

# Collaborative Adversarial Networks for Joint Synthesis and Segmentation of X-ray Breast Mass Images

1<sup>st</sup> Tianyu Shen

State Key Laboratory of Management  
and Control for Complex Systems  
Institute of Automation, Chinese Academy of Sciences  
Beijing, China  
School of Artificial Intelligence  
University of Chinese Academy of Sciences  
Beijing, China  
shentianyu2016@ia.ac.cn

3<sup>rd</sup> Jiangong Wang

State Key Laboratory of Management  
and Control for Complex Systems  
Institute of Automation, Chinese Academy of Sciences  
Beijing, China  
School of Artificial Intelligence  
University of Chinese Academy of Sciences  
Beijing, China  
wangjiangong2018@ia.ac.cn

2<sup>nd</sup> Chao Gou\*

School of Intelligent Systems Engineering  
Sun Yat-sen University  
Guangzhou, China  
gouchao@mail.sysu.edu.cn

4<sup>th</sup> Fei-Yue Wang

State Key Laboratory of Management  
and Control for Complex Systems  
Institute of Automation, Chinese Academy of Sciences  
Beijing, China  
Qingdao Academy of Intelligent Industries,  
Qingdao, China  
feiyue.wang@ia.ac.cn

**Abstract**—In this paper, we propose Collaborative Adversarial Networks (CAN) to enable simultaneous forward synthesis and backward segmentation of X-ray breast mass image. The proposed CAN consists of a generator ( $G$ ), an inverter ( $I$ ) and a discriminator ( $D$ ).  $G$  aims to reconstruct mass images from corresponding annotated masks, while  $I$  is trained for mapping images back to accurate segmentation masks. All the obtained mask-image pairs are fed to  $D$  trained in an adversarial learning scheme. Through the collaborative adversarial training using a joint loss function,  $G$  synthesizes realistic mass images consistent with provided masks and  $I$  effectively segments the tumor regions from the images. Qualitative and quantitative evaluations on publicly available INbreast database demonstrate the effectiveness of our model. Furthermore, different from conventional GANs-based methods that can only perform either image synthesis or segmentation, the proposed model can be generalized to other bidirectional image-to-image translation of multimodal medical data.

**Index Terms**—generative adversarial network, medical image synthesis, mass segmentation, X-ray breast mass

## I. INTRODUCTION

Breast cancer is one of the most common cancers among women worldwide. Digital X-ray mammography is the most

Research supported by National Natural Science Foundation of China under Grant No.61806198, No.61533019 and the Key Research and Development Program of Guangzhou under Grant 202007050002.

\*Corresponding author: Chao Gou.

effective tool for early diagnosis of breast cancer, which is essential for the survival of patients. Currently, deep learning (DL) models have led to a significant breakthrough in automatic detection and segmentation of mammograms, which provides desirable assistance in the accurate diagnosis and clinical treatment. However, for DL models, obtaining precisely annotated training dataset remains a challenging work in medical domain. This is mainly because of the protection of patient privacy, the scarcity of corresponding disease and the expense of expert annotation. Therefore, an effective synthesis of realistic X-ray breast mass appearances and a precise segmentation of tumor regions in mammograms have immediate practical significance.

Nevertheless, synthesizing realistic X-ray breast mass images is still challenging due to the variety of mass in terms of texture and shape as well as the presence of intricate and diverse breast tissue surrounding the masses. Recent development of mass image synthesis can be divided into transformation-based and generative model-based approaches. Transformation-based approaches generate new samples by applying operations such as mathematical affine transformation [1] or feature transformation [2] to the existing samples. However, these small mathematical modifications can only provide a little information alteration and the obtained texture distribution can be inconsistent with real mass. On

the other hand, generative model-based approaches leverage on massive training data to acquire the potential distribution of realistic samples and produce new samples from the same distribution [3]. The advantage is that the synthetic samples are similar in appearance to the real ones but possess a variety of properties different from the existing ones. At present, generative adversarial networks (GANs) [4] along with various stabilization methods produce state-of-the-art results in natural and biomedical image synthesis tasks [5].

The main difficulties in X-ray breast mass segmentation are caused by the intense noise and low contrast of mammogram, as well as strong irregularities and ambiguous boundaries of tumor regions. Traditional methods such as region growing, active contour and Markov random field (MRF) for automatic mass segmentation have been gradually replaced by DL models like convolutional neural networks (CNNs) [1], which overcome the shortcomings of artificial designed features. Some mostly used CNN models, such as fully convolutional network (FCN) [6], SegNet [7] and U-Net [8], have produced state-of-the-art performance in biomedical segmentation tasks.

There is little work on a unified model for joint mask-to-image synthesis and image-to-mask segmentation despite their close correlation. To achieve this goal, we present a novel Collaborative Adversarial Networks (CAN) inspired by two recent researches Pix2Pix network [9] and bidirectional-GAN [10]. The motivation is two-fold. Firstly, we aim to synthesize X-ray breast mass images with pixel-wise masks, different from conventional class-specific synthesis that the generated samples are confined to a specified category. To address this issue, we consider the mask-to-image reconstruction inherited from Pix2Pix network [9]. It is worth noticing that mask-to-image is more versatile compared with label-to-image synthesis in medical domain, since it preserves the semantic information and allows the generalization of the framework to different category fields. Secondly, we aim to achieve accurate segmentation of lesion region at the same time. Therefore, we consider the bidirectional idea from BiGAN [10] to generate the masks from images. Our work is different from [9] and [10] because we not only alter the whole network architecture, but also propose an algorithm of parameter learning. In addition, different from BiGAN [10] that decodes the images into latent representation  $z$  through feature learning, our model is designed to acquire the masks from images.

To summarize, our main contributions include: (1) Propose a novel Collaborative Adversarial Network (CAN) for joint forward synthesis and backward segmentation of X-ray breast mass image, which is capable of being applied beyond mass images, to any bidirectional image-to-image related fields. (2) Enable synthesizing large amount of mask-annotated X-ray breast mass images with variations in appearance. (3) Enable effective segmentation of the tumor regions from the mass images and achieve better performance compared with existing deep segmentation models. (4) Introduce a joint loss function for the collaborative adversarial training of proposed model.

## II. COLLABORATIVE ADVERSARIAL NETWORKS

As shown in Fig.1, the CAN is composed of a generator ( $G$ ), an inverter ( $I$ ) and a discriminator ( $D$ ), and achieves joint synthesis and segmentation tasks through a collaborative adversarial training. In this section, we first briefly review the Pix2Pix and BiGAN model, then introduce the architecture and learning algorithm of proposed CAN in details.

### A. Preliminaries

**Pix2Pix** Pix2Pix network [9] generalizes the cGANs [11] for various image-to-image translation tasks by altering the condition  $c$  to a modality of image  $x$  and adding a traditional L1 loss for ensuring the consistency between expected output  $y$  and input  $x$ . As a result, the generator plays a role in fooling the discriminator, as well as producing output images close to the corresponding ground truth in terms of L1 loss metric. The objective of Pix2Pix can be expressed:

$$\min_G \max_D V(D; G) = \lambda \mathcal{L}_{L1}(G) + \mathbb{E}_{x, y \sim p_{\text{data}}(x, y)} [\log D(x, y)] + \mathbb{E}_{x \sim p_{\text{data}}(x), z \sim p_{\text{data}}(z)} [\log(1 - D(x, G(x, z)))] \quad (1)$$

where  $z$  represents a random input noise vector.

