# *OpenUS*: A Fully Open-Source Foundation Model for Ultrasound Image Analysis via Self-Adaptive Masked Contrastive Learning

Xiaoyu Zheng[1,*]    Xu Chen[1]    Awais Rauf[1]    Qifan Fu[1]

Benedetta Monosi[2]    Felice Rivellese[2]    Myles J. Lewis[2]

Shaogang Gong[3]    Gregory Slabaugh[1,*]

[1]Digital Environment Research Institute (DERI), Queen Mary University of London
[2]The William Harvey Research Institute, Queen Mary University of London
[3]School of Electronic Engineering and Computer Science, Queen Mary University of London
London, UK

[*]zhengxiaoyu0427@gmail.com/g.slabaugh@qmul.ac.uk

## Abstract

*Ultrasound (US) is one of the most widely used medical imaging modalities, thanks to its low cost, portability, real-time feedback, and absence of ionizing radiation. However, US image interpretation remains highly operator-dependent and varies significantly across anatomical regions, acquisition protocols, and device types. These variations, along with unique challenges such as speckle, low contrast, and limited standardized annotations, hinder the development of generalizable, label-efficient ultrasound AI models. In this paper, we propose OpenUS, the first reproducible, open-source ultrasound foundation model built on a large collection of public data. OpenUS employs a vision Mamba backbone, capturing both local and global long-range dependencies across the image. To extract rich features during pre-training, we introduce a novel self-adaptive masking framework that combines contrastive learning with masked image modeling. This strategy integrates the teacher's attention map with student reconstruction loss, adaptively refining clinically-relevant masking to enhance pre-training effectiveness. OpenUS also applies a dynamic learning schedule to progressively adjust the difficulty of the pre-training process. To develop the foundation model, we compile the largest to-date public ultrasound dataset comprising over 308K images from 42 publicly available datasets, covering diverse anatomical regions, institutions, imaging devices, and disease types. Our pre-trained OpenUS model can be easily adapted to specific downstream tasks by serving as a backbone for label-efficient fine-tuning. Code is available at* https://github.com/XZheng0427/ OpenUS.
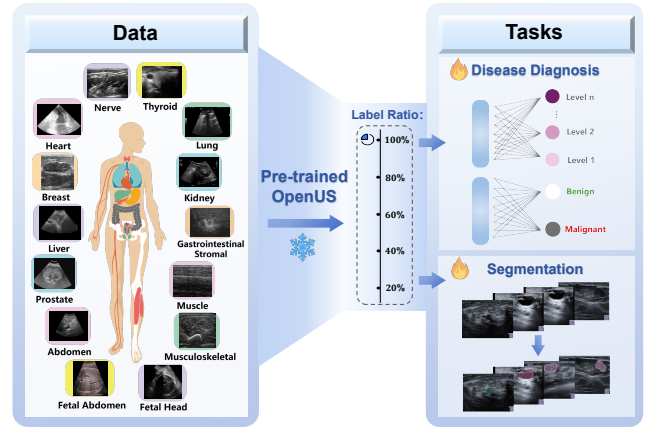
Figure 1. Overview of Universal US Foundation Model.

## 1. Introduction

Ultrasound (US) imaging is widely used in healthcare due to its affordability, portability, real-time feedback, and safety [15]. With applications ranging from fetal monitoring to cardiac, abdominal, and musculoskeletal imaging, US plays a central role in both diagnostics and image-guided interventions. Recent advances in computer-aided US analysis have enabled automated classification, segmentation, and disease detection [11, 38, 72], yet the clinical adoption of supervised deep learning methods remains limited by the scarcity of large, diverse, and well-annotated datasets. Moreover, US images vary considerably across organs, acquisition protocols, and devices, complicating the generalization of task-specific models. These challenges have created growing interest in US foundation models that can

learn transferable representations across anatomical regions and clinical tasks, facilitating label-efficient adaptation to new applications [37]. Such models have the potential to unlock scalable, intelligent US workflows and support broader integration into real-world healthcare systems.

Recently, visual foundation models [81] have attracted significant attention in the field of natural image analysis, demonstrating strong potential for developing generalized and robust representations. These models are typically developed using self-supervised pre-training methods on large-scale, unlabeled datasets [25, 75, 79]. To effectively learn meaningful representations in self-supervised learning (SSL), researchers have developed various visual pretext tasks [33, 44, 74], broadly categorized into contrastive and generative approaches. Motivated by successes in natural image analysis, extending foundation models to US imaging offers promising potential for efficiently developing versatile models applicable across various tasks and anatomical structures. However, transferring foundation models from natural images to US images remains challenging due to fundamental differences in imaging principles [84]. Therefore, there is a critical need to develop a foundation model specifically designed for US imaging, accompanied by comprehensive evaluation of its versatility and adaptability across diverse downstream tasks. However, current US foundation models [37, 39] rely on proprietary US images for pre-training, hindering reproducibility. Addressing this issue is essential for developing a fully reproducible, open-source US foundation model.

In this paper, we propose *OpenUS*, a foundation model developed for universal US imaging analysis with label efficiency and downstream task adaptability, as shown in Fig. 1. *OpenUS* is built solely on publicly available data, making it the first reproducible ultrasound foundation model. To achieve self-supervised pre-training on clinically relevant and generalizable US features, we use vision Mamba, based on VMamba [45, 86], to capture long-range spatial dependencies. Second, we propose a novel self-adaptive masking strategy that combines the teacher's attention map with the student's reconstruction loss. This approach enables the model to more accurately target and mask challenging, clinically relevant regions by using student feedback to correct for potential inaccuracies in the teacher's attention. Additionally, we incorporate multi-view masked modeling into the contrastive learning framework, which not only encourages the model to learn more discriminative representations but also compels it to capture more refined, pixel-level features. Extensive experiments confirm that *OpenUS* exhibits strong generalizability, superior performance, and enhances label efficiency for downstream task fine-tuning. Overall, our contributions are summarized as follows:

- We are the first to develop a fully reproducible, open-source US foundation model, *OpenUS*. [1] *OpenUS* is designed as a plug-and-play module to improve the label-efficiency and performance of various downstream US image analysis tasks. *OpenUS* is trained on large-scale, fully open-access datasets spanning diverse organs, institutions, and devices.
- We propose a multi-view masked contrastive representation learning framework that integrates global-local-view masked modeling with a contrastive approach for pre-training. This framework effectively overcomes the limitations of conventional contrastive methods in capturing dense pixel-level dependencies by enabling the prediction of pixel intensities within masked regions.
- We introduce a novel self-adaptive masking strategy to extract clinically relevant and generalizable US features. Teacher attention-guided masking and student reconstruction-loss masking synergistically enhance the robust extraction of effective US image features.
- We conduct extensive experiments on 38 fully public datasets for pre-training and 4 datasets for classification and segmentation downstream tasks to evaluate the performance of *OpenUS* compared to existing methods.

## 2. Related Work

### 2.1. Contrastive Learning

SSL adopts unlabeled data to extract meaningful representations without the need for explicit supervision, leading to considerable progress in pattern recognition [5, 17]. In this paradigm, contrastive learning (CL) initially attracted significant interest in both computer vision [10, 12, 30] and medical imaging domains [4, 82]. CL captures rich semantic representations by bringing similar image pairs closer together and pushing dissimilar pairs further apart in the feature domain. However, CL depends on specifically designed data augmentation techniques to generate paired images for learning [13]. In applications that demand accurate pixel-level embeddings for downstream tasks, i.e., classification, segmentation and detection, CL approaches may overlook the subtle details essential for identifying small organs or lesion regions as they tend to emphasize the broad, discriminative features rather than focusing on the fine-grained features.

