# MAMBO-NET: Multi-Causal Aware Modeling Backdoor-Intervention Optimization for Medical Image Segmentation Network

Ruiguo Yu[a,b,c,d], Yiyang Zhang[d], Yuan Tian[a,b,c], Yujie Diao[a,b,c], Di Jin[a], Witold Pedrycz[f,g,h,*]

[a]*College of Intelligence and Computing, Tianjin University, Tianjin, 300350, China*
[b]*Tianjin Key Laboratory of Cognitive Computing and Application, Tianjin, 300350, China*
[c]*Tianjin Key Laboratory of Advanced Networking, Tianjin, 300350, China*
[d]*Tianjin International Engineering Institute, Tianjin University, Tianjin, 300350, China*
[e]*Tianjin Central Hospital of Gynecology Obstetrics, Tianjin, 300100, China*
[f]*Department of Measurement and Control Systems, Silesian University of Technology, Akademicka 2 Gliwice, 44-100, Poland*
[g]*Department of Electrical and Computer Engineering, University of Alberta, Edmonton, AB, T6G 2R3, Canada*
[h]*Research Center of Performance and Productivity Analysis, Istinye University, Istanbul, 34010, Turkiye*

## Abstract

Medical image segmentation methods generally assume that the process from medical image to segmentation is unbiased, and use neural networks to establish conditional probability models to complete the segmentation task. This assumption does not consider confusion factors, which can affect medical images, such as complex anatomical variations and imaging modality limitations. Confusion factors obfuscate the relevance and causality of medical image segmentation, leading to unsatisfactory segmentation results. To address this issue, we propose a multi-causal aware modeling backdoor-intervention optimization (MAMBO-NET) network for medical image segmentation. Drawing insights from causal inference, MAMBO-NET utilizes self-modeling with multi-Gaussian distributions to fit the confusion factors and introduce causal intervention into the segmentation process. Moreover, we design appropriate posterior probability constraints to effectively train the distributions of confusion factors. For the distributions to effectively guide the segmentation and mitigate and eliminate the impact of confusion factors on the segmentation, we introduce classical backdoor intervention techniques and analyze their feasibility in the segmentation task. Experiments on five medical image datasets demonstrate a maximum improvement of 2.28% in Dice score on three ultrasound datasets, with false discovery rate reduced by 1.49% and 1.87% for dermatoscopy and colonoscopy datasets respectively, indicating broad applicability.

*Keywords:* Medical Image Segmentation, Causal Inference, Backdoor Model, Gaussian Modeling

## 1. Introduction

Medical image segmentation is a vital component of computer-aided diagnosis, playing a key role in assisting clinicians with treatment planning and decision-making. In recent years, deep learning methods such as UNet [1] and its variants [2, 3, 4] have successfully been employed in various segmentation tasks in pathology and imaging modalities, including colon polyp segmentation, skin lesion segmentation, and breast nodules analysis. These methods demonstrate the ability to classify regions of interest or lesion pixels accurately. However, they still face challenges in achieving high accuracy and reliability in pixel classification at the boundaries of the segmented areas. As illustrated in Fig. 1, different segmentation methods exhibit varying degrees of under-segmentation and over-segmentation at the boundary.

Many existing segmentation methods utilize networks that directly couple feature extraction with segmentation prediction.
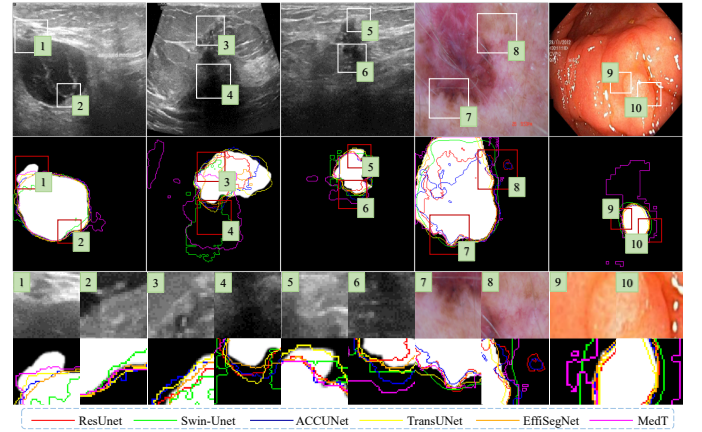


Figure 1: The segmentation results of multiple networks. The third row corresponds to a zoomed-in image of the selected area. The curves of different colors in the second and last rows represent the corresponding models' predictions of the lesion area.

However, this approach introduces confusion factors where the features become intertwined with target information, resulting in biased model predictions. Causal learning theory [5] suggests that modeling and incorporating causal interventions can

---

*Corresponding author
*Email addresses:* `rgyu@tju.edu.cn` (Ruiguo Yu),
`zyy_0203@tju.edu.cn` (Yiyang Zhang), `tiany@tju.edu.cn` (Yuan Tian),
`yujiediao@tju.edu.cn` (Yujie Diao), `jindi@tju.edu.cn` (Di Jin),
`wpedrycz@ualberta.ca` (Witold Pedrycz)
[1]These authors contributed equally to this work.

alleviate the negative impact of confusion factors. For instance, JointMCMC [6] proposes a Monte Carlo sampling embedding model that probabilistically describes noise patterns and reduces speckle noise in ultrasound image lesion segmentation. SwinHR [7] demonstrates effective modeling of kinetic features at different acquisition times to attenuate the influence of lesion heterogeneity on segmentation. While causal inference methods help mitigate confusion factors, existing approaches often rely on manual hierarchical modeling [8, 9] or oversimplify by focusing on specific causal concepts, neglecting other potential critical factors. Additionally, some abstract concepts within the confusion factors (e.g., fuzzy boundaries or underlying pathological features) are complex to describe quantitatively, thus increasing the complexity of modeling.

To comprehensively model confusion factors and mitigate their influence on the segmentation process, we propose the Multi-causal Aware Modeling Backdoor-intervention Optimization for Medical Image Segmentation Network (MAMBO-NET). For complete modeling of confusion factors, MAMBO-NET introduces a latent space modeling approach based on Gaussian self-modeling to address the challenges posed by complex confusion factors. Through intervention optimization, the method effectively reduces the confounding effects on model decision-making, thereby achieving significant improvements in segmentation performance. Specifically, our contribution can be summarised as follows:

- We classify the adverse factors that affect segmentation as confusion factors and employ the construction of a structured causal diagram to analyze their underlying mechanisms.

- The implicit modeling approach of abstract and concrete confusion factors using Gaussian Self-modeling is proposed. This approach enables the self-modeling of a more comprehensive latent space of confusion concepts, eliminating the need for dependency on existing manual confusion modeling methods.

- To reduce the bias of confusion factors on segmentation decisions, we propose combining implicit modeling with the segmentation process through backdoor interventions.

- Our method achieved a maximum segmentation index improvement of 2.28% on three ultrasound datasets, and reduced lesion area false alarm rates by 1.49% and 1.87% in the dermatoscopy and colonoscopy datasets, respectively, indicating that the method has broad applicability.

## 2. Related Work

### 2.1. Latent Space Modeling

Latent space modeling effectively characterizes high-dimensional lesion features and their implicit relationships to enhance segmentation performance. Current approaches predominantly follow three paradigms: feature decoupling, prototype learning, and probabilistic modeling. Feature decoupling methods, exemplified by KDFD [10] for content-style separation and CDDSA [11] for contrast domain disentanglement, isolate confusion factors from target features, though potential information loss remains a limitation. Li et al. [12] proposed a foundational work in image representation, which has been widely used in recent decades. It achieved exceptional robustness in handling outliers, having a profound impact on the development of image processing.

Prototype-based techniques like SSMIS [13] with boundary-aware prototypes and PROCNS [14] using progressive calibration offer more compact representations through feature clustering, yet face challenges in capturing complex confusion relationships. Probabilistic approaches, including GMM-SDF [15] for morphological variations and BayeSeg [16] for domain adaptation, provide flexible distribution modeling, though their potential for segmentation tasks warrants deeper investigation.

### 2.2. Medical Image Segmentation Based on Causal Inference

Recent advances in causal inference have spurred its integration into medical image analysis, primarily focusing on causal interpretation, intervention, and counterfactual prediction.

Causal interpretation methods analyze relationships between anatomical structures and predictions. C-CAM [17] constructs interpretable causal models by examining category-anatomical causality, while P-CSS [18] models disease co-occurrence effects to mitigate bias in radiology reports. Causal intervention addresses prediction bias through frontdoor or backdoor adjustments. CaMIL [19] reduces spurious disease-color associations in WSI classification via frontdoor adjustment, whereas CausalCLIPSeg [20] suppresses confusion bias in segmentation using backdoor adjustment. Counterfactual prediction frameworks estimate potential outcomes under interventions. Wang et al. [21] synthesize counterfactual mammographic features, while Richens et al. [22] demonstrate how counterfactual reasoning improves diagnostic decisions. Medical images contain entangled confusion factors, complicating their separation in high-dimensional feature space. We propose a latent space modeling approach with backdoor adjustment to mitigate their adverse effects on segmentation.

## 3. Method

### 3.1. MAMBO-NET Architecture

MAMBO-NET proposes a causal intervention framework with three core components: causal relationship modeling, Gaussian-based implicit modeling (GSm), and bias-aware intervention (CIBM). As illustrated in Fig. 2, the architecture integrates: (1) the UNeXt [4] segmentation backbone(Encoder and Decoder), (2) GSm module with Gaussian Distribution Extraction Backbone (GDEB) and posterior constraint backbone (PCB) for confusion factors modeling, and (3) CIBM for intervention fusion. GSm employs reparameterization to sample prior distributions while PCB provides training-phase constraints. CIBM performs weighted fusion of bias-mitigated features across decoder stages. The framework implements
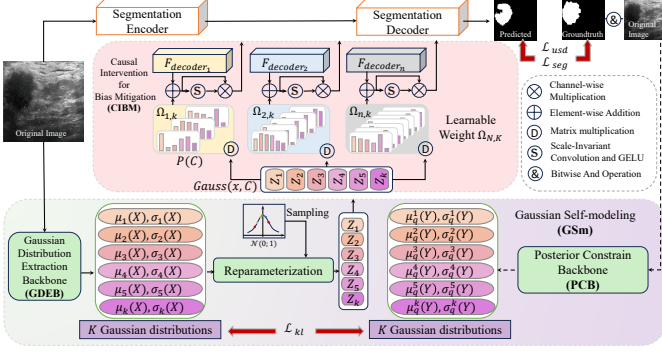
Figure 2: The architecture of the proposed **MAMBO-NET**. Dashed lines indicate that the data stream is disabled during model inference. **Segmentation Encoder** and **Segmentation Decoder** are the encoder and the decoder in UNeXt. **Gaussian Backbone** and **Posterior Constrain Backbone** will use the global average pooling(GAP) and linear mapping for scale alignment.

backdoor adjustment to suppress confusion effects through: explicit causal modeling of hidden confusion factors (Sec. 3.2), GSm's learnable Gaussian priors with posterior regularization (Sec. 3.3), and CIBM's adaptive feature conditioning (Sec. 3.4).

### 3.2. Causal Relationship Modeling

The backdoor model is one of the most classical models in causal relationships. Previous works have attempted to incorporate the backdoor model into medical imaging tasks and model and intervene in the confusion factors that affect model decisions, such as imaging artifacts and scattered noise. The modeling of confusion factors in these works is based on hierarchical assumptions made by humans about known concepts. Therefore, the confusion factors usually only include visible and concrete concepts. However, in this paper, we unify the known and unknown, abstract and concrete factors that affect segmentation boundary decisions into confusion factors, which are distributed at the segmentation boundaries. For example, the comet tail sign in ultrasound images is a known, concrete concept, while factors such as physician acquisition habits are unknown, abstract concepts.
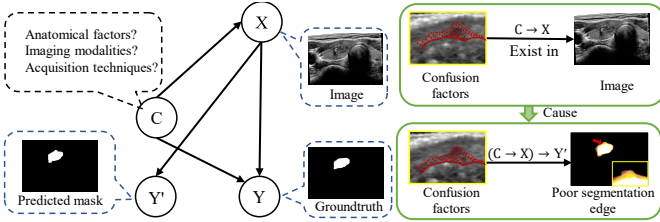


Figure 3: The process of causal relationship modeling. X denotes the original image; Y denotes the predicted mask, and c denotes the confusion factors.

Fig. 3 represents the causal backdoor model, where $X$ denotes the input image and $Y'$ represents the segmentation mask produced by the network. The variable $C$ represents the confusion factors described in this paper. Since the confusion factors exist within the acquired medical image $X$, we consider that the confusion factors have a causal effect on the acquired medical

image $X$. Simultaneously, due to the complex causal relationship between the segmentation mask $Y'$ and the complex confusion factors, the confusion factors C also have a causal effect on the segmentation mask $Y'$. Traditional segmentation models utilize a prediction process where the encoder extracts feature representations $\phi(X; \theta)$ from the medical image $X$, and the decoder generates the corresponding segmentation mask $Y'$, i.e., the process of obtaining $P(Y'|X)$. However, these models typically overlook the impact of confusion factors $C$. As $C$ becomes the confounder acting on both $X$ and $Y'$, the learned $P(Y'|X)$ by the network becomes biased.

To mitigate the bias introduced by the confounder C on the network's learned decisions and enforce the network to learn an unbiased decision rule $P(Y'|do(X))$, where the do operator represents causal intervention, we aim to adjust the backdoor using a hierarchical modeling of the confusion factors $C$, which can be expressed as Eq. (1).

$$P(Y'|do(X)) = \sum_{c_i}^{C} P(Y' = y'|X = x, C = c_i)P(C = c_i) \quad (1)$$

where $x$ represents the current image characteristics and $y'$ represents the prediction result of the network. $c_i$ is a concept of the set of confusion factors in the current image.

### 3.3. Gaussian-based implicit modeling

The confusion factors summarized by existing methods often only cover observable and known concepts, oversimplifying the nature of the confusion factors. Since the confusion factors we define encompass a broader range of abstract concepts, we introduce a Gaussian mixture model to represent the Gaussian distributions of the generalized and unknown abstract concepts.

In the Gaussian Self-Modeling (GSm) module, we associate the features corresponding to the potential $K$ types of confusion concepts with $K$ Gaussian distributions in Eq. (2).

$$\mathcal{N}(\mu_P(X;\theta), \sigma_P^2(X;\theta)) = \left\{ \mathcal{N}(\mu_p^i(X;\theta), \sigma_p^{i\,2}(X;\theta)) \right\}_{i=1}^{K} \quad (2)$$

where $\mathcal{N}(\mu_p^i(X;\theta), \sigma_p^{i\,2}(X;\theta))$ represents the prior Gaussian distribution of the $i$-th group of confusion factors, $\mu_p^i$ and $\sigma_p^i$ represent the mean and standard deviation of the Gaussian distribution.

Since the number of confusion factors varies across tasks, we avoided manually setting concept layers and instead used a large constant KK, allowing the network to learn relevant concepts implicitly. Unlike existing methods that average concept features, treating confounders as static, our approach accounts for their randomness and variability across samples.

To maintain the stochasticity and variability, we utilize the Gaussian distribution sampling strategy to represent statistical information while modeling the statistical characteristics of the latent space for each confusion factors. Moreover, we utilize the reparameterization technique[23] to make distribution sampling models trainable. Specifically, the sampled feature represented by the $i$-th Gaussian distribution can be expressed as Eq. (3).

$$z_i = \epsilon \cdot \sigma_i(X;\theta) + \mu_i(X;\theta). \quad (3)$$

3

where $\epsilon$ is sampled from the standard normal distribution $\epsilon \sim \mathcal{N}(0, 1)$ and $\mu_i, \sigma_i$ are the mean and standard deviation extracted from the GDEB.

To utilize posterior information to constrain the Gaussian modeling process, we construct a posterior constraint backbone by using the real mask $Y$ corresponding to the input image $X$ to learn the distribution space of confusion factors. We employ another backbone network to extract the posterior probability distribution regarding $K$ confusion factors, which can be expressed as Eq. (4).

$$\mathcal{N}(\mu_Q(Y; \theta'), \sigma_Q{}^2(Y; \theta')) = \left\{ \mathcal{N}(\mu_q^i(Y; \theta'), \sigma_q^{i\,2}(Y; \theta')) \right\}_{i=1}^K, \quad (4)$$

where $\mathcal{N}(\mu_q^i(X; \theta), \sigma_q^{i\,2}(X; \theta))$ represents the posterior Gaussian distribution of the $i$-th group of confusion factors, $\mu_q^i$ and $\sigma_q^i$ represent the mean and standard deviation.

We utilize $KL$ divergence to minimize the distance between the prior probability distribution of confusion factors and the posterior-constrained probability distribution. $KL$ divergence loss is designed for multiple independent Gaussian distributions, which can be expressed as Eq. (5).

$$\mathcal{L}_{kl} = KL(\mathcal{N}(\mu_P, \sigma_P^2) \parallel \mathcal{N}(\mu_Q, \sigma_Q^2))$$
$$= \frac{1}{K} \sum_{i=1}^K \left[ \log \frac{\sigma_q^i}{\sigma^i} + \frac{\sigma^{i\,2} + (\mu^i - \mu_q^i)^2}{2\sigma_q^{i\,2}} - \frac{1}{2} \right], \quad (5)$$

where $\mu^i, \sigma^i$ and $\mu_q^i, \sigma_q^i$ represent the $i$-th set mean and standard deviation in $K$-set distribution features.

To add strong constraints to the distribution from the segmentation process and to limit the modeling space of the confusion factors, we only focus on the segmentation boundary regions that are susceptible to bias interference. To achieve this, we propose $\mathcal{L}_{usd}$ that limits the uncertainty spatial distribution, which can be expressed as Eq. (6).

$$\mathcal{L}_{usd} = -\frac{1}{N} \sum_{i=1}^N (1 + V_i) \cdot [b_i \log \hat{b}_i + (1 - b_i) \log (1 - \hat{b}_i)], \quad (6)$$

where $b_i$ and $\hat{b}_i$ denote the ground truth and predicted values of the pixel $i$ within within the range of the boundary. To compute $b_i$, we extract the edge regions of the mask using the Sobel operator in our implementation. $V_i$ represents the degree of uncertainty of pixel $i$. We calculate the mean value $P$ of the pixels and compute the variance $V_i$ to quantify the uncertainty of the pixel $i$. This process can be described as Eq. (7).

$$V_i = (p_i - P)^2, \quad where \quad P = \frac{1}{N} \sum_{i=1}^N p_i, \quad (7)$$

where $N$ is the number of pixels, $p_i$ is the predicted pixel value before the thresholding process. Finally, we define the loss for implicit modeling of confusion factors as Eq. (8).

$$\mathcal{L}_{Gaus} = \mathcal{L}_{usd} + \mathcal{L}_{kl}, \quad (8)$$

### 3.4. Bias-Aware Causal Intervention

To implement an active intervention strategy for confusing features during segmentation, the CIBM is proposed, which integrates confusion semantics and decoding semantics to achieve bias intervention of confusion factors on prediction results. CIBM aims to implement the confusion factors intervention term $P(Y'|do(X))$ in Eq. (1). CIBM uses a normalized weighted geometric mean to project the probability prediction of confusion factors in the image onto the semantic feature space, which can be expressed as Eq. (9).

$$P(Y' \mid do(X)) = \mathbb{E}_{\hat{C} \sim P(\hat{C}|X)}[P(Y' \mid X, \hat{C})]$$
$$\approx P(y' \mid x, \sum_{\hat{c}_i}^{\hat{C}} \hat{c}_i P(\hat{c}_i)) \quad (9)$$
$$\approx P(y' \mid x \odot \sum_{\hat{c}_i}^{\hat{C}} (Gauss(x, \hat{c}_i; \theta'') \cdot P(\hat{c}_i))),$$

where $\odot$ denotes the fusion operation of confusion factors modeling and original image mapping. $Gaus(x, c_i; \theta'')$ represents the implicit modeling of the confusion factors $c_i$, which appearing in the image $x$.

CIBM assumes that the occurrence of various confusion factors follows a learnable random variable distribution. This modeling assumption is based on class differences observed in actual scenarios, such as in thyroid nodule imaging, where there is no equal relationship between edge blur and the frequency of specular reflection artifacts. CIBM introduces a learnable probability parameter $\Omega$, which is used to construct Gaussian mixture distributions and control the contribution strength of each confusion factor to semantic features as a mixing coefficient. The weighted modeling of the Gaussian mixture distribution is shown in Eq. (10).

$$\sum_{\hat{c}_i}^{\hat{C}} (Gauss(x, \hat{c}_i; \theta'') \cdot P(\hat{c}_i)) = \Omega \times Z, \quad s.t. \sum_{k=1}^K \Omega_{n,k} = 1 \quad (10)$$

where $\Omega \in \mathbb{R}^{n \times K}$ represents the weighted probability value of $K$ dimensional confusion features, $Z \in \mathbb{R}^K$ is the set of confusion features sampled in Eq. (3), and $n$ represents the number of feature channels corresponding to the decoding stage. The matrix multiplication operation $\Omega \times Z$ achieves feature fusion of confusion semantics, allowing the features of each channel in the corresponding decoding stage to adaptively adjust the interference level of different confusion factors.

The semantic feature $\Omega \times Z$ in the CIBM module is replicated and expanded into a two-dimensional feature map through the channel dimension, and is matched with the corresponding decoding feature $F_{decoder_i}$ to perform cascading splicing. Splicing features utilize convolution and Gaussian Error Linear Unit (GELU) based on Gaussian distribution to enhance expressive power, in the form of channel weights $\mathbb{S}$, compared with the original decoding feature $F_{decoder_i}$ performs channel wise weighting to form an enhanced decoding output, and the fusion calculation method is as Eq. (11).

$$F = \mathbb{S} \otimes (F_{decoder_i} \oplus \text{Repeat}(\Omega \times Z)) \quad (11)$$

where $\oplus$ represents element wise addition, $\otimes$ is channel level multiplication operation, $\mathbb{S}$ represents the channel weight of the fused feature after convolution and activation operation, and Repeat($\cdot$) represents expanding the one-dimensional vector in the spatial dimension into a two-dimensional feature matrix that matches $F_{decoder_i}$.

4

The MAMBO-NET model guides end-to-end network learning through multiple loss functions. For image segmentation tasks, a combination of Binary Cross Entropy Loss (BCELoss) and Dice loss is used to simultaneously optimize pixel level classification accuracy and region overlap; For the modeling task of confusion factors bias intervention, the Eq.(5) and (6) are used to model the confusion factors loss, guiding the model to learn the distribution characteristics of the confusion factors and optimize the segmentation ability of the confusion factors region. BCELoss improves pixel-level classification accuracy by calculating the difference between the predicted mask and the true mask pixel by pixel. The calculation of BCELoss is shown in Eq. (12).

$$L_{BCE} = -\frac{1}{N}\sum_{i=1}^{N}[y_i log\hat{y}_i + (1 - y_i)log(1 - \hat{y}_i)] \qquad (12)$$

where $y_i$ represents the true mask value of the $i$ th pixel, $\hat{y}_i$ represents the predicted value of the $i$th pixel, and $N$ is the number of pixels.

Dice loss measures segmentation performance from the perspective of region overlap and spatial consistency between predicted and annotated regions, which is shown as Eq. (13), and the meaning of the variables is the same as Eq. (12).

$$L_{Dice} = 1 - \frac{2\sum_{i=1}^{N}y_i\hat{y}_i}{\sum_{i=1}^{N}y_i + \sum_{i=1}^{N}\hat{y}_i} \qquad (13)$$

The total loss of the MAMBO-NET model consists of segmentation loss and confusion factors modeling loss, and the total loss expression is shown in Eq. (14), where $L_{Gaus}$ is Eq. (8).

$$L_{Total} = L_{Gaus} + L_{BCE} + L_{Dice} \qquad (14)$$

## 4. Experiments

### 4.1. Datasets and Settings

**Datasets.** We evaluated on five medical datasets: BUSI [24] (665 breast ultrasound images), DDTI [25] (3,644 thyroid ultrasound images), TUI [26] (15,233 thyroid ultrasound images), ISIC2018 [27] (3,694 dermoscopy images), and KVASIR [28] (1,000 colonoscopy images).

**Settings.** Using PyTorch on RTX 3090, we trained with SGD (momentum=0.9, weight decay=0.01), lr=1e-3, batch=32 for 200 epochs with cosine decay. Data was split 70%/30% (train/test) with standard augmentation (flipping/rotation/cropping).

### 4.2. Comparison with state-of-the-art methods

Tab. 1 shows MAMBO-NET's superior performance on ultrasound datasets, outperforming UNeXt by 3.66% in Dice and 2.33% in IoU. On BUSI, it improved Dice by 0.8%, IoU by 1.54%, and AUC by 0.35%, while on DDTI, gains were 2.28% in Dice, 4.27% in IoU, and 0.03% in AUC, confirming enhanced segmentation consistency.

Tab. 2 highlights MAMBO-NET's advantages on dermatoscopy and colonoscopy datasets. Compared to UNeXt, it

Table 1: The experimental results were obtained from three ultrasound datasets. The upward arrow ↑ denotes that a higher value is preferable, whereas the downward arrow ↓ indicates that a lower value is preferable. The best performance is highlighted in bold, and a horizontal line is used to mark the second-best performance.

| Methods | BUSI | | | | DDTI | | | | TUI | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ |
| UNet[1](2015) | 66.67 | 50.01 | 32.41 | 91.41 | 79.99 | 66.92 | 25.47 | 96.49 | 87.00 | 77.20 | 12.03 | 97.83 |
| UNet++[3](2018) | 74.44 | 59.35 | 20.00 | 91.07 | 84.29 | 72.99 | 13.36 | 96.89 | 89.01 | 80.34 | 8.56 | 97.73 |
| ResUNet[2](2020) | 70.69 | 55.17 | 44.86 | 92.26 | 82.94 | 71.21 | 53.61 | 95.11 | 86.31 | 76.31 | 17.15 | 95.18 |
| MedT[29](2021) | 67.90 | 51.55 | 31.83 | 94.21 | 73.66 | 58.73 | 25.93 | 97.61 | 77.52 | 63.82 | 20.94 | 98.38 |
| GGNet[30](2021) | 70.20 | 54.10 | 29.92 | 93.98 | 71.54 | 55.94 | 28.45 | 96.88 | 80.39 | 67.44 | 17.66 | 97.71 |
| TransUNet[31](2021) | 74.48 | 59.66 | 21.45 | 93.76 | 86.25 | 76.12 | 12.48 | 98.33 | 89.09 | 80.38 | 10.91 | 98.90 |
| AAU-net[32](2022) | 77.14 | 62.86 | 19.48 | 91.47 | 87.04 | 77.31 | 12.74 | 98.41 | 87.54 | 78.02 | 12.14 | 98.83 |
| Swin-Unet[33](2022) | 65.03 | 48.50 | 33.28 | 93.80 | 66.02 | 49.82 | 35.54 | 96.04 | 56.36 | 39.78 | 44.79 | 95.44 |
| UNeXt[4](2022) | 74.28 | 59.19 | 22.98 | 93.31 | 87.22 | 77.46 | 24.92 | 95.24 | 91.04 | 83.58 | 8.50 | 99.46 |
| Acc-Unet[34](2023) | 74.63 | 59.89 | 16.53 | 93.79 | 85.41 | 74.83 | 12.21 | 98.69 | 89.94 | 82.55 | 10.68 | 99.56 |
| BUSSeg[35](2023) | 75.24 | 60.51 | 18.49 | 93.61 | 87.38 | 77.63 | 13.46 | 97.62 | 89.23 | 80.78 | 9.27 | 99.76 |
| EffiSegNet[36](2024) | 73.45 | 58.22 | 20.80 | 94.35 | 84.56 | 73.42 | 13.87 | 98.34 | 88.39 | 79.43 | 11.44 | 99.29 |
| LMNet[37](2024) | 74.31 | 59.28 | 26.02 | 94.01 | 83.05 | 71.27 | 16.38 | 94.29 | 86.38 | 76.71 | 21.72 | 98.37 |
| MAMBO-NET(OURS) | 77.94 | 64.40 | 18.36 | 94.70 | 89.66 | 81.90 | 12.50 | 98.72 | 91.10 | 83.96 | 7.95 | 99.63 |

Table 2: The experimental results on ISIC2018 and KVASIR datasets. The best performance is highlighted in bold, and a horizontal line is used to mark the second-best performance.

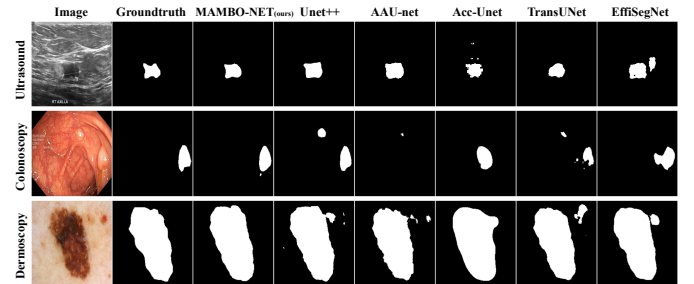| | ISIC2018 | | | | KVASIR | | | |
|---|---|---|---|---|---|---|---|---|
| | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ | Dice ↑ | IoU ↑ | FDR ↓ | AUC↑ |
| UNet[1](2015) | 86.21 | 76.12 | 15.22 | 95.83 | 75.63 | 61.27 | 30.25 | 95.02 |
| UNet++ [3](2018) | 88.65 | 79.52 | 21.01 | 95.37 | 81.87 | 69.35 | 18.99 | 95.21 |
| ResUNet [2](2020) | 88.14 | 79.03 | 19.83 | 96.78 | 73.98 | 58.97 | 22.71 | 92.76 |
| MedT [29](2021) | 87.69 | 78.37 | 17.22 | 96.31 | 62.18 | 45.44 | 37.44 | 91.79 |
| GGNet [30](2021) | 87.91 | 78.54 | 14.24 | 95.27 | 68.83 | 52.55 | 33.54 | 95.81 |
| TransUNet [31](2021) | 88.71 | 79.84 | 25.64 | 94.28 | 83.22 | 71.36 | 16.75 | 93.92 |
| AAU-Unet [32](2022) | 89.16 | 80.47 | 12.98 | 96.36 | 82.87 | 71.06 | 19.23 | 96.39 |
| Swin-Unet [33](2022) | 87.53 | 78.14 | 16.47 | 94.43 | 61.98 | 45.13 | 40.22 | 92.83 |
| UNeXt [4](2022) | 88.66 | 79.65 | 19.73 | 97.71 | 81.71 | 68.22 | 12.11 | 96.33 |
| Acc-unet [34](2023) | 85.67 | 75.31 | 14.39 | 92.74 | 79.49 | 65.72 | 23.21 | 95.94 |
| BUSSeg [35](2023) | 86.53 | 76.84 | 15.37 | 95.27 | 80.19 | 67.01 | 18.24 | 96.77 |
| EffiSegNet [36](2024) | 89.26 | 80.85 | 12.88 | 95.94 | 84.27 | 72.98 | 13.49 | 93.28 |
| LMNet[37](2024) | 87.50 | 77.87 | 17.47 | 95.32 | 80.38 | 67.24 | 24.28 | 91.54 |
| MAMBO-NET(OURS) | 89.30 | 80.97 | 11.39 | 97.23 | 83.96 | 72.11 | 10.24 | 97.88 |



Figure 4: Visualization results of segmentation of multiple models on ultrasound, dermatoscopy, and colonoscopy images.

achieved higher Dice (0.64%) and IoU (1.32%) with an 8.34% reduction in FDR. On colonoscopy data, improvements were 2.25% in Dice, 3.89% in IoU, and 1.89% in FDR. Although KVASIR's Dice and IoU were slightly lower than EffiSegNet, it reduced FDR by 3.25% and improved AUC by 1.11%, demonstrating better false-positive suppression. Joint analysis reveals MAMBO-NET's particularly significant improvements on low-resolution ultrasound data. Fig. 4 visually confirms segmentation quality enhancement across modalities.

### 4.3. Analysis of K sets Gaussian Distributions

As shown in Tab. 3, the results show a non-linear relation between the Gaussian distribution quantity K and model performance. Increasing K from 16 to 128 boosts segmentation metrics, but expanding to 512 causes saturation: Dice declines by

Table 3: Ablation experiments on *K* conducted on the BUSI dataset. The best performance is highlighted in bold. The model uses *K* = 128.

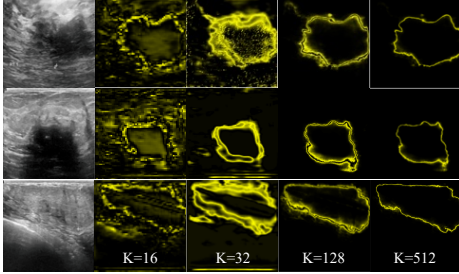| K | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ |
|---|--------|-------|-------|-------|
| 16 | 76.72 | 62.51 | 24.15 | 94.04 |
| 32 | 77.43 | 63.74 | 20.46 | 94.22 |
| 128 | 77.94 | 64.49 | **18.36** | **94.70** |
| 512 | **77.98** | **64.60** | 20.74 | 94.67 |



Figure 5: Feature entropy map generated by the decoder layer, where *K* represents the number of distributions in the GSm.

Table 4: The ablation experiment results of GSm and CIBM(%)

| | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ |
|---|--------|-------|-------|-------|
| Backbone | 74.28 | 59.19 | 22.98 | 93.31 |
| Backbone + GSm | 76.91 | 62.79 | 18.76 | 92.26 |
| Backbone + CIBM | 76.20 | 61.62 | 19.47 | 92.69 |
| Backbone + GSm + CIBM | **77.94** | **64.49** | **18.36** | **94.70** |

Table 5: The results of the model selection for GDEB and PCB (%)

| GDEB | PCB | Dice ↑ | IoU ↑ | FDR ↓ | AUC ↑ |
|------|-----|--------|-------|-------|-------|
| VGG16BN(w/o pre) | VGG16BN(w/o pre) | 76.63 | 62.21 | 20.11 | 93.87 |
| VGG16BN(pre) | VGG16BN(w/o pre) | 76.36 | 61.77 | **17.29** | 92.07 |
| VGG16BN(pre) | VGG16BN(pre) | 77.94 | 64.49 | 18.36 | **94.70** |
| DenseNet(pre) | DenseNet(pre) | 76.22 | 61.64 | 17.38 | 93.33 |
| ResNet(pre) | ResNet(pre) | 76.01 | 61.44 | 18.56 | 91.19 |

0.04% and IoU by 0.11% before leveling off, while FDR rises 2.38%, indicating diminishing returns from excessive components and feature interference. Decoding feature entropy visualizations confirms K's role in modulating feature bias.

Fig. 5 shows that lighter areas mean higher pixel uncertainty. At *K* = 16, large light regions in the entropy map suggest poor encoding of confounders, leading to high decoding uncertainty, especially in artifact areas. As K rises to 32 and 128, light areas shrink, showing better confounder representation and reduced uncertainty. At *K* = 512, uncertain regions over-converge, implying noise in bias modeling and loss of confounder simulation ability.

### 4.4. Ablation of GSm and CIBM

Two control schemes were set up in the experiment: (1) feature concatenation and linear mapping were used to replace the Causal Intervention for Bias Mitigation(CIBM); (2) replace the sampling features of GSm with decoder features. The module ablation experiment will be conducted on the BUSI dataset, and the experimental results are shown in Tab. 4. Results demonstrate that GSm alone improved Dice by 2.63 while reducing FDR by 4.22%, though it caused a slight 1.05% AUC drop due to increased feature complexity. CIBM alone enhanced Dice by 2.08% and IoU by 2.43% while lowering FDR by 3.51%, but its effectiveness was constrained by limited confounder modeling. The combined use of both modules achieved optimal performance across all metrics, showing their complementary nature - GSm provides confounder distributions that enable more effective bias intervention through CIBM.

### 4.5. Model selection for GDEB and PCB in GSm

The VGG16BN[38], DenseNet[39], and ResNet[40] were used as alternative solutions for Gaussian Distribution Extraction Backbone(GDEB) and Posterior Constraint Backbone(PCB). The pre-training involved in the experiment was conducted on ImageNet21k, and global average pooling and linear mapping of the same dimension were used to ensure size alignment across features. The impact of the selection of GDEB and PCB on the model results on the BUSI dataset is shown in Tab.5.

Both GDEB and PCB used the untrained VGG16BN to achieve a Dice of 76.63% and an IoU of 62.21%. However, when only PCB used pre-trained VGG16BN, Dice and IoU decreased by 0.27% and 0.44%, respectively. However, when both GDEB and PCB used pre-trained VGG16BN, Dice, IoU, and AUC reached their optimal levels, with improvements of 1.31%, 2.28%, and 0.83% compared to suboptimal levels. Experiments on VGG16BN have shown that ensuring that the prior and posterior constraint distributions are in the same feature space is beneficial for the convergence of KL divergence loss and alignment of distribution. A symmetric architecture design is advantageous for optimizing end-to-end segmentation. The performance of the pre-trained backbone of DenseNet and ResNet lags behind that of VGG16BN, indicating that the architecture of VGG16BN is more suitable for Gaussian feature extraction and posterior constraint modeling.

## 5. Conclusion

We propose MAMBO-NET, a novel framework for medical image segmentation that systematically addresses the challenges posed by confusion factors using causal inference theory. By integrating causal relationship modeling, Gaussian Self-modeling, and causal intervention for bias mitigation, MAMBO-NET effectively disentangles confusion factors from target features and mitigates their impact on segmentation decisions. Experimental results demonstrate its superiority over state-of-the-art methods, with Gaussian Self-modeling proving crucial for comprehensive latent space representation. The further analysis highlights the importance of posterior constraints and pretrained backbones in optimizing performance. These contributions establish MAMBO-NET as an effective approach for enhancing segmentation accuracy and robustness in complex medical imaging scenarios.

# References

[1] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18, Springer, 2015, pp. 234–241.

[2] Z. Zhang, Q. Liu, Y. Wang, Road extraction by deep residual u-net, IEEE Geoscience and Remote Sensing Letters 15 (5) (2018) 749–753.

[3] Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, J. Liang, Unet++: A nested u-net architecture for medical image segmentation, in: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4, Springer, 2018, pp. 3–11.

[4] J. M. J. Valanarasu, V. M. Patel, Unext: Mlp-based rapid medical image segmentation network, in: International conference on MICCAI, Springer, 2022, pp. 23–33.

[5] M. Hernan, J. Robins, Causal inference: What if chapman hall/crc, boca raton (2020).

[6] N. Zhao, A. Basarab, D. Kouamé, J.-Y. Tourneret, Joint segmentation and deconvolution of ultrasound images using a hierarchical bayesian model based on generalized gaussian priors, IEEE transactions on Image Processing 25 (8) (2016) 3736–3750.

[7] Z. Zhao, S. Du, Z. Xu, Z. Yin, X. Huang, X. Huang, C. Wong, Y. Liang, J. Shen, J. Wu, J. Qu, L. Zhang, Y. Cui, Y. Wang, L. Wee, A. Dekker, C. Han, Z. Liu, Z. Shi, C. Liang, Swinhr: Hemodynamic-powered hierarchical vision transformer for breast tumor segmentation, Computers in Biology and Medicine 169 (2024) 107939.

[8] B. Liu, D. Wang, X. Yang, Y. Zhou, R. Yao, Z. Shao, J. Zhao, Show, deconfound and tell: Image captioning with causal inference, in: Proceedings of the IEEE/CVF Conference on CVPR, 2022, pp. 18041–18050.

[9] X. Yang, H. Zhang, J. Cai, Deconfounded image captioning: A causal retrospect, IEEE Transactions on Pattern Analysis and Machine Intelligence 45 (11) (2021) 12996–13010.

[10] J. Wang, C. Zhong, C. Feng, Y. Zhang, J. Sun, Y. Yokota, Disentangled representation for cross-domain medical image segmentation, IEEE Transactions on Instrumentation and Measurement 72 (2022) 1–15.

[11] R. Gu, G. Wang, J. Lu, J. Zhang, W. Lei, Y. Chen, W. Liao, S. Zhang, K. Li, D. N. Metaxas, et al., Cddsa: Contrastive domain disentanglement and style augmentation for generalizable medical image segmentation, Medical Image Analysis 89 (2023) 102904.

[12] X. Li, Y. Pang, Y. Yuan, L1-norm-based 2dpca, IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics) 40 (4) (2010) 1170–1175.

[13] Y. Wang, B. Xiao, X. Bi, W. Li, X. Gao, Boundary-aware prototype in semi-supervised medical image segmentation, IEEE Transactions on Image Processing (2024).

[14] Y. Liu, L. Lin, K. K. Wong, X. Tang, Procns: Progressive prototype calibration and noise suppression for weakly-supervised medical image segmentation, IEEE Journal of Biomedical and Health Informatics (2024).

[15] L. Zhou, L. Wang, W. Li, B. Lei, J. Mi, W. Yang, Multi-stage liver segmentation in ct scans using gaussian pseudo variance level set, IEEE Access 9 (2021) 101414–101423.

[16] S. Gao, H. Zhou, Y. Gao, X. Zhuang, Bayeseg: Bayesian modeling for medical image segmentation with interpretable generalizability, Medical Image Analysis 89 (2023) 102889.

[17] Z. Chen, Z. Tian, J. Zhu, C. Li, S. Du, C-cam: Causal cam for weakly supervised semantic segmentation on medical image, in: Proceedings of the IEEE/CVF Conference on CVPR, 2022, pp. 11676–11685.

[18] X. Song, J. Liu, Y. Liu, Y. Li, W. Lei, R. Wang, Rethinking radiology report generation via causal inspired counterfactual augmentation, in: Proceedings of the 15th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, 2024, pp. 1–10.

[19] K. Chen, S. Sun, J. Zhao, Camil: Causal multiple instance learning for whole slide image classification, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 38, 2024, pp. 1120–1128.

[20] Y. Chen, M. Wei, Z. Zheng, J. Hu, Y. Shi, S. Xiong, X. X. Zhu, L. Mou, Causalclipseg: Unlocking clip's potential in referring medical image segmentation with causal intervention, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2024, pp. 77–87.

[21] C. Wang, J. Li, F. Zhang, X. Sun, H. Dong, Y. Yu, Bilateral asymmetry guided counterfactual generating network for mammogram classification, IEEE Transactions on Image Processing 30 (2021) 7980–7994.

[22] J. G. Richens, Improving the accuracy of medical diagnosis with causal machine learning, Nature communications 11 (1) (2020) 3923.

[23] D. P. Kingma, Auto-encoding variational bayes, arXiv preprint arXiv:1312.6114 (2013).

[24] W. Al-Dhabyani, M. Gomaa, H. Khaled, A. Fahmy, Dataset of breast ultrasound images, Data in brief 28 (2020) 104863.

[25] L. Pedraza, C. Vargas, F. Narváez, O. Durán, E. Muñoz, E. Romero, An open access thyroid ultrasound image database, in: 10th International symposium on medical information processing and analysis, Vol. 9287, SPIE, 2015, pp. 188–193.

[26] Z. Bai, L. Chang, R. Yu, X. Li, X. Wei, M. Yu, Z. Liu, J. Gao, J. Zhu, Thyroid nodules risk stratification through deep learning based on ultrasound images, Medical Physics 47 (12) (2020) 6355–6365.

[27] M. A. Al-Masni, D.-H. Kim, Multiple skin lesions diagnostics via integrated deep convolutional networks for segmentation and classification, Computer methods and programs in biomedicine 190 (2020) 105351.

[28] D. Jha, P. H. Smedsrud, M. A. Riegler, P. Halvorsen, T. De Lange, D. Johansen, H. D. Johansen, Kvasir-seg: A segmented polyp dataset, in: MultiMedia modeling: 26th international conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, proceedings, part II 26, Springer, 2020, pp. 451–462.

[29] J. M. J. Valanarasu, P. Oza, I. Hacihaliloglu, V. M. Patel, Medical transformer: Gated axial-attention for medical image segmentation, in: MICCAI 2021: 24th international conference, Springer, 2021, pp. 36–46.

[30] C. Xue, L. Zhu, H. Fu, X. Hu, X. Li, H. Zhang, P. A. Heng, Global guidance network for breast lesion segmentation in ultrasound images, Medical Image Analysis 70 (2021) 101989.

[31] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, Y. Zhou, Transunet: Transformers make strong encoders for medical image segmentation, arXiv preprint arXiv:2102.04306 (2021).

[32] G. Chen, L. Li, Y. Dai, J. Zhang, M. H. Yap, Aau-net: an adaptive attention u-net for breast lesions segmentation in ultrasound images, IEEE Transactions on Medical Imaging 42 (5) (2022) 1289–1300.

[33] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, Swin-unet: Unet-like pure transformer for medical image segmentation, in: European conference on computer vision, Springer, 2022, pp. 205–218.

[34] N. Ibtehaz, D. Kihara, Acc-unet: A completely convolutional unet model for the 2020s, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2023, pp. 692–702.

[35] H. Wu, X. Huang, X. Guo, Z. Wen, J. Qin, Cross-image dependency modeling for breast ultrasound segmentation, IEEE Transactions on Medical Imaging 42 (06) (2023) 1619–1631.

[36] I. Vezakis, K. Georgas, D. Fotiadis, G. Matsopoulos, Effisegnet: Gastrointestinal polyp segmentation through a pre-trained efficientnet-based network with a simplified decoder, in: Proceedings of the International Conference of the IEEE Engineering in Medicine and Biology Society, IEEE, 2024, pp. 1–4.

[37] Z. Lu, C. She, W. Wang, Q. Huang, Lm-net: A light-weight and multiscale network for medical image segmentation, Computers in Biology and Medicine 168 (2024) 1–12.

[38] K. Simonyan, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556 (2014).

[39] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on CVPR, 2017, pp. 4700–4708.

[40] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on CVPR, 2016, pp. 770–778.