**BiGAN** Bidirectional generative adversarial networks [10] provide a bidirectional mapping between the random noise vector  $z$  and generated image  $y$  by adding an encoder ( $E$ ) into the original  $G - D$  architecture. The  $G$  is defined the same as the original GAN, while  $E$  models a mapping  $E(z; y) : \Omega_Y \rightarrow \Omega_Z$  and induces a distribution  $z \sim p_E(z|y)$ . The  $D$  is modified to take two kind of data pairs,  $(z; G(z))$  obtained from  $G$  and  $(E(y); y)$  from  $E$ , as input to make the predictions  $D(z; G(z))$  or  $D(E(y); y)$  of real or fake.  $G$  and  $E$  are inverse of each other and are trained in an adversarial scheme trying to fool the discriminator. The training objective of BiGAN is defined as a new min-max objective:

$$\min_{G; E} \max_D V(D; G; E) = \mathbb{E}_{y \sim p_y} [\mathbb{E}_{E(y) \sim p_E(z|y)} [\log D(E(y); y)]] + \mathbb{E}_{z \sim p_z} [\mathbb{E}_{G(z) \sim p_G(y|z)} [\log(1 - D(z; G(z)))] \quad (2)$$

### B. CAN Architecture and Learning

In our research, the problems we expect to address are two-fold: (1) How to ensure the authenticity of mass images as well as consistency between masks and images; (2) How to obtain the segmentation mask from mass images simultaneously while achieving better performance compared with a conventional encode-decode segmentation model.

We present the architecture as shown in Fig.1 and a collaborative adversarial learning for solving the above-mentioned problems. Firstly, in order to solve problem (1), we set the segmentation masks as input conditional information, the black pixels of which correspond to the normal breast tissue and white pixels to the lesion region. Thus, the resulted generation is a function over the input mask pixels, which is learned by  $G$ . And we design the

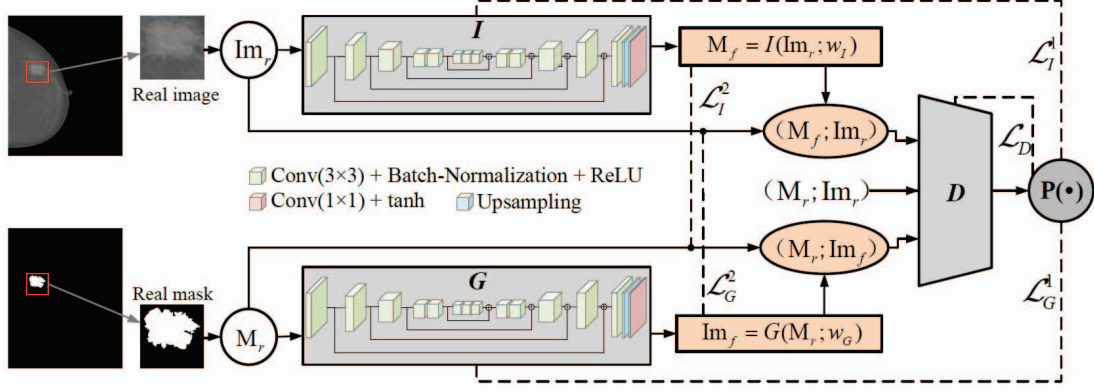


Fig. 1. Proposed CAN architecture.

**Algorithm 1:** The learning of CAN

**Input:**  $(M_r; Im_r)$ , batchsize, epoch, learning rate  $\alpha$ , weight factors  $\lambda_1, \lambda_2$

- 1 **Initialization:**  $\{w_G\}; \{w_I\}; \{w_D\};$
- 2 **while not converge do**
- 3   - mass image synthesis:  $Im_f = G(M_r; w_G);$
- 4   - mass segmentation:  $M_f = I(Im_r; w_I);$
- 5   - discriminative probability for mask-image pairs:  
 $P_r = D(M_r; Im_r; w_D); \quad P_f^1 = D(M_r; Im_f; w_D);$   
 $P_f^2 = D(M_f; Im_r; w_D);$
- 6   -  $G$  learning:  
 $\mathcal{L}_G = \log(1 - P_f^1) + \lambda_1 \mathbb{E}[\|Im_r - G(M_r)\|_1];$   
 $w_G = w_G - \alpha \frac{\partial \mathcal{L}_G}{\partial w_G};$
- 7   -  $I$  learning:  
 $\mathcal{L}_I = \log(1 - P_f^2) + \lambda_2 \mathbb{E}[\|M_r - I(Im_r)\|_2];$   
 $w_I = w_I - \alpha \frac{\partial \mathcal{L}_I}{\partial w_I};$
- 8   -  $D$  learning:  
 $\mathcal{L}_D = \log(P_r) + \log(1 - P_f^1) + \log(1 - P_f^2);$   
 $w_D = w_D - \alpha \frac{\partial \mathcal{L}_D}{\partial w_D};$

discriminator  $D$  to guarantee not only the authenticity of mass images but also the consistency of the output image with the associated mask. Secondly for problem (2), we incorporate an inverter  $I$ , which learns the inverse mapping from mass image samples to condition factors, in the CAN framework.  $I$  is trained jointly with  $G$  and  $D$  using a collaborative adversarial loss function to enable effective learning, as well as improve the generation and segmentation performance through better mode coverage and robustness against mode collapse.

**Network Description** For  $G$  and  $I$  network of the CAN, we choose a U-Net-based [8] architecture, whose capability has been proved on biomedical segmentation tasks. As shown in Fig.1, the architecture is composed of encoding and decoding network with skip connections between mirrored layers. The skip connections enable the encoder and the decoder to share

information. For the problem we consider, a mass image and its corresponding mask differ in surface appearance and texture, but share the same underlying structure and shape. In the conventional encoder-decoder net, all the information passing through the layers directly results in the loss of low-level features.

$D$  is designed as a Patch-based network inspired from [9] with  $16 \times 16$  patch. It only needs to determine if each  $16 \times 16$  patch in an input mask-image pair is real or fake.  $D$  is run across the whole mask-image pair by computing a convolution integral. Then all the responses are averaged to produce the ultimate prediction. This operation assumes each patch to be independent in discrimination process and emphasizes the perception on the authenticity of local texture/style instead of the content of whole image.

**Parameter Learning** As shown in Fig.1,  $G$  takes a mask as input condition  $M_r$  ( $r$  represents *real*) and outputs a generated mass image  $Im_f$  ( $f$  represents *fake*):  $G(M_r; Im_f) : \Omega_M \rightarrow \Omega_{Im}$ . On the contrary,  $I$  outputs a segmentation mask from an input mass image:  $I(Im_r; M_f) : \Omega_{Im} \rightarrow \Omega_M$  and induces a new distribution  $Im_r \sim p_I(Im_r | M_f)$ . Therefore, we construct the mask-image pairs including three forms, which are:  $(M_r; Im_f)$  obtained from  $G$ ,  $(M_f; Im_r)$  from  $I$ , and  $(M_r; Im_r)$  from the matched real mask and image. The obtained mask-image pairs are fed to the  $D$  and a probability  $P$  of real or fake is predicted.

For  $D$ , only the third form  $(M_r; Im_r)$  is real while the other two forms are fake.  $D$  is trained by maximizing the ability of distinguishing real and fake pairs ( $\mathcal{L}_D$ ). For  $G$ , the first form  $(M_r; Im_f)$  is desired to be realistic and the generated image is expected as close as possible to the corresponding ground truth. Thus,  $G$  is learned through maximizing the authenticity, expressed as the adversarial loss term ( $\mathcal{L}_G^1$ ), and minimizing the distance between  $Im_f$  and corresponding  $Im_r$ , expressed as L1 loss term ( $\mathcal{L}_G^2$ ). For  $I$ , similar to  $G$ , the second form  $(M_f; Im_r)$  is desired to be realistic ( $\mathcal{L}_I^1$ ) and the distance between  $M_f$  and corresponding  $M_r$  is expected to be close defined by L2 loss term ( $\mathcal{L}_I^2$ ).

The final joint loss function for the collaborative adversarial



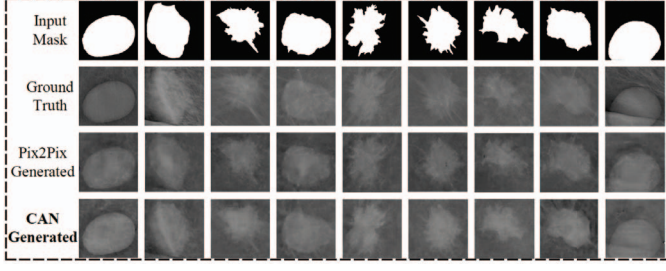


Fig. 2. Visual comparison for real and synthesized images of Pix2Pix and proposed CAN on the INbreast dataset: improvement of the normal synthesized mass images.

training is defined as:

$$\begin{aligned} \mathcal{L}_{CAN}(G; I; D) = & \mathbb{E}_{M, Im \sim p_{data}(M, Im)} [\log D(M, Im)] \\ & + \mathbb{E}_{M \sim p_{data}(M)} [\log(1 - D(M; G(M)))] + \lambda_1 \mathcal{L}_{L1}(G) \\ & + \mathbb{E}_{Im \sim p_{data}(Im)} [\log(1 - D(I(Im); Im))] + \lambda_2 \mathcal{L}_{L2}(I) \end{aligned} \quad (3)$$

Overall, the parameter learning of CAN is summarized in Alg.1. When the final convergence is reached,  $G$  can map a given mask to the corresponding mass image with the realistic distribution:  $Im = G(M; w_G)$ . And  $I$  can acquire the segmentation mask from an input mass image:  $M = I(Im; w_I)$ .

### III. EXPERIMENTS

The INbreast dataset [12] is a public mammogram dataset created by the Breast Research Group, INESC Porto, Portugal. 107 X-ray breast images with their masks are cropped into  $256 \times 256$  pixels size containing the mass ROIs. Original dataset is augmented by rotating and flipping to meet the requirement of a large amount of data. A total of 850 samples is obtained. Thereinto, 700 samples are used for training the model and 150 samples are used for testing. The CAN model is implemented using the Keras framework. The parameter settings are as follows: the batchsize is 1, the training epoch is 300, an Adam optimizer is adopted with a learning rate of 0.0002. The weight  $\lambda_1$  in objective is 10 and  $\lambda_2$  is 100 through experiment and observation. Since there is little work focusing on simultaneous mass image synthesis and segmentation, separate comparisons with existing methods are conducted in this work.

#### A. Mask-to-Image Synthesis

For mask-to-image synthesis, we employ qualitative evaluations and compare the results with Pix2Pix [9] under the same experimental settings. It can be observed that the generated images have shape features resembling to the input masks while their texture features are similar to the ground truth images. We achieve comparable results as Pix2Pix as shown in Fig.2. The reason is that both methods adopt similar architecture and training objective for the generator. However, it is observed that CAN presents some slight improvements at the edge of the synthesized tumor region because of the collaborative optimization of  $D$ . In addition, our model reduces the tissue

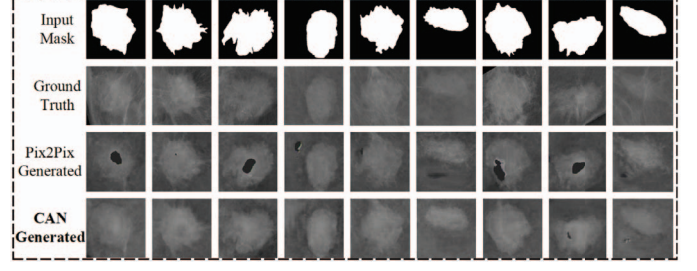


Fig. 3. Visual comparison of Pix2Pix and proposed CAN: improvement of the synthesized mass images with tissue defects.

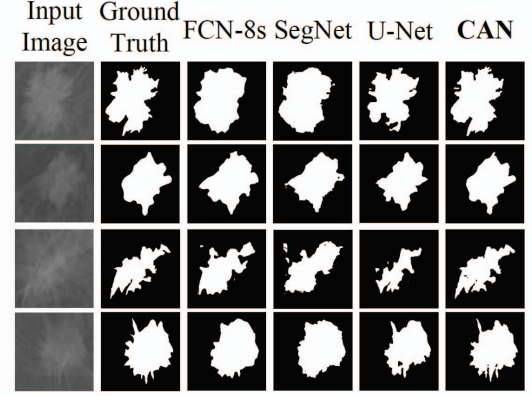


Fig. 4. Qualitative comparison of image-to-mask segmentation.

defects of synthesized images observably as shown in Fig.3, although there are still slight deficiencies in a few ones.

#### B. Image-to-Mask Segmentation

The segmentation performance is evaluated quantitatively as shown in Table I and qualitatively as shown in Fig.4. Our method outperforms other segmentation networks and is capable of obtaining better segmentation boundaries and masks due to the adversarial training for  $I$ .

More specifically, the CAN gives comparable performance with other models in benign mass segmentation, but gives more desirable performance in malignant tumors, especially for the malignant samples with lobulated or spiculated borders. It is more significant for the accurate segmentation of malignant tumor regions in the clinical applications, because the treatment for malignant breast tumor usually requires an operation to remove the lesions.

TABLE I  
QUANTITATIVE EVALUATION FOR IMAGE-TO-MASK SEGMENTATION.

Models	Measurement metrics (%)			
	Acc.	Dice	Jac.	MCC
FCN-8s [6]	93.4	91.0	83.8	85.8
SegNet [7]	93.2	90.4	82.9	85.1
U-Net [8]	93.1	90.5	83.0	85.2
<b>Proposed CAN</b>	<b>94.1</b>	<b>91.7</b>	<b>85.1</b>	<b>87.2</b>

#### IV. CONCLUSION

To summarize, we propose a Collaborative Adversarial Network with joint training for simultaneous end-to-end synthesis and segmentation of X-ray breast mass image. For mask-to-image synthesis, we achieve desirable results with fewer defects in synthesized images using a more generalized framework, as compared to Pix2Pix network. For image-to-mask segmentation, our model outperforms other segmentation models due to the addition of adversarial loss. Furthermore, the proposed model can be generalized beyond mass images to other bidirectional image-to-image translation fields.

#### REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] L. Maaten, M. Chen, S. Tyree, and K. Weinberger, "Learning with marginalized corrupted features," in *International Conference on Machine Learning*, 2013, pp. 410–418.
- [3] D. Nie, R. Trullo, J. Lian, C. Petitjean, S. Ruan, Q. Wang, and D. Shen, "Medical image synthesis with context-aware generative adversarial networks," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2017, pp. 417–425.
- [4] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [5] E. Wu, K. Wu, D. Cox, and W. Lotter, "Conditional infilling gans for data augmentation in mammogram classification," in *Image Analysis for Moving Organ, Breast, and Thoracic Images*. Springer, 2018, pp. 98–106.
- [6] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [7] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [8] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [9] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [10] J. Donahue, P. Krähenbühl, and T. Darrell, "Adversarial feature learning," *arXiv preprint arXiv:1605.09782*, 2016.
- [11] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [12] I. C. Moreira, I. Amaral, I. Domingues, A. Cardoso, M. J. Cardoso, and J. S. Cardoso, "Inbreast: toward a full-field digital mammographic database," *Academic radiology*, vol. 19, no. 2, pp. 236–248, 2012.