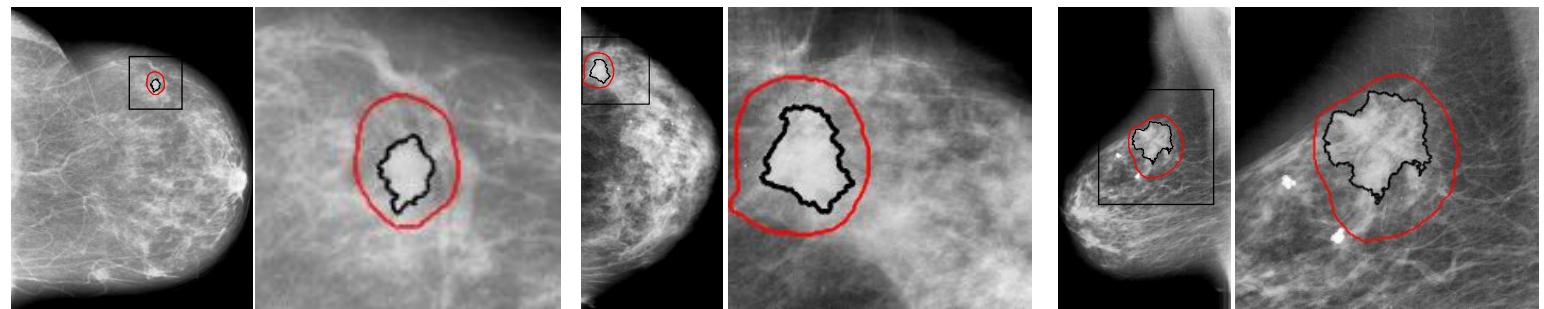


Figure



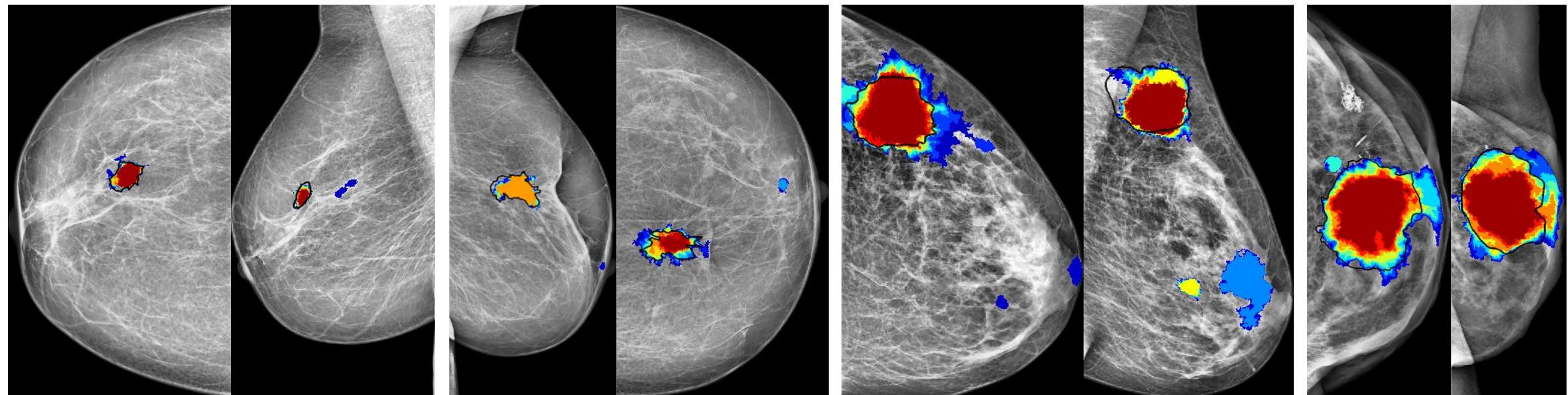
Figure



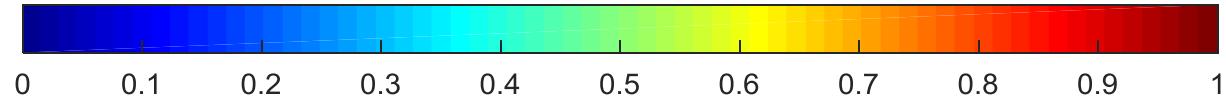


(f) AD, spiculated

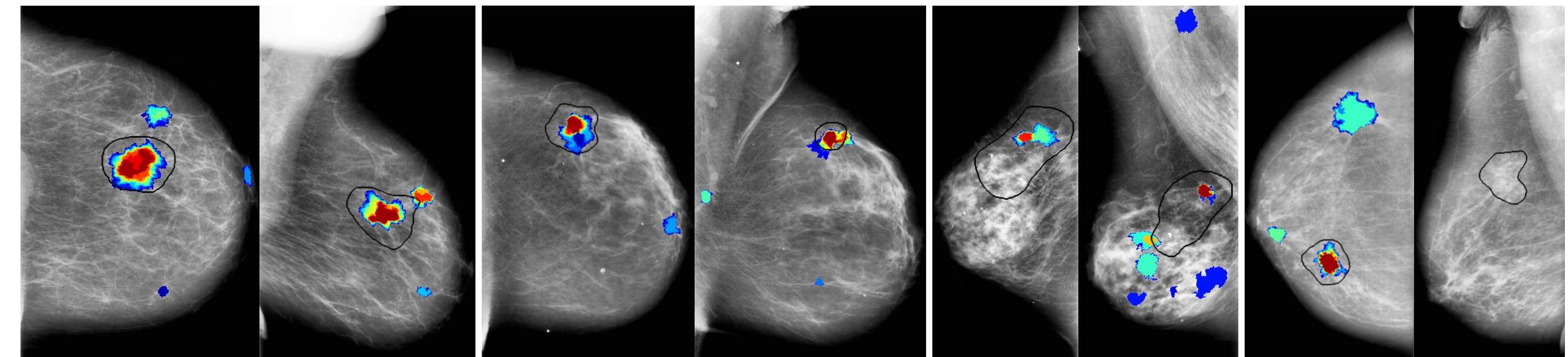
Figure



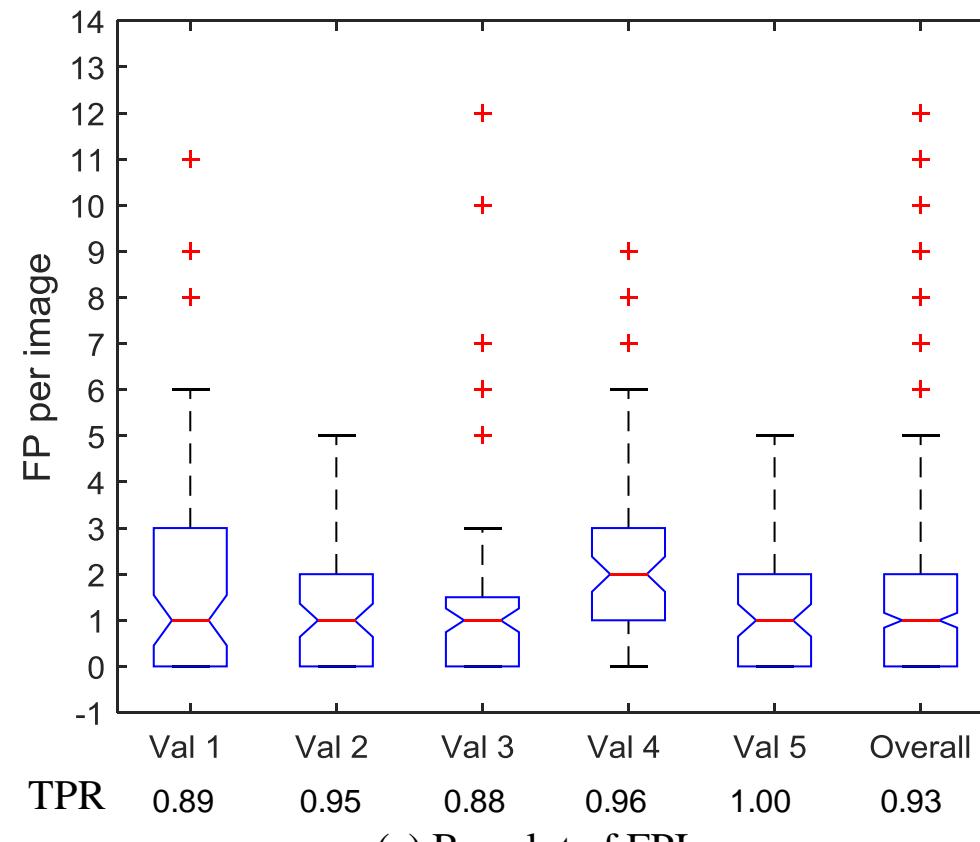
(a) Examples of heat maps from INbreast



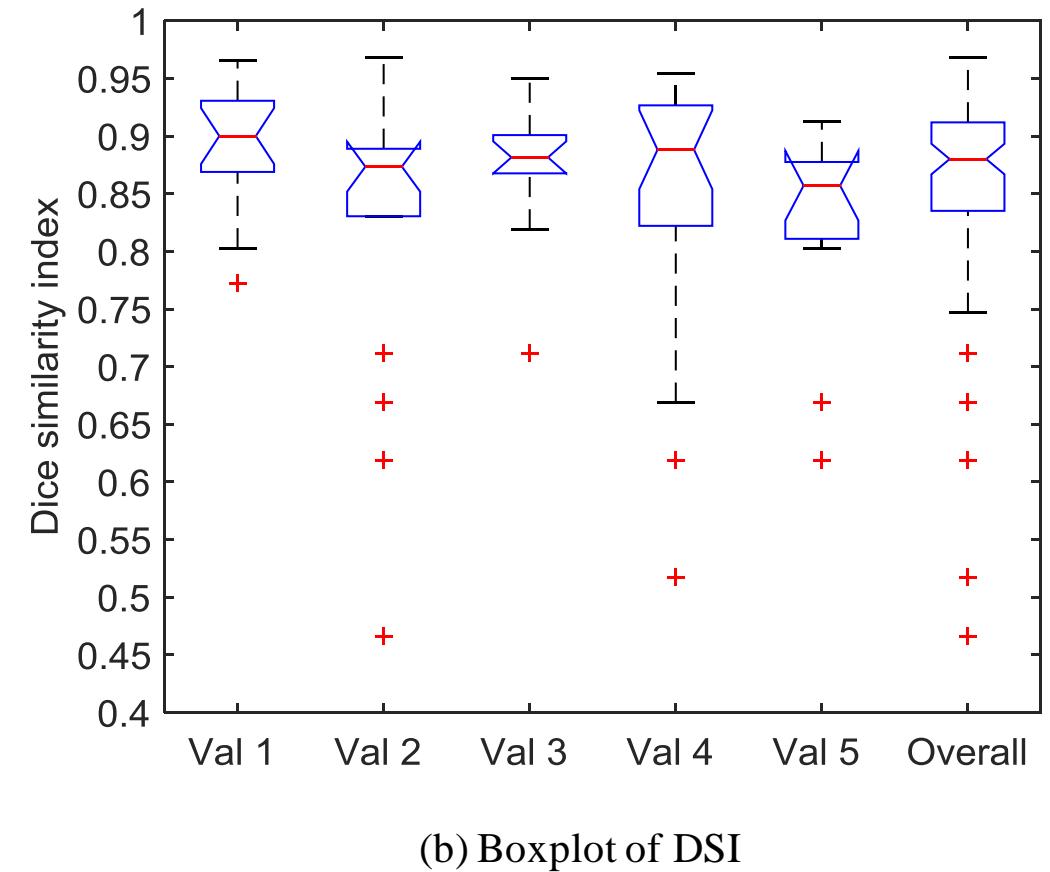
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1



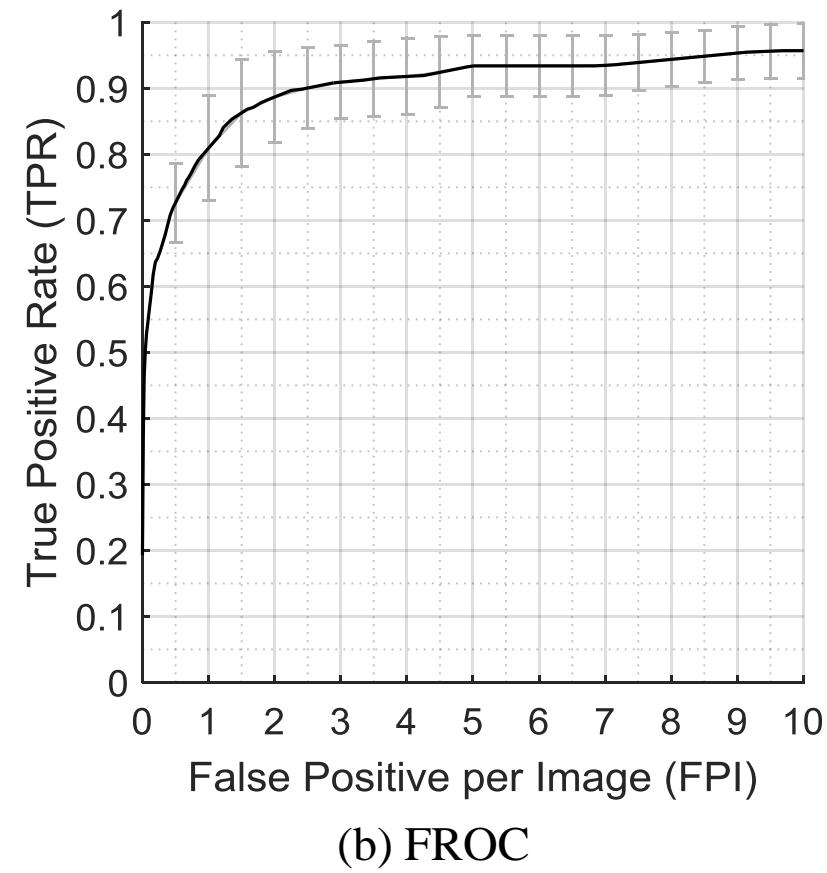
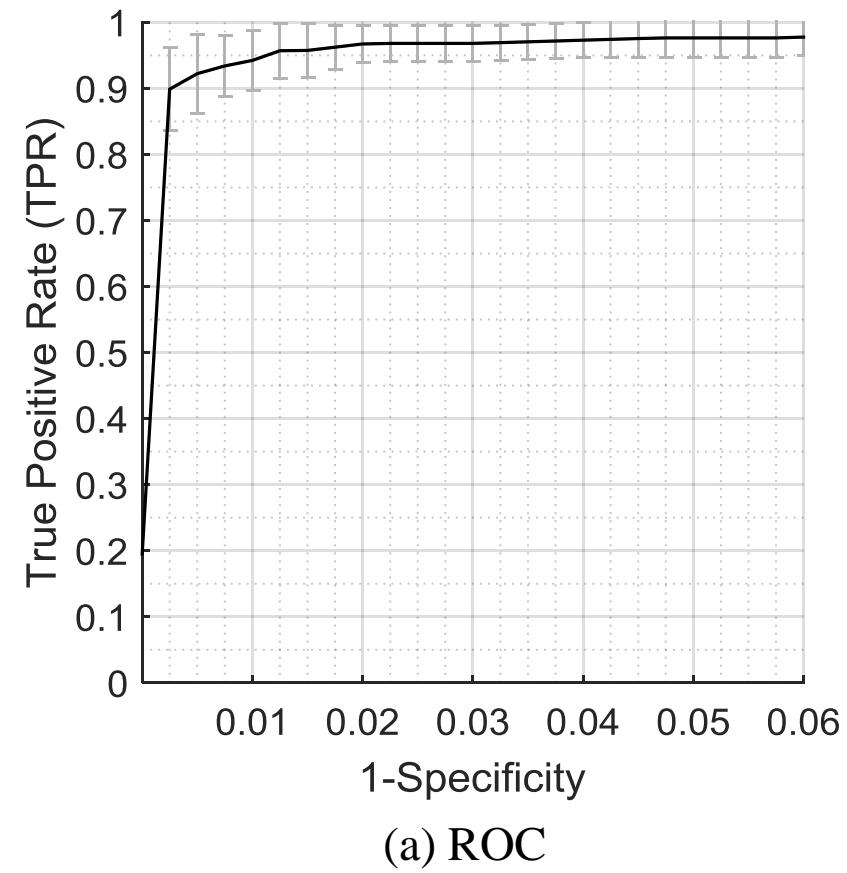
(b) Examples of heat maps from DoD BCRP



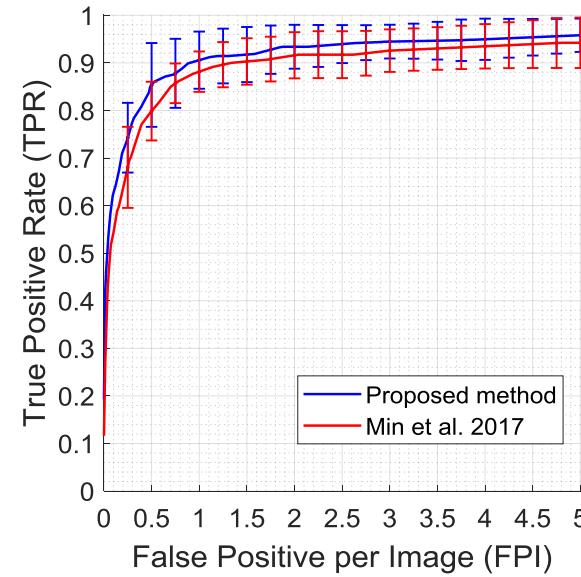
(a) Boxplot of FPI



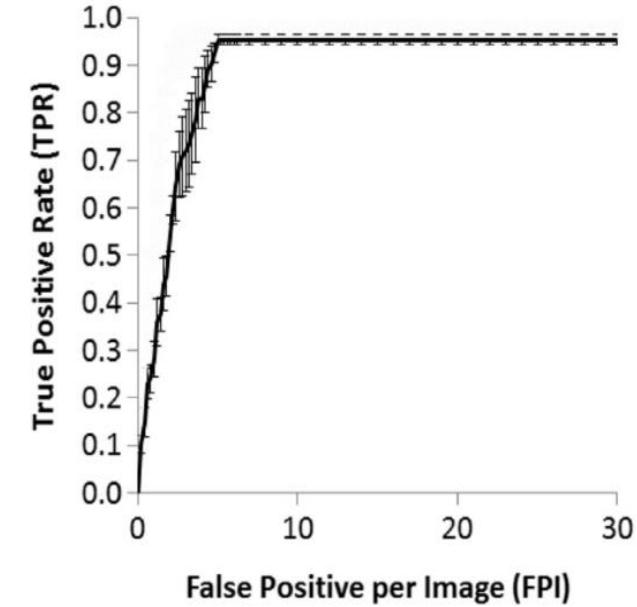
(b) Boxplot of DSI



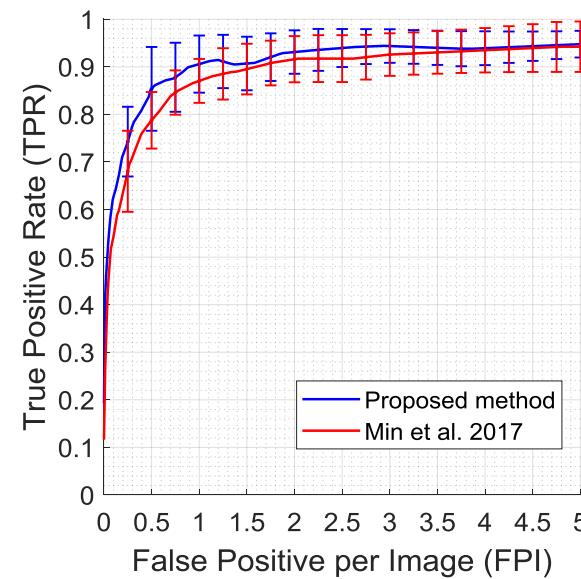
# Figure



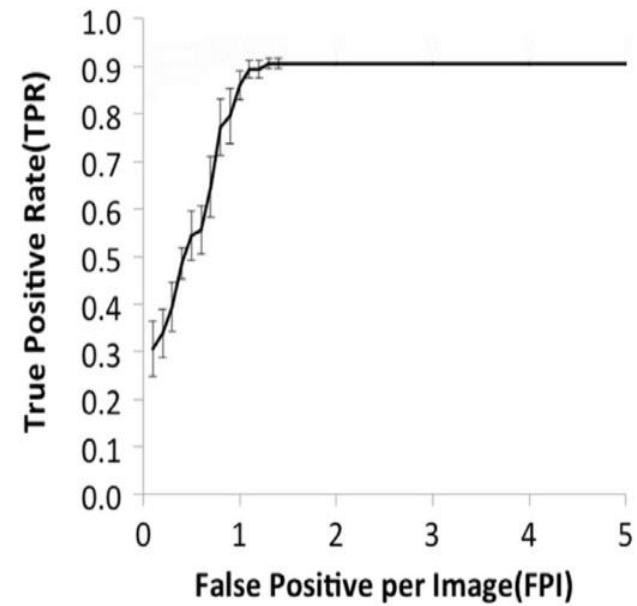
(a) FROC of the proposed method and Min et al. 2017 ( $DT=0.2$ )



(b) FROC of Dhungel et al. 2017 ( $DT=0.2$ )

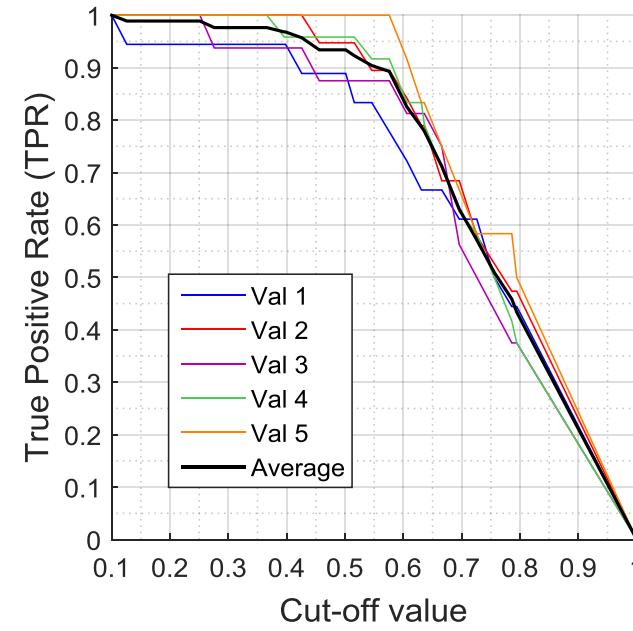


(c) FROC of the proposed method and Min et al. 2017 ( $DT=0.5$ )

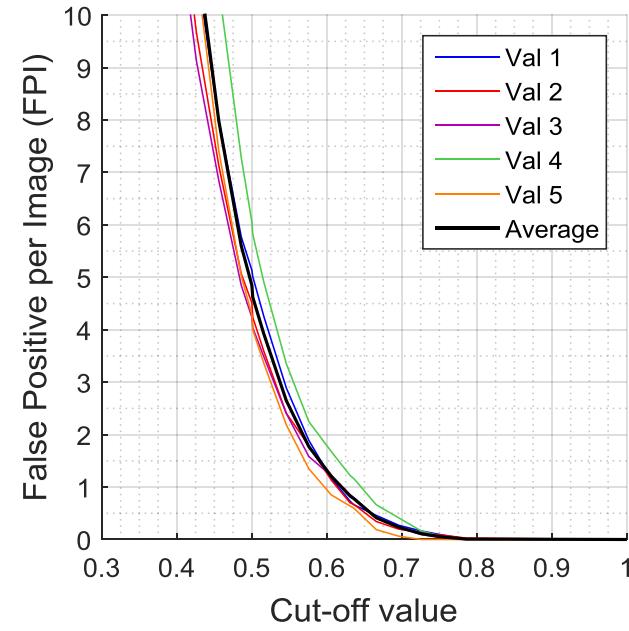


(d) FROC of Dhungel et al. 2017 ( $DT=0.5$ )

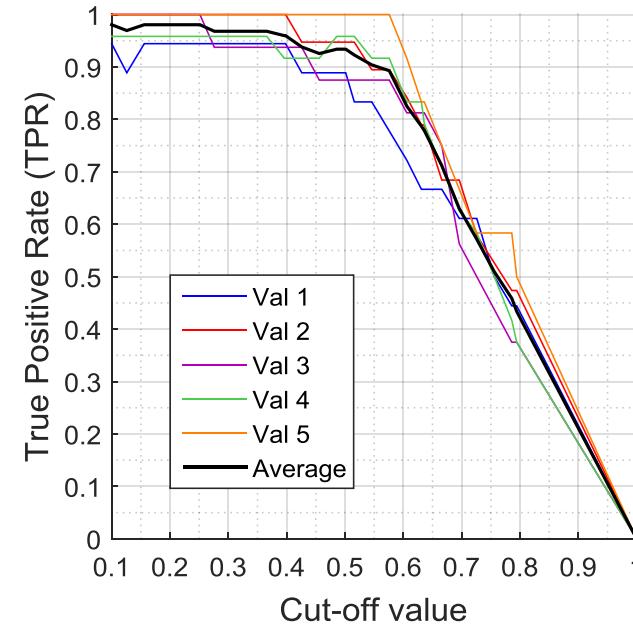
# Figure



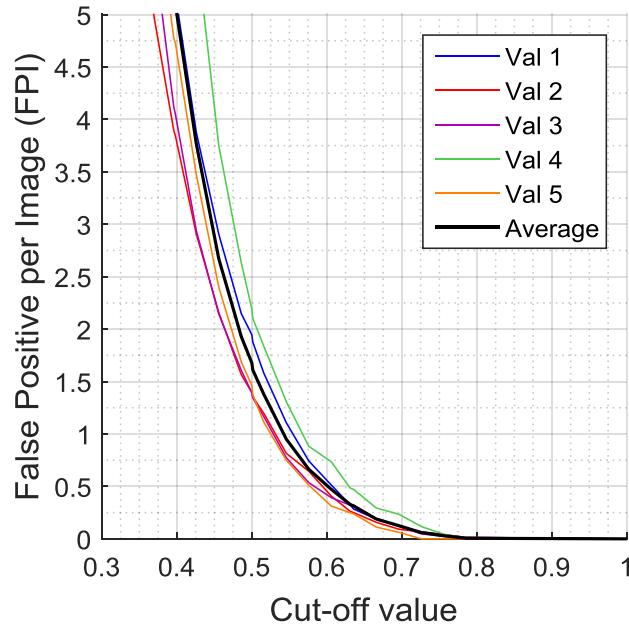
(a) TPR-Cut-off value curve before outline fusion



(b) FPI-Cut-off value curve before outline fusion

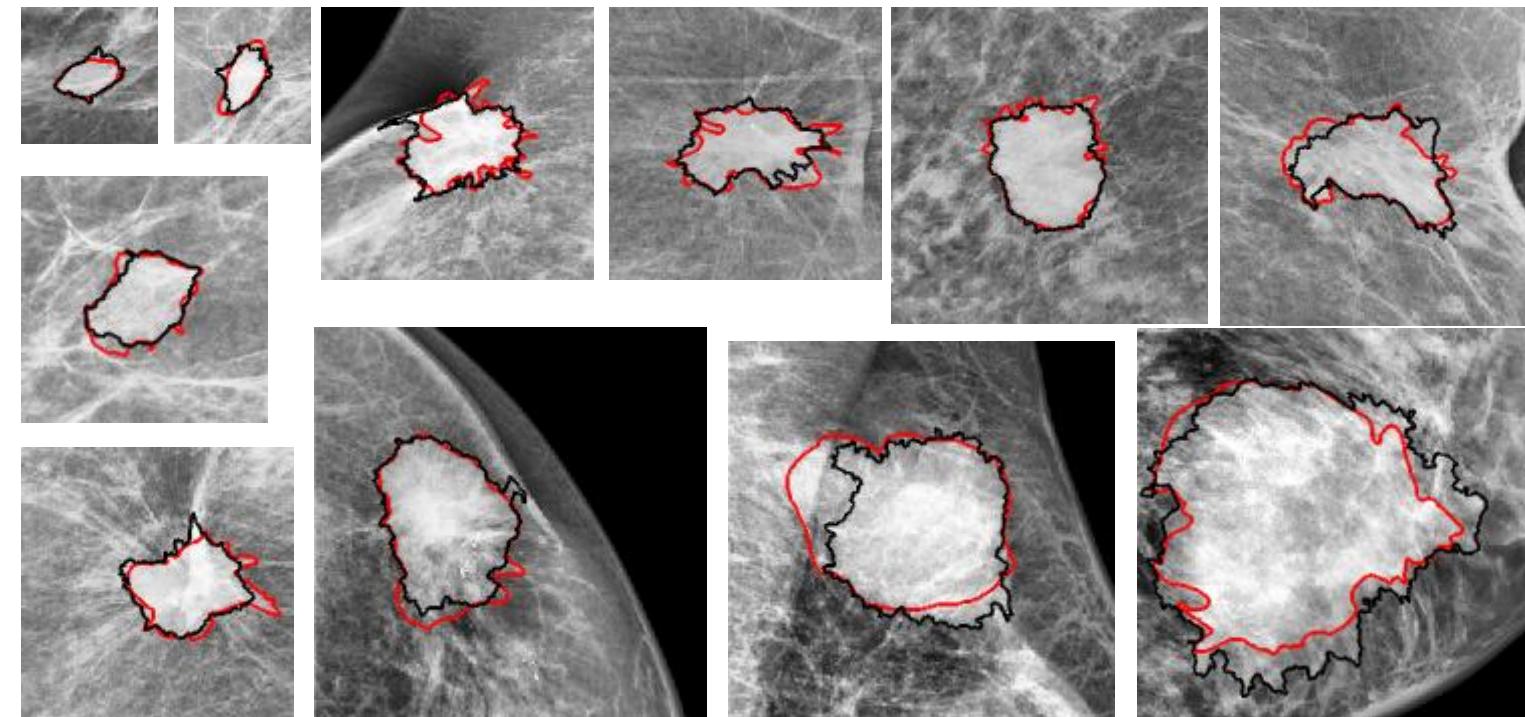


(c) TPR-Cut-off value curve after outline fusion

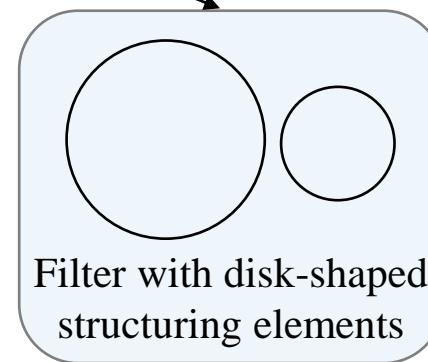
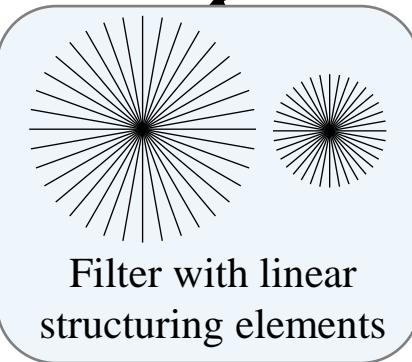
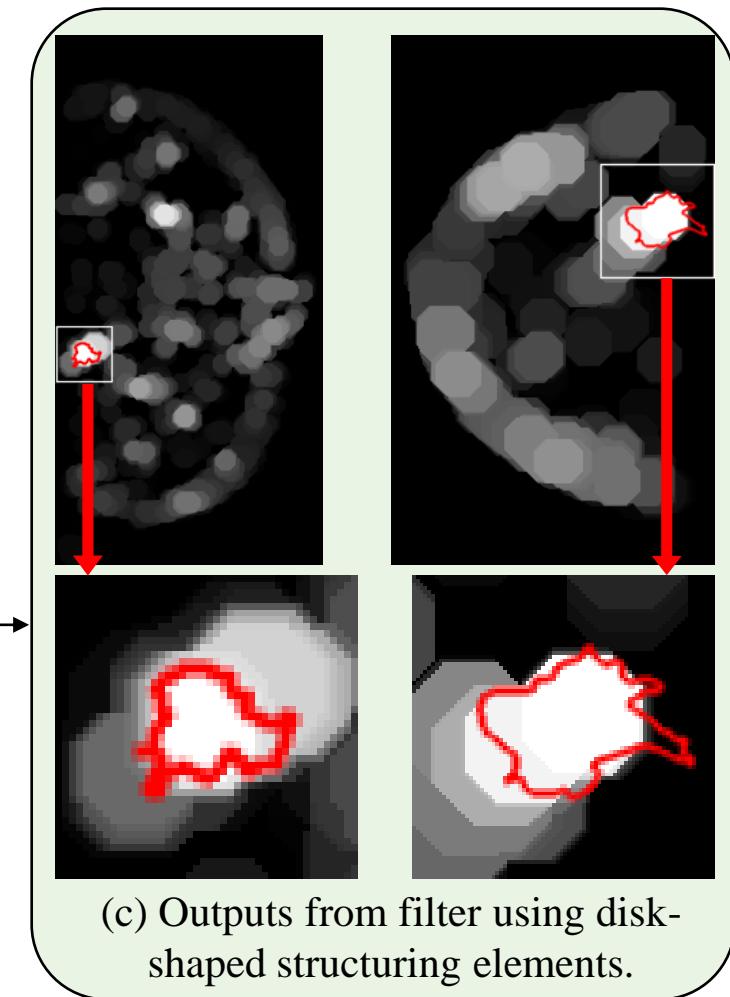
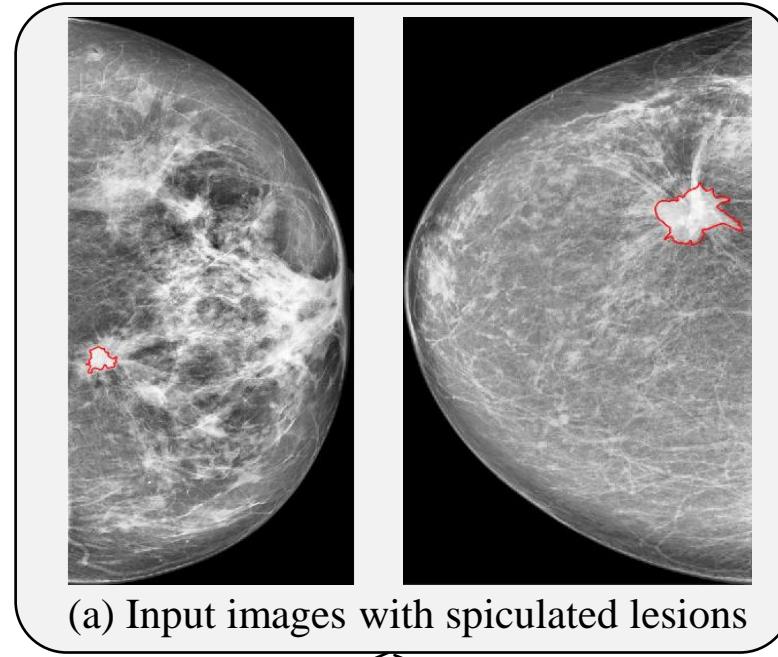
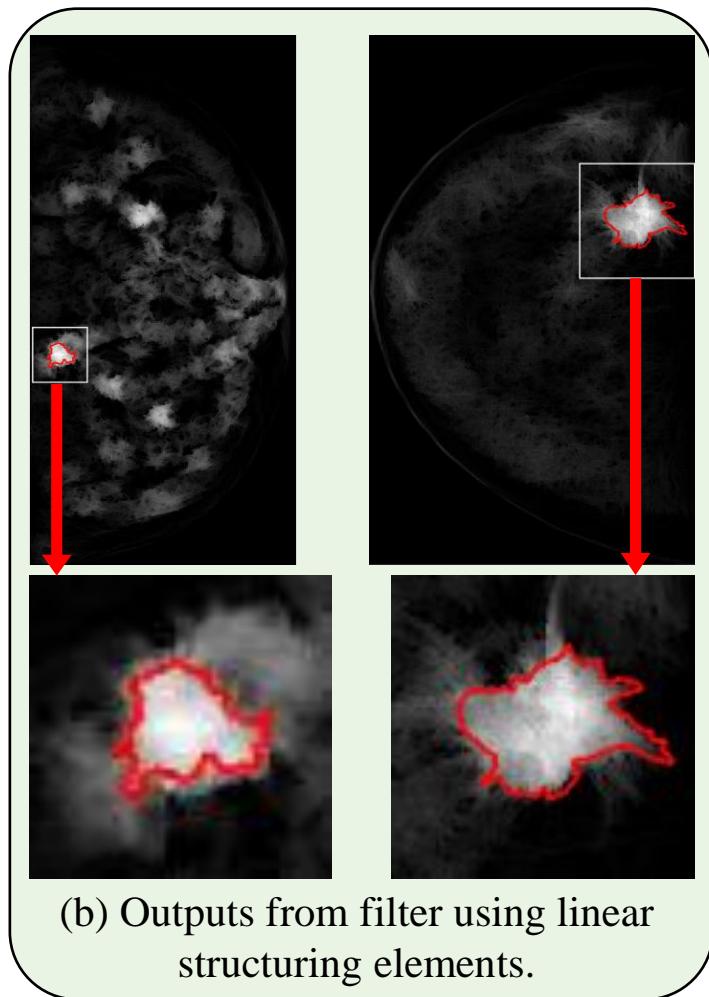


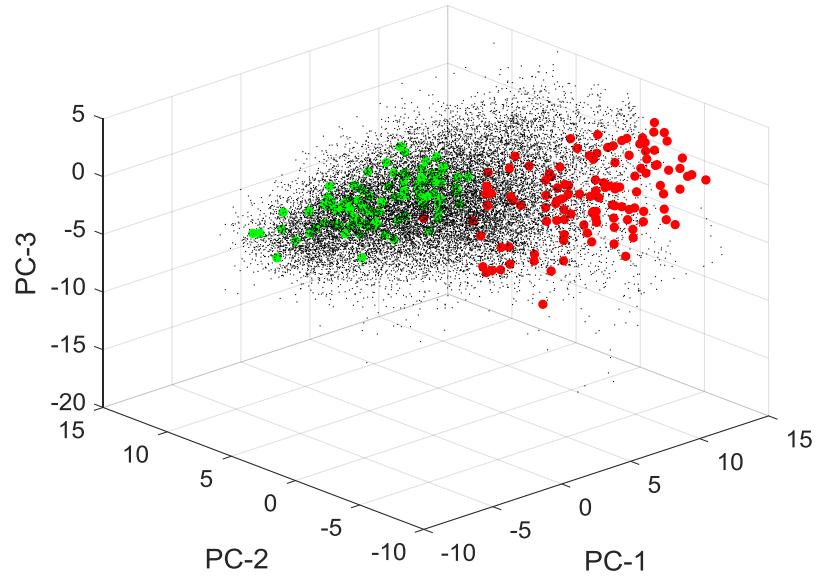
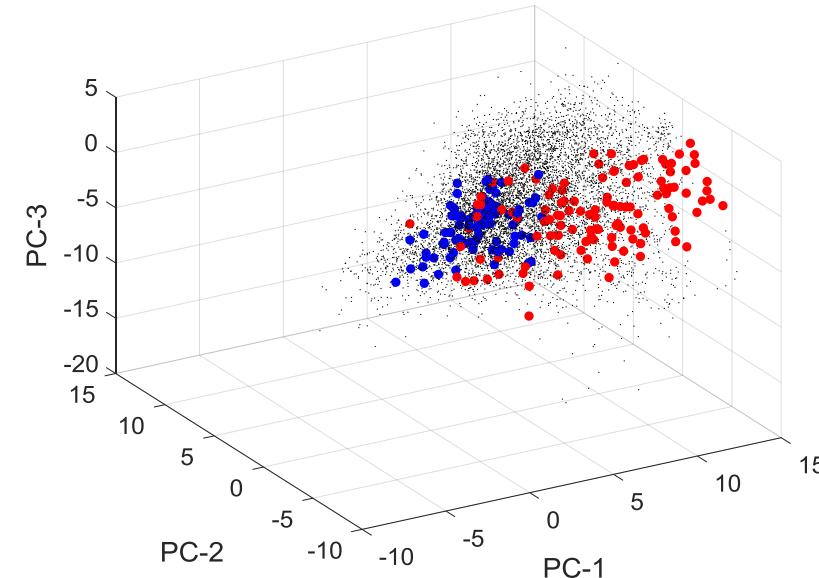
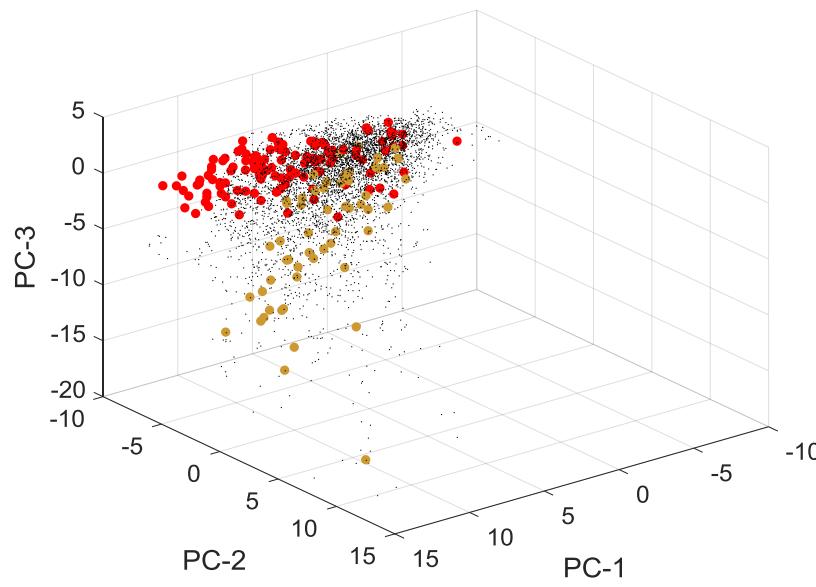
(d) FPI-Cut-off value curve after outline fusion

Figure



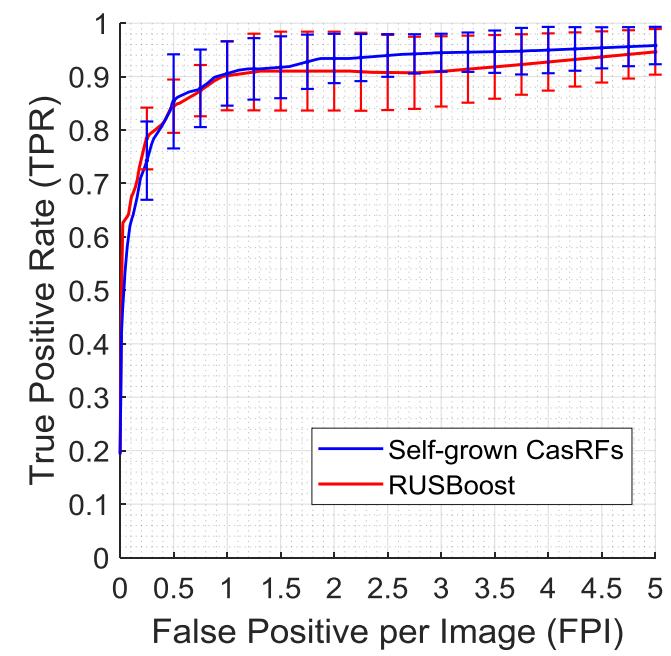
Figure



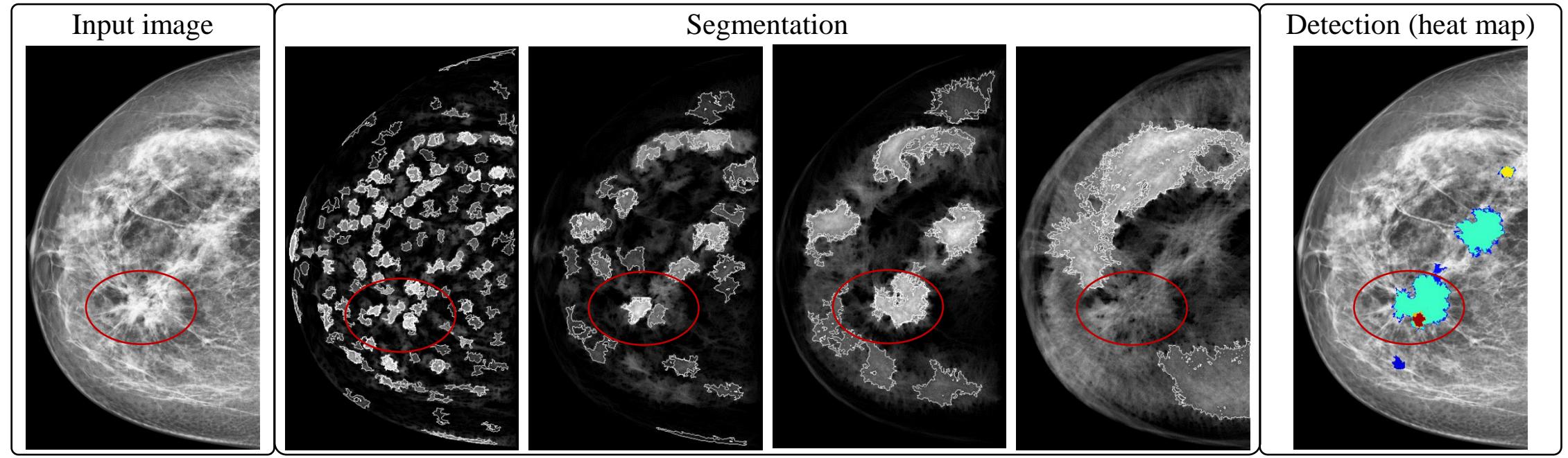
(a) Sample distribution in the 1<sup>st</sup> layer(b) Sample distribution in the 2<sup>nd</sup> layer(c) Sample distribution in the 3<sup>rd</sup> layer

- All negative samples in current layer
- Positive samples
- Negative samples chosen in Layer1
- Negative samples chosen in Layer2
- Negative samples chosen in Layer3

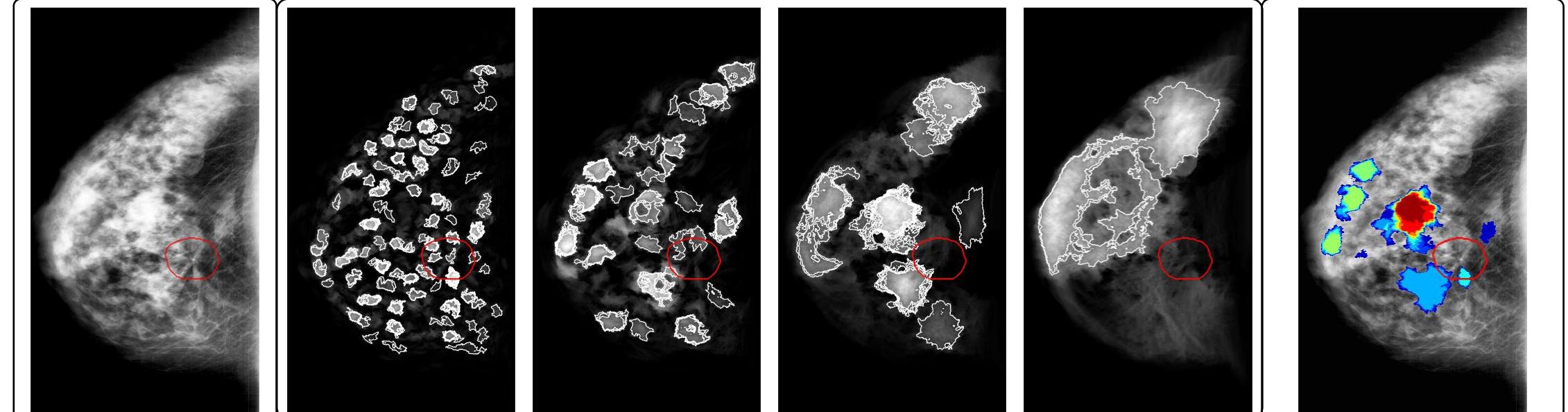
Figure



Figure



(a) A case from INbreast



(b) A case from DoD BCRP

Figure