### 2.2. Masked Image Modeling

Inspired by the breakthroughs achieved with masked language modeling in natural language processing, masked image modeling (MIM) has garnered much attention in computer vision, i.e., the bidirectional encoder representation from image transformers (BEiT) [7], masked autoencoder

---

[1] Concurrent to this work, Zhang et al. introduced *EchoCare* [83]. At the time of writing, this EchoCare appears on ArXiv, and its code and pretrained weights have not been released.
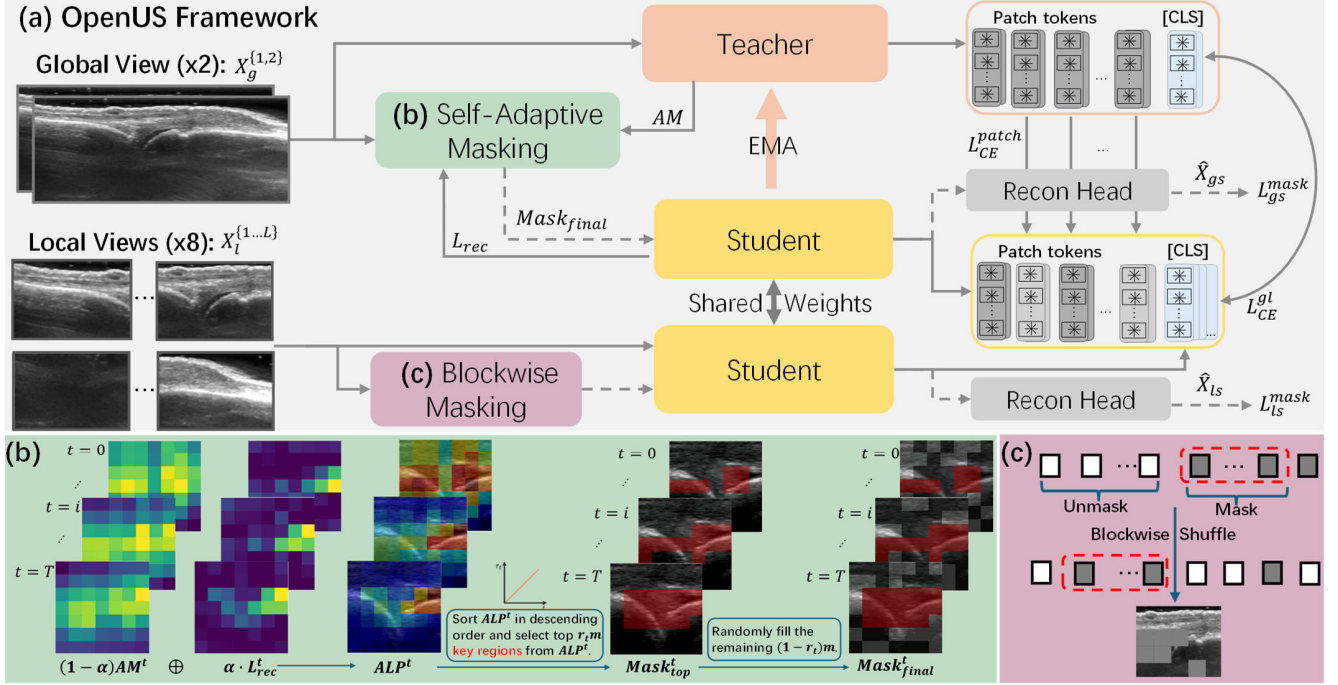
Figure 2. (a) The *OpenUS* pipeline including the Teacher and Student models, masking approaches and reconstructions heads. For global and local views, we design two distinct masking strategies: (b) self-adaptive masking and (c) random block-wise masking. Both are integrated with masked image reconstruction and contrastive learning.

(MAE) [31] and the simple framework for MIM (SimMIM) [77]. Both He et al. [31] and Xie et al. [77] propose to mask random image patches and then predict the corresponding RGB pixel values. Following the introduction of BEiT and MAE, iBOT [87] emerged as an online tokenizer adaptation of BEiT designed to overcome its shortcoming of capturing only low-level semantics within local details. While MIM strategies have been successful in natural images, they can overlook the detailed and clinically meaningful patterns in medical images, such as anatomical boundaries or diagnostically relevant regions, which can limit the effectiveness of self-supervised pre-training. Recent advancements in MIM have shifted focus towards novel reconstruction targets [65, 73] and masking strategies [43, 46, 71]. Specifically, Liu et al. [46] proposed an attention-driven masking approach to overcome the limitations of random masking in fully utilizing semantic information. However, novel MIM strategies specifically designed for medical image analysis remain underexplored.

### 2.3. Foundation Model for US Image Analysis

In recent years, foundation models based on SSL have attracted significant attention in medical image analysis, including applications in ultrasound (US) imaging. Specifically, US imaging introduces unique challenges compared to other modalities, including speckle patterns, motion

artefacts, and geometric discrepancies across probe types and manufacturers [69]. USFM [37] employs a spatial-frequency dual-MIM method to extract effective features from low-quality US images. Kang et al. [39] proposed deblurring MIM that integrated the US-specific deblurring task into the standard MAE framework. Instead of only reconstructing masked patches, the model also learns to restore the characteristic low signal-to-noise US texture and capture finer detail. Basu et al. [8] introduced FocusMAE, a region-focused masked modeling approach for US videos aimed at detecting gallbladder carcinoma. However, this method exhibits certain limitations when applied to general purpose US imaging. While MIM has become the dominant paradigm for self-supervised ultrasound pre-training, relying on MIM alone can limit the ability to learn rich, discriminative representations needed for downstream tasks.

## 3. Method

Fig. 2 illustrates the overall pipeline of *OpenUS*. Here, we first discuss the preliminaries of vision Mamba and self-distillation MIM. Then, we introduce our *OpenUS* pipeline and provide a detailed explanation of our proposed self-adaptive masking strategy.
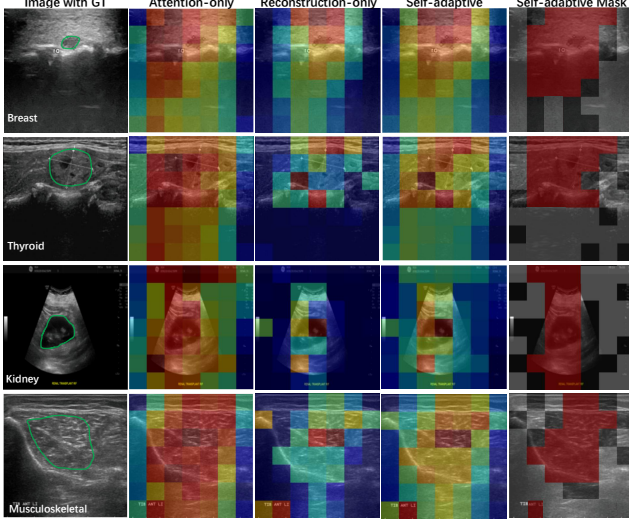
3

Figure 3. Visual comparison with segmentation ground truths, attention-only, reconstruction loss and self-adaptive $ALP$ scores. In the last column, the red areas tend to overlap with clinically relevant regions, while the dark grey regions represent the remaining randomly masked areas.

## 3.1. Preliminaries

### 3.1.1. Vision Mamba (VMamba).

Here we describe the common backbone VMamba [45], which has a hierarchical encoder design. Specifically, an input image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, is first partitioned into patches by a stem module, yielding a 2D feature map with the spatial dimension of $\frac{H}{4} \times \frac{W}{4}$, where $(H, W)$ denotes the height and width of the original image, $C$ is the number of the channels. Hierarchical representations are then created through multiple network stages, achieving resolutions of $\frac{H}{8} \times \frac{W}{8}$, $\frac{H}{16} \times \frac{W}{16}$, and $\frac{H}{32} \times \frac{W}{32}$. Each stage includes a down-sampling layer followed by the Visual State Space (VSS) block that is derived from the zero-order hold (ZOH) method [24],

$$\mathbf{h}_b = e^{\mathbf{A}(\Delta_a + \cdots + \Delta_{b-1})} \left( \mathbf{h}_a + \sum_{i=a}^{b-1} \mathbf{B}_i u_i e^{-\mathbf{A}(\Delta_a + \cdots + \Delta_i)} \Delta_i \right) \tag{1}$$

where $[\mathbf{a}, \mathbf{b}]$ is the corresponding discrete step interval. $\mathbf{A} \in \mathbb{R}^{N \times N}$, $\mathbf{B} \in \mathbb{R}^{N \times N}$, $\mathbf{h}$ denotes the hidden state, and $\Delta$ is the time-scale parameter.

### 3.1.2. Self-distillation masked image modeling.

Self-distillation, which employs a moving average of the student model as a teacher, has been explored for SSL in BYOL [23]. This approach was subsequently extended to vision transformers in DINO [10], where the distillation loss is applied globally to the [cls] token. iBOT [87] addresses this task through MIM, where the loss is computed densely over the masked tokens. Following iBOT, the transformer encoder is succeeded by a head, comprising a Multi-Layer Perceptron (MLP) and scaled softmax, enabling the

output token embeddings to be interpreted as probabilities. The teacher parameters $\theta'$ are updated from the student parameters $\theta$ using an exponential moving average (EMA), following the update rule $\theta' \leftarrow \lambda \theta' + (1 - \lambda)\theta$, where $\lambda \in [0, 1)$ controls the momentum of the update. For each input image, two augmented global views of standard resolution are generated, resulting in tokenized images $\mathbf{P}^a$, $\mathbf{P}^b$ and corresponding mask vectors $\mathbf{m}^a$, $\mathbf{m}^b$. For each view $v \in \{a, b\}$ and its respective masked token, the MIM objective is to minimize the reconstruction loss between the student's output ($f_\theta$ applied to the masked input $\tilde{\mathbf{P}}^v$) and the teacher's output ($f_{\theta'}$ applied to the non-masked input $\mathbf{P}^v$),

$$\mathcal{L}_{\text{MIM}} = -\sum_{\mathbf{v} \in \mathbf{V}} \sum_{i=1}^{n} m_i^{\mathbf{v}} f_{\theta'}(\mathbf{P}_i^{\mathbf{v}}) \log \left( f_\theta(\mathbf{P}_i^{\mathbf{v}}) \right) \tag{2}$$

## 3.2. OpenUS via Global-local Masked Contrastive Representation Learning

The pipeline of our proposed *OpenUS* is illustrated in Fig. 2(a). It employs contrastive learning combined with global-local mask modeling to simultaneously learn representations that are both fine-grained and discriminative. The contrastive learning component of our *OpenUS* model follows iBOT [87], which also adopts a self-distillation approach to facilitate representation learning.

Given a US image, two types of views are created under random data augmentation, i.e., global views $\mathbf{X}_g^i$ and local views $\mathbf{X}_l^j$. During pre-training, the global views are fed into both the teacher and student networks, while the local views are only fed into the student network. The network output is normalized using a softmax function with temperature $\tau$ to yield the probability distribution $p = \text{softmax}(f/\tau)$. Furthermore, the cross-entropy loss function is employed to compute the loss $L_{CE}^{GL}$ between the teacher's global-view [cls] tokens and the student's local-view [cls] tokens, as well as the loss $L_{CE}^{patch}$ between the student network outputs for the masked-view patch tokens and the teacher network outputs for the non-masked-view patch tokens. Building upon iBOT, we incorporate a global-local view mask reconstruction pre-training task into the contrastive learning framework to enhance dense pixel dependencies, which are critical for dense prediction tasks such as segmentation. Specifically, the reconstruction head leverages the contextual information from unmasked patches to regress the dense pixel intensities of the masked patches. The loss for global view reconstruction is $\mathcal{L}_{gs}^{mask} = \frac{1}{N_{gs}^{mask}} \sum_{i=1}^{G} \left\| \hat{\mathbf{X}}_g^i - \mathbf{X}_g^i \right\|$.

Due to the low contrast between lesions and normal tissues in US images, global views may have difficulty accurately capturing local details. Hence, we employ the random blockwise masking strategy [7] on local views to capture more fine-grained and detailed information, as shown in Fig. 2(c). The loss for local view reconstruction is $\mathcal{L}_{ls}^{mask} = \frac{1}{N_{ls}^{mask}} \sum_{i=1}^{L} \left\| \hat{\mathbf{X}}_l^i - \mathbf{X}_l^i \right\|$.

*OpenUS* leverages both a contrastive loss and a reconstruction loss during optimization, with the total loss defined as: $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}}^{gl} + \mathcal{L}_{\text{CE}}^{patch} + \mathcal{L}_{gs}^{mask} + \mathcal{L}_{ls}^{mask}$. The student network updates its parameters $\theta_s$ by minimizing the total loss $\mathcal{L}_{\text{total}}$, while the teacher network parameters $\theta_t$ are updated using an EMA of the student weights, following the rule: $\theta_t \leftarrow \lambda\theta_t + (1 - \lambda)\theta_s$, where $\lambda$ defines the momentum coefficient. By integrating contrastive learning with MIM, the model learns both holistic discriminative features and fine-grained pixel details, enhancing its capability to process complex, clinically relevant features from US images.

### 3.3. Self-adaptive Masking Strategy

MIM employs random and attention-only mask strategies to extract meaningful representations from pre-training images. While these strategies are effective for general image datasets with stable and well-curated distributions, they do not adequately address the unique characteristics inherent in medical images. US images often exhibit strong speckle, low contrast between tissue boundaries, and appearance variability caused by differences in probe type, acquisition protocol, and operator technique. Furthermore, diagnostically important structures may occupy only small or subtle regions of the image, making them harder to capture with purely random or attention-only masking. Therefore, to address the challenges outlined, we propose a self-adaptive masking strategy utilizing a teacher-student feedback mechanism, as depicted in Fig. 2(b). The specifics are detailed below.

#### 3.3.1. Semantic Information Extractor.

Our network employs a self-attention mechanism known as 2D Selective Scan (SS2D) [45], which boosts the learning capacity of the network. Specifically, we take the US image, from which we sample the global views $\{\mathbf{X}_g^i \in \mathbb{R}^{3 \times H_g \times W_g}\}_{i=1}^G$. Each image is then divided into $N = H_gW_g/P^2$ patches, which are mapped into patch tokens and fed into the multiple teacher VSS blocks with the down-sampling layer to create hierarchical representations. The network incorporates the learnable class token and patch tokens, which represent the global and local features extracted by the network across the spatial dimension. Unlike Vision Transformers (ViTs) that incorporate a learnable class token (cls), our network utilizes average pooling ($avgpool$) to derive the class token from the patch tokens. Given the intermediate token $\mathbf{x}^n \in \mathbb{R}^{(N+1) \times D}$ from block $n$, the token for the subsequent block is computed as follows:

$$\begin{cases} \mathbf{x}^n = SS2D\left(LN(\mathbf{x}^{(n-1)})\right) + \mathbf{x}^{n-1}, \\ \mathbf{x}^{n+1} = MLP\left(LN(\mathbf{x}^n)\right) + \mathbf{x}^n \end{cases} \tag{3}$$

where $SS2D$, $LN$ and $MLP$ represent the 2D selective scan attention, layer normalization and multi-layer perception, respectively. We can obtain the attention map of the

last layer of the VSS blocks by computing the correlation between the query embeddings $Q$ and the key embeddings $K$, a mechanism similar to self-attention of ViT. Specifically, based on Eq. 1, the discretized solution of time-varying SSMs can be defined as:

$$\mathbf{h}_b = \mathbf{w}_T \odot \mathbf{h}_a + \sum_{i=1}^T \frac{\mathbf{w}_T}{\mathbf{w}_i} \odot \left(\mathbf{K}_i^\top \mathbf{V}_i\right) \tag{4}$$

where $\odot$ represents the element-wise product between matrices. Building upon the hidden state $\mathbf{h_b}$ (Eq. 4), the output of SSM, i.e., $\mathbf{Y}^{(i)}$, can be defined as:

$$\mathbf{Y}^{(i)} = \left[\mathbf{Q} \odot \mathbf{w}^{(i)}\right]\mathbf{h}_a^{(i)} + \left[\left(\mathbf{Q} \odot \mathbf{w}^{(i)}\right)\left(\frac{\mathbf{K}}{\mathbf{w}^{(i)}}\right)^\top \odot \mathbf{M}\right]\mathbf{V}^{(i)} \tag{5}$$

where $\mathbf{h_a} \in \mathbb{R}^{(D_k)}$ is the hidden state at step $a$, $\mathbf{M}$ defines the temporal mask matrix. Therefore, the attention map ($AM$) is calculated by $\mathbf{QK}^\top$ and $\left(\mathbf{Q} \odot \mathbf{w}\right)\left(\frac{\mathbf{K}}{\mathbf{w}}\right)^\top$, which is capable of obtaining an approximation of the clinically relevant lesion region of the US image.

#### 3.3.2. Reconstruction Loss Predictor.

Our network incorporates a reconstruction head employing a lightweight decoder composed of a few linear layers to reconstruct the masked regions. Concretely, let $\mathbf{Mask}_{final}^t$ denote the newly generated mask, the reconstruction loss can be formulated as:

$$\mathcal{L}_{\text{rec}}^t = \mathcal{M}\left(d_s(f_s(\mathbf{x} \odot \mathbf{Mask}_{\text{final}}^t)), \mathbf{x} \odot (1 - \mathbf{Mask}_{\text{final}}^t)\right) \tag{6}$$

where $\odot$ denotes element-wise dot product, and $\mathbf{x} \odot (1 - \mathbf{Mask}_{final}^t)$ represents masked (invisible) patches and vice versa. $\mathcal{M}(\cdot, \cdot)$ represents the similarity measurement, i.e., $l_2$-distance. We utilize the EMA of the student's reconstruction loss ($\mathcal{L}_{\text{rec}}^t$), which provides a more stable loss map highlighting regions that consistently pose challenges for the student. Furthermore, we perform an *argsort* operation on $\mathcal{L}_{\text{rec}}^t$ to rank the masked regions based on their reconstruction losses in descending order. This ranking directs the student network's attention toward reconstructing the most difficult regions, thereby effectively predicting clinically relevant features.

#### 3.3.3. Mask Generator: Teacher-Student Feedback.

An attention-only approach can become repetitive, repeatedly masking regions the student already finds easy, while a reconstruction-only strategy risks focusing on noisy or irrelevant patches that are difficult but not informative. Our method avoids these pitfalls by integrating the teacher's "top-down" attention to critical information with the student's "bottom-up" feedback regarding challenging aspects. This synergy ensures the model consistently targets the most valuable and important regions for feature learning, accelerating convergence and leading to more robust representations.

5

Specifically, the first is the teacher attention map ($\mathbf{AM}_t$), where higher values indicate regions considered important by the teacher network. The second is the student reconstruction loss map ($\mathbf{L}_{rec_s}$), where higher values indicate patches the student finds challenging to reconstruct. These two inputs are integrated into a unified metric, which we term the *Adaptive Learning Priority* (**ALP**) score. For each image patch, the **ALP** score is calculated as a weighted sum:

$$\mathbf{ALP} = (1 - \alpha) \cdot \mathbf{AM}_t + \alpha \cdot \mathbf{L}_{rec_s} \qquad (7)$$

Here, the hyper-parameter $\alpha \in [0, 1]$ is a crucial balancing factor that we dynamically schedule to adapt to the model's learning stage. Following a cosine decay schedule, $\alpha$ gradually increases from a low value at the beginning of training to a higher value towards the end. This curriculum ensures that the masking strategy initially prioritizes the stable guidance from the teacher and later transitions to emphasize the more nuanced difficulty feedback from the student, creating a more robust and efficient training process.

Additionally, in the initial training stages, a high $ALP$ score may not always correlate with clinically relevant importance, as the model is still learning to effectively target the informative regions. To this end, we also adopted a dynamic easy-to-hard mask generation strategy [42, 71], wherein the proportion of important regions within the masked areas gradually increases, thereby providing progressive guidance that directs the model's $ALP$ score incrementally toward important regions. As presented in Fig. 2(b), for each training epoch $t$, $r_t$ of the mask patches are generated by $\mathbf{ALP}^t$, and the remaining $1 - r^t$ of the masked regions in $\mathbf{Mask}^t_{final}$ are randomly selected. Specifically, $r_t = r_0 + \frac{t}{T}(r_t - r_0)$, where $T$ is the total training epochs, and $r_0, r_t \in [0, 1]$ are two tunable hyper-parameters which we empirically set to 0.1 and 0.9, respectively. As demonstrated in Fig. 3, we identify the clinically relevant regions often demonstrate high $ALP$ scores, indicative of rich semantic information for the network to learn.

## 4. Experiments

### 4.1. Datasets and Experimental Settings

We conducted experiments on 42 publicly available ultrasound (US) image datasets, which have a total of 308,584 images covering 12 human organs. The first 38 datasets are utilized for pre-training, where we sample $G = 2$ global views and $L = 8$ local views, with image sizes set to $224 \times 224$ and $96 \times 96$, respectively. We take VMamba-S [45] as the backbone for both the teacher and student networks, and perform 150 epochs of pre-training. In downstream tasks, we perform classification on the BUSI [1] (breast cancer) and Fetal Planes [9] datasets, and segmentation on the BUSBRA [21] (breast lesions) and TN3K [22] (thyroid nodule) datasets. Our implementation is based on

**Algorithm 1** Pseudo-Code of SAM in a PyTorch-like Style.

```
# m_s, m_t: networks for student and teacher
# rat_m: masking ratio; thr_m: masking threshold
# mask_pre: mask from previous epoch output

# teacher attention
t_out, AM = model_t.Mamba(x)
# student reconstruction loss
s_out, rec_x = model_s.Mamba.recon_head(x * mask_pre)
L_rec  = (rec_x - x[~mask_pre]) ** 2
mask = sam_gen(AM, L_rec, rat_m, thr_m)

# self-adaptive masking (SAM) generator
def sma_gen(AM, L_rec, rat_m, thr_m):
    # compute ALP score
    ALP = (1 - α) · AM + α · L_rec
    N = ALP.size(1)
    n_mask_patches = ceil(N * rat_m)
    top_k_mask = floor(N * thr_m)
    topIdx = topk_indices(ALP, top_k_mask)
    n = n_mask_patches - top_k_mask
    if n > 0:
        nonTop = all_indices(N) \ topIdx
        addIdx = random_sample(nonTop, n)
        idx = concat(topIdx, addIdx)
    else:
        idx = topIdx[:, :n_mask_patches]
    mask = zeros(B, N); mask[B, idx] = 1

    return mask
```

iBOT, and more experimental details can be found in *Supplementary Material*.

### 4.2. Comparison with Prior Work

We compare our method with the recent state-of-the-art (SOTA) US image pre-training model, USFM [37], the first pre-training foundation model (FM) developed on two million *private* US images spanning multiple organs and devices. We also report results from other pre-training methods, including DeblurringMAE [39], which introduces a deblurring-based MIM strategy specifically designed for thyroid nodule ultrasound pre-training, and MedSAM [49], which is developed using a large-scale medical image dataset comprising 1.5 million image-mask pairs. We also compare against recent SOTA SSL frameworks, including SimMIM [77], DINO [10], DINOv2 [57], and iBOT [87]. Additionally, we compare the latest supervised learning methods for natural or medical image analysis, including U-Net [62], U-Net++ [88], nnUNet [35], SegFormer [76], UMamba [50], VMamba [45], ResNet50 [29], and ViT [10], for downstream segmentation and classification tasks.

#### 4.2.1. Quantitative Evaluation.

We observe that our method outperforms existing SOTA methods, as shown in Tab. 1 and 2. Particularly, compared to the FM-based USFM, our *OpenUS* achieves improvements of 7.9% DSC and 9.0% IoU on the TN3K segmentation task, as well as 6.5% DSC and 8.7% IoU on the BUS-BRA segmentation task. Additionally, compared to the FM-based DeblurringMAE pretrained on 280K US im-

Table 1. Quantitative results on classification tasks.

| Model | Pretraining | Fetal Planes | | BUSI | |
|---|---|---|---|---|---|
| | | ACC(%) | F1(%) | ACC(%) | F1(%) |
| Supervised | | | | | |
| ResNet50 | ImageNet | $82.6 \pm 0.3$ | $79.5 \pm 0.5$ | $76.3 \pm 0.6$ | $77.1 \pm 0.4$ |
| ViT | ImageNet | $80.3 \pm 0.7$ | $75.4 \pm 0.8$ | $69.8 \pm 0.4$ | $49.2 \pm 0.7$ |
| VMamba | ImageNet | $85.6 \pm 0.5$ | $84.6 \pm 0.3$ | $83.7 \pm 0.5$ | $82.8 \pm 0.3$ |
| SSL/FM | | | | | |
| SimMIM | ImageNet | $87.7 \pm 0.8$ | $86.1 \pm 0.6$ | $80.2 \pm 0.3$ | $78.3 \pm 0.4$ |
| DINO+ViT | ImageNet | $87.8 \pm 0.6$ | $85.3 \pm 0.7$ | $82.3 \pm 0.5$ | $81.7 \pm 0.7$ |
| DINOv2+ViT | ImageNet | $88.4 \pm 0.3$ | $86.1 \pm 0.4$ | $85.2 \pm 0.4$ | $83.5 \pm 0.5$ |
| iBOT+ViT | ImageNet | $88.3 \pm 0.2$ | $86.5 \pm 0.5$ | $85.1 \pm 0.6$ | $83.3 \pm 0.7$ |
| DeblurringMAE | US-Esaote-280K | $88.2 \pm 0.4$ | $85.9 \pm 0.2$ | $84.9 \pm 0.7$ | $82.4 \pm 0.3$ |
| USFM | 3M-US-2M | $90.2 \pm 0.5$ | $89.1 \pm 0.6$ | $87.5 \pm 0.9$ | $85.8 \pm 0.5$ |
| OpenUS | US-308K | $\mathbf{90.3 \pm 0.4}$ | $\mathbf{89.4 \pm 0.3}$ | $\mathbf{87.8 \pm 0.4}$ | $\mathbf{86.3 \pm 0.3}$ |

Table 2. Quantitative results on segmentation tasks.

| Model | Pretraining | BUS-BRA | | TN3K | |
|---|---|---|---|---|---|
| | | DSC(%) | IoU(%) | DSC(%) | IoU(%) |
| Supervised | | | | | |
| U-Net | ImageNet | $87.0 \pm 1.3$ | $79.3 \pm 1.4$ | $75.4 \pm 0.6$ | $63.3 \pm 0.6$ |
| U-Net++ | ImageNet | $86.7 \pm 1.0$ | $79.2 \pm 0.8$ | $76.7 \pm 0.7$ | $66.4 \pm 0.4$ |
| nnUNet | / | $89.8 \pm 1.6$ | $81.2 \pm 1.3$ | $80.4 \pm 0.9$ | $71.4 \pm 0.9$ |
| SegFormer | ImageNet | $85.4 \pm 1.9$ | $74.7 \pm 1.6$ | $78.9 \pm 0.9$ | $68.6 \pm 1.1$ |
| UMamba | / | $89.9 \pm 1.7$ | $81.3 \pm 1.5$ | $80.2 \pm 1.0$ | $70.6 \pm 1.0$ |
| SSL/FM | | | | | |
| SimMIM | ImageNet | $88.1 \pm 0.9$ | $80.4 \pm 1.0$ | $80.3 \pm 1.2$ | $70.5 \pm 1.8$ |
| MedSAM | MedImage-1.57M | $85.6 \pm 1.1$ | $74.7 \pm 1.2$ | $72.8 \pm 2.0$ | $62.2 \pm 1.5$ |
| DeblurringMAE | US-Esaote-280K | $88.8 \pm 1.3$ | $81.3 \pm 1.1$ | $\mathbf{82.9 \pm 2.5}$ | $\mathbf{73.8 \pm 1.0}$ |
| USFM | 3M-US-2M | $84.5 \pm 2.2$ | $74.8 \pm 2.5$ | $74.8 \pm 2.1$ | $64.1 \pm 1.9$ |
| OpenUS | US-308K | $\mathbf{91.0 \pm 0.9}$ | $\mathbf{83.5 \pm 1.0}$ | $82.7 \pm 1.2$ | $73.1 \pm 1.1$ |

ages of the same modality as TN3K, our *OpenUS* achieves comparable performance, with only 0.2% lower DSC and 0.7% lower IoU on the TN3K segmentation task. Moreover, *OpenUS* attains notable improvements of 2.2% in DSC and 2.2% in IoU on the BUS-BRA segmentation task, 2.1% in ACC and 3.5% in F1-score on the Fetal Planes classification task, and 2.9% in ACC and 3.9% in F1-score on the BUSI classification task. These improvements can be attributed to our proposed pre-training approach, which integrates self-adaptive masking and global-local-view mask modeling with contrastive learning. Clearly, compared to other FM/SSL-based methods such as SimMIM and MedSAM, our *OpenUS* demonstrates substantial advantages in downstream segmentation tasks. Although performance gains on the classification task are modest compared to the latest USFM, our *OpenUS* achieves substantial progress using significantly fewer pre-training ultrasound images (296K vs. 2M), while still obtaining slightly better results on the Fetal Planes (by 0.1% ACC and 0.3% F1) and BUSI (by 0.3% ACC and 0.5% F1) classification tasks.

### 4.2.2. Qualitative Evalutation.

We visualize segmentation results in Fig. 4. Clearly, *OpenUS* outperforms other methods in accurately recognizing hazy tissue borders and demonstrates greater robustness to speckle in US images. Specifically, *OpenUS* demonstrates superior visual results in segmenting both large and small thyroid nodule regions ($1^{st}$ rows in Fig. 4). We attribute these improvements to our self-adaptive masking strategy, which effectively captures more precise and clinically relevant details from US images. Similarly, *OpenUS* also exhibits strong performance in segmenting breast tumour ($2^{nd}$ row in Fig. 4).

### 4.2.3. Label Efficiency Analysis.

The label efficiency results for the downstream segmentation tasks are presented in Fig. 5. When developing downstream tasks on training sets with varying label ratios, *OpenUS* consistently outperformed both supervised

and most SSL methods in terms of label efficiency. Even with only a 20% label available, *OpenUS* achieves satisfactory performance, with DSC of 73.2% and 82.7% for segmentation downstream tasks of thyroid nodule and breast tumour. With the label ratio increasing from 20% to 40%, *OpenUS* still demonstrates robust performance, achieving DSC of 77.1% and 85.8% in the thyroid nodule and breast tumor segmentation. When the label ratio increased to 60%, *OpenUS* achieves segmentation performance on thyroid and breast comparable to that with a 100% label ratio. Compared to *OpenUS*, DeblurringMAE, pre-trained on 280K thyroid-only US images, achieved a 0.2%-1.4% higher DSC score across all label ratios on the TN3K segmentation task, likely due to its enhanced feature extraction capability for thyroid nodule images across all label ratios. However, for breast tumor segmentation, DeblurringMAE consistently underperformed relative to *OpenUS* across all label ratios. More results for label efficiency in the classification task can be found in *Supplementary Material*.

### 4.3. Ablation Studies

#### 4.3.1. Global-Local-View Mask Strategies.

Tab. 3 presents the impact of global-view masking and combined global-local-view masking strategies. All experiments use the same masking ratio. For the single-view, different masking strategies are applied to global views. We evaluate our approach against two baseline strategies. The random blockwise mask strategy [7] applies masking uniformly across the spatial domain, while the attention mask selects tokens to mask based on attention maps from global views. In contrast, our proposed self-adaptive mask employs a teacher-student feedback mechanism. This mechanism calculates an $ALP$ score from global views to guide the masking process, significantly outperforming these baselines with 82.1% on segmentation and 90.2% on classification tasks, respectively. In the multi-view setting, applying a self-adaptive mask to global views and a random blockwise mask to local views yields superior performance
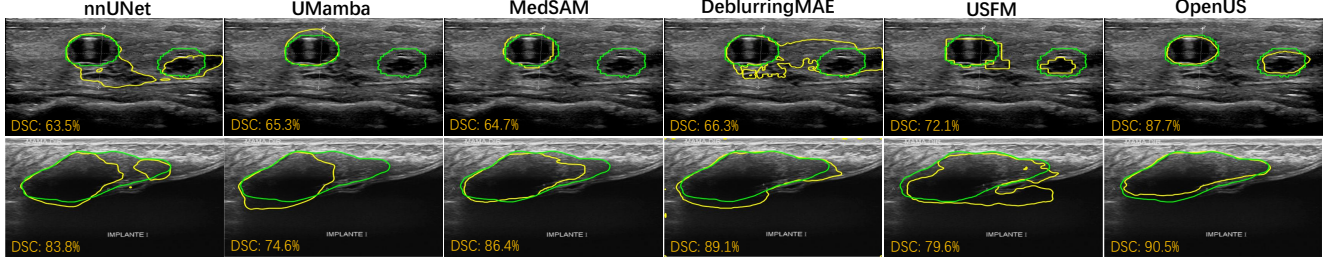
Figure 4. Visualization of US segmentation results on TN3K and BUS-BRA. The ground truth is depicted in green, and the prediction is shown in yellow.
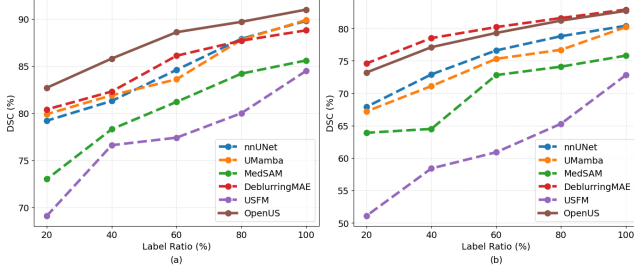


Figure 5. Label efficiency experiments of the downstream segmentation tasks: (a) BUS-BRA and (b) TN3K. The DSC scores (%) of the model trained at different label ratios are reported.

on both classification (90.3% in ACC) and segmentation (82.7% in DSC) tasks. These complementary strategies enable the model to learn more comprehensive details from the US features.

Table 3. Comparison of mask strategies on different views. RBW - Random Blockwise.

| Views | Mask Strategies | | Cla. | Seg. |
|---|---|---|---|---|
| | Global | Local | ACC(%)↑ | DSC(%)↑ |
| Single-view | RBW | ✗ | 88.6 ± 0.2 | 79.1 ± 1.1 |
| | Attention | ✗ | 89.1 ± 0.4 | 80.9 ± 1.0 |
| | Self-adaptive | ✗ | 90.2 ± 0.6 | 81.6 ± 1.3 |
| Multi-view | RBW | RBW | 89.4 ± 0.3 | 81.2 ± 1.4 |
| | Attention | RBW | 89.7 ± 0.5 | 81.7 ± 1.1 |
| | Self-adaptive | RBW | **90.3 ± 0.4** | **82.7 ± 1.2** |

### 4.3.2. Analysis of Components.

We conducted an ablation study to analyse the framework's components using two downstream tasks: Fetal Planes and TN3K. As demonstrated from the first row in Tab. 4, the results indicate that the contrastive learning framework performs reasonably well on classification but lacks robustness in pixel-level tasks, particularly segmentation. Specifically, combining contrastive learning with global and local-view MIM boosted segmentation performance, increasing

the DSC score from 79.1% to 82.7%. Additionally, we observed that the global-view MIM was superior to the local-view MIM on the segmentation task, resulting in a 1.8% increase in DSC. These results highlight the complementary strengths of the two methods: the powerful instance discrimination from contrastive learning and the robust, global and local pixel-level feature acquisition from MIM.

Table 4. Ablation study on components analysis.

| CL | MIM | | Cla. | Seg. |
|---|---|---|---|---|
| | Global | Local | ACC(%) | DSC(%) |
| ✓ | ✗ | ✗ | 89.7 ± 0.3 | 79.1 ± 1.1 |
| ✓ | ✓ | ✗ | 90.2 ± 0.6 | 81.6 ± 1.3 |
| ✓ | ✗ | ✓ | 89.9 ± 0.5 | 80.3 ± 1.4 |
| ✓ | ✓ | ✓ | **90.3 ± 0.4** | **82.7 ± 1.2** |

### 4.3.3. Comparison with Other SSL Methods.

We evaluated the *OpenUS* model against SOTA SSL methods on a downstream breast tumor classification task. To ensure a fair comparison, all frameworks utilized the same VMamba teacher and student networks, were pre-trained on all 38 pre-trained US datasets, and were then tested on the held-out BUSI dataset. As shown in Tab. 5, *OpenUS* outperforms the other SSL frameworks under identical experimental settings.

Table 5. Comparison of different frameworks on the BUSI downstream classification task.

| Framework | Teacher/Student | ACC(%)↑ | F1 (%)↑ |
|---|---|---|---|
| DINO | VMamba | 84.3 ± 0.5 | 83.6 ± 0.6 |
| iBOT | VMamba | 85.6 ± 0.4 | 84.3 ± 0.5 |
| DINOv2 | VMamba | 85.9 ± 0.3 | 84.5 ± 0.4 |
| OpenUS | VMamba | **87.8 ± 0.3** | **86.3 ± 0.5** |

## 5. Conclusion

In this study, we developed a universal US image foundation model named *OpenUS*, characterized by label

efficiency and downstream task adaptability. We propose a self-adaptive masking strategy to capture the clinically relevant and generalizable US features for global and local views. Importantly, *OpenUS* was developed using 42 publicly available US datasets to promote open access, reproducibility and ensure wide applicability.

# References

[1] Walid Al-Dhabyani, Mohammed Gomaa, et al. Dataset of breast ultrasound images. *Data in brief*, 28:104863, 2020. 6, 1

[2] Ali Abbasian Ardakani, Afshin Mohammadi, Mohammad Mirza-Aghazadeh-Attari, and U Rajendra Acharya. An open-access breast lesion ultrasound image database: Applicable in artificial intelligence studies. *Computers in Biology and Medicine*, 152:106438, 2023. 2

[3] Vahid Ashkani Chenarlogh, Mostafa Ghelich Oghli, Ali Shabanzadeh, Nasim Sirjani, Ardavan Akhavan, Isaac Shiri, Hossein Arabi, Morteza Sanei Taheri, and Mohammad Kazem Tarzamni. Fast and accurate u-net model for fetal ultrasound image segmentation. *Ultrasonic imaging*, 44(1):25–38, 2022. 2

[4] Shekoofeh Azizi, Basil Mustafa, et al. Big self-supervised models advance medical image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3478–3488, 2021. 2

[5] Alexei Baevski, Wei-Ning Hsu, et al. Data2vec: A general framework for self-supervised learning in speech, vision and language. In *International conference on machine learning*, pages 1298–1312. PMLR, 2022. 2

[6] Jing Bai, Ilyas Khobo, Yixuan Lu, Dong Ni, Muhammad Yaqub, Karim Lekadir, Jun Ma, and Shuo Li. Landmark detection challenge for intrapartum ultrasound measurement meeting the actual clinical assessment of labor progress. In *Medical Image Computing and Computer Assisted Intervention 2025 (MICCAI)*. Zenodo, 2025. 2

[7] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: Bert pre-training of image transformers. In *Proc. Int. Conf. Learn. Representations*, pages 11–13, 2022. 2, 4, 7

[8] Soumen Basu, Mayank Gupta, et al. Focusmae: Gallbladder cancer detection from ultrasound videos with focused masked autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11715–11725, 2024. 3

[9] Xavier P Burgos-Artizzu, David Coronado-Gutiérrez, et al. Evaluation of deep convolutional neural networks for automatic classification of common maternal fetal ultrasound planes. *Scientific Reports*, 10(1):10200, 2020. 6, 1

[10] Mathilde Caron, Hugo Touvron, et al. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 2, 4, 6

[11] Guang-Hui Chen, Li-Ming Li, Yu-Fan Dai, Jia-Xin Zhang, and Mooi-Hooi Yap. AAU-Net: An Adaptive Attention U-Net for Breast Lesions Segmentation in Ultrasound Images. *IEEE Transactions on Medical Imaging*, 42(5):1289–1300, 2023. 1

[12] Ting Chen, Simon Kornblith, et al. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 2

[13] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning, 2020. 2

[14] C Da et al. Fetal abdominal structures segmentation dataset using ultrasonic images. *Mendeley Data*, 2023. 2

[15] Maria-Antonietta D'Agostino, Lene Terslev, Philippe Aegerter, et al. Scoring ultrasound synovitis in rheumatoid arthritis: a eular-omeract ultrasound taskforce—part 1: definition and development of a standardised, consensus-based scoring system. *RMD open*, 3(1):e000428, 2017. 1

[16] Yi Ding, IEEE Member, Qiqi Yang, Yiqian Wang, Dajiang Chen, Zhiguang Qin, and Jian Zhang. Mallesnet: A multi-object assistance based network for brachial plexus segmentation in ultrasound images. *Medical Image Analysis*, 80: 102511, 2022. 2

[17] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 2

[18] Hanae Elmekki, Ahmed Alagha, Hani Sami, Amanda Spilkin, Antonela Mariel Zanuttini, Ehsan Zakeri, Jamal Bentahar, Lyes Kadem, Wen-Fang Xie, Philippe Pibarot, et al. Cactus: An open dataset and framework for automated cardiac assessment and classification of ultrasound images using deep transfer learning. *Computers in Biology and Medicine*, 190:110003, 2025. 2

[19] Kim-Ann Git. Organ classification on abdominal ultrasound using javascript. https://github.com/ftsvd/USAnotAI, 2020. Accessed: 2020-12-04. 2

[20] Alexander D Gleed, D Mishra, V Chandramohan, Z Fu, Alice Self, S Bhatnagar, Aris T Papageorghiou, and J Alison Noble. Towards multi-sweep ultrasound video understanding: application in detection of breech position using statistical priors. In *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2023. 2

[21] W. Gómez-Flores et al. Bus-bra: a breast ultrasound dataset for assessing computer-aided diagnosis systems. *Medical Physics*, 51(4):3110–3123, 2024. 6, 1

[22] Haifan Gong, Jiaxin Chen, et al. Thyroid region prior guided attention for ultrasound segmentation of thyroid nodules. *Computers in biology and medicine*, 155:106389, 2023. 6, 1

[23] Jean-Bastien Grill, Florian Strub, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 4

[24] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 4

[25] Jie Gui, Tuoliang Chen, et al. A Survey on Self-Supervised Learning: Algorithms, Applications, and Future Trends.

*IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9052–9071, 2024. 2

[26] Yanjun Guo, Xingguang Duan, Chengyi Wang, and Huiqin Guo. Segmentation and recognition of breast ultrasound images based on an expanded u-net. *PLOS ONE*, 16(6): e0253202, 2021. 2

[27] Sami Hamid, Ian A Donaldson, Yipeng Hu, Rachael Rodell, Barbara Villarini, Ester Bonmati, Pamela Tranter, Shonit Punwani, Harbir S Sidhu, Sarah Willis, et al. The smart-target biopsy trial: a prospective, within-person randomised, blinded trial comparing the accuracy of visual-registration and magnetic resonance imaging/ultrasound image-fusion targeted biopsies for prostate cancer risk stratification. *European urology*, 75(5):733–740, 2019. 2

[28] Alexander Hann, Lucas Bettac, Mark M Haenle, Tilmann Graeter, Andreas W Berger, Jens Dreyhaupt, Dieter Schmal-stieg, Wolfram G Zoller, and Jan Egger. Algorithm guided outlining of 105 pancreatic cancer liver metastases in ultra-sound. *Scientific Reports*, 7(1):12779, 2017. 2

[29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6

[30] Kaiming He, Haoqi Fan, et al. Momentum contrast for un-supervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 2

[31] Kaiming He, Xinlei Chen, et al. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3

[32] Qi He, Sophia Bano, Jing Liu, Wentian Liu, Danail Stoy-anov, and Siyang Zuo. Query2: Query over queries for im-proving gastrointestinal stromal tumour detection in an en-doscopic ultrasound. *Computers in Biology and Medicine*, page 106424, 2022. 2

[33] Longhui Huang, Shusheng You, et al. Green hierarchical vi-sion transformer for masked image modeling. In *Advances in Neural Information Processing Systems*, pages 19997–20010, 2022. 2

[34] Ahmed Iqbal and Muhammad Sharif. Memory-efficient transformer network with feature fusion for breast tumor seg-mentation and classification task. *Engineering Applications of Artificial Intelligence*, 127:107292, 2024. 2

[35] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Pe-tersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmen-tation. *Nature methods*, 18(2):203–211, 2021. 6

[36] Hongxu Jiang, Muhammad Imran, Preethika Muralidharan, Anjali Patel, Jake Pensa, Muxuan Liang, Tarik Benidir, Joseph R Grajo, Jason P Joseph, Russell Terry, et al. Mi-crosegnet: A deep learning approach for prostate segmen-tation on micro-ultrasound images. *Computerized Medical Imaging and Graphics*, 112:102326, 2024. 2

[37] Jia Jiao, Jian-jiang Zhou, et al. Usfm: A universal ultrasound foundation model generalized to tasks and organs towards label efficient image analysis. *Medical Image Analysis*, 96: 103202, 2024. 2, 3, 6

[38] Qiang Kang, Qiongjie Lao, et al. Thyroid nodule segmen-tation and classification in ultrasound images through intra- and inter-task consistent learning. *Medical image analysis*, 79:102443, 2022. 1

[39] Qiang Kang, Weijian Li, et al. Deblurring masked image modeling for ultrasound image analysis. *Medical Image Analysis*, 97:103256, 2024. 2, 3, 6

[40] Markus Krönke, Christine Eilers, Desislava Dimova, Melanie Köhler, Gabriel Buschner, Lilit Schweiger, Lemo-nia Konstantinidou, Marcus Makowski, James Nagarajah, Nassir Navab, et al. Tracked 3d ultrasound and deep neu-ral network-based thyroid segmentation reduce interobserver variability in thyroid volumetry. *Plos one*, 17(7):e0268550, 2022. 2

[41] Sarah Leclerc, Erik Smistad, Joao Pedrosa, Andreas Østvik, Frederic Cervenansky, Florian Espinosa, Torvald Espeland, Erik Andreas Rye Berg, Pierre-Marc Jodoin, Thomas Gre-nier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE transac-tions on medical imaging*, 38(9):2198–2210, 2019. 2

[42] Yuheng Li, Tianyu Luan, et al. Anatomask: Enhancing med-ical image segmentation with reconstruction-guided self-masking. In *European Conference on Computer Vision*, pages 146–163. Springer, 2024. 6

[43] Zhaowen Li, Zhiyang Chen, et al. MST: Masked Self-Supervised Transformer for Visual Representation. In *Ad-vances in Neural Information Processing Systems*, pages 13165–13176. Curran Associates, Inc., 2021. 3

[44] Xiao Liu, Fuyuan Zhang, et al. Self-Supervised Learning: Generative or Contrastive. *IEEE Transactions on Knowledge and Data Engineering*, 35(1D):857–876, 2023. 2

[45] Yue Liu, Yunjie Tian, et al. Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063, 2024. 2, 4, 5, 6

[46] Zilu Liu, Jie Gui, and Hai Luo. Good Helper Is around You: Attention-Driven Masked Image Modeling. In *Pro-ceedings of the AAAI Conference on Artificial Intelligence*, pages 1799–1807, 2023. 3

[47] Loshchilov et al. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 3

[48] Ilya Loshchilov et al. Decoupled weight decay regulariza-tion. In *International Conference on Learning Representa-tions*, 2019. 3

[49] Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, and Bo Wang. Segment anything in medical images. *Nature Communications*, 15(1):654, 2024. 6

[50] Jun Ma et al. U-mamba: Enhancing long-range depen-dency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024. 6

[51] Francesco Marzola, Nens van Alfen, Jonne Doorduin, and Kristen Meiburger. DATASET for "Deep learning segmen-tation of transverse musculoskeletal ultrasound images for neuromuscular disease assessment", 2021. 2

[52] James R McLaughlan et al. Lung ultrasound covid phantom dataset used for training machine learning model, 2024. 2

[53] Adrien Meyer, Aditya Murali, Didier Mutter, and Nicolas Padoy. Ultrasam: a foundation model for ultrasound us-

ing large open-access segmentation datasets. *arXiv preprint arXiv:2411.16222*, 2024. 1

[54] Hugo Michard, Bertrand Luvison, Quoc-Cuong Pham, Antonio J Morales-Artacho, and Gaël Guilhem. Aw-net: Automatic muscle structure analysis on b-mode ultrasound images for injury prevention. In *Proceedings of the 12th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, pages 1–9, 2021. 2

[55] Agata Momot. Common carotid artery ultrasound images, 2022. 2

[56] Anna Montoya, Hasnin, kaggle446, shirzad, Will Cukierski, and yffud. Ultrasound nerve segmentation. https://www.kaggle.com/competitions/ultrasound-nerve-segmentation, 2016. Kaggle Competition. 2

[57] Maxime Oquab, Timothée Darcet, et al. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. Featured Certification. 6

[58] David Ouyang, Bryan He, Amirata Ghorbani, Neal Yuan, Joseph Ebinger, Curtis P Langlotz, Paul A Heidenreich, Robert A Harrington, David H Liang, Euan A Ashley, et al. Video-based ai for beat-to-beat assessment of cardiac function. *Nature*, 580(7802):252–256, 2020. 2

[59] Anna Pawłowska, Anna Ćwierz-Pieńkowska, Agnieszka Domalik, Dominika Jaguś, Piotr Kasprzak, Rafał Matkowski, Łukasz Fura, Andrzej Nowicki, and Norbert Żołek. Curated benchmark dataset for ultrasound based breast lesion analysis. *Scientific Data*, 11(1):148, 2024. 2

[60] Lina Pedraza, Carlos Vargas, Fabián Narváez, Oscar Durán, Emma Muñoz, and Eduardo Romero. An open access thyroid ultrasound image database. In *10th International symposium on medical information processing and analysis*, pages 188–193. SPIE, 2015. 2

[61] Charitha D Reddy, Leo Lopez, David Ouyang, James Y Zou, and Bryan He. Video-based deep learning for automated assessment of left ventricular ejection fraction in pediatric patients. *Journal of the American Society of Echocardiography*, 36(5):482–489, 2023. 2

[62] Olaf Ronneberger et al. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 6

[63] Rohit Singla, Cailin Ringstrom, Grace Hu, Victoria Lessoway, Janice Reid, Christopher Nguan, and Robert Rohling. The open kidney ultrasound data set. In *International Workshop on Advances in Simplifying Medical Ultrasound*, pages 155–164. Springer, 2023. 2

[64] Stanford Center for AIMI. Thyroid Ultrasound Cine-clip Dataset. AIMI Shared Dataset, 2021. Contains radiologist-annotated ultrasound cine clips and clinical metadata. 2

[65] Chenxin Tao, Xizhou Zhu, et al. Siamese image modeling for self-supervised vision representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2132–2141, 2023. 3

[66] Abhishek Tyagi, Abhay Tyagi, Manpreet Kaur, Richa Aggarwal, Kapil D Soni, Jayanthi Sivaswamy, and Anjan Trikha. Nerve block target localization and needle guidance for

autonomous robotic ultrasound guided regional anesthesia. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5867–5872. IEEE, 2024. 2

[67] Tamas Ungi, Hastings Greer, Kyle R Sunderland, Victoria Wu, Zachary MC Baum, Christopher Schlenger, Matthew Oetgen, Kevin Cleary, Stephen R Aylward, and Gabor Fichtinger. Automatic spine ultrasound segmentation for scoliosis visualization and measurement. *IEEE Transactions on Biomedical Engineering*, 67(11):3234–3241, 2020. 2

[68] Noelia Vallez, Gloria Bueno, Oscar Deniz, Miguel Angel Rienda, and Carlos Pastor. Bus-uclm: Breast ultrasound lesion segmentation dataset. *Scientific Data*, 12(1):242, 2025. 2

[69] Bram VanBerlo, Tom van Sonsbeek, et al. A survey of the impact of self-supervised pretraining for diagnostic tasks in medical x-ray, ct, mri, and ultrasound. *BMC Medical Imaging*, 24(1):79, 2024. 3

[70] Santiago Vitale, José Ignacio Orlando, Emmanuel Iarussi, and Ignacio Larrabide. Improving realism in patient-specific abdominal ultrasound simulation using cyclegans. *International journal of computer assisted radiology and surgery*, 15(2):183–192, 2020. 2

[71] Haochen Wang, Kaiyou Song, et al. Hard Patches Mining for Masked Image Modeling, 2023. 3, 6

[72] Yixiong Wang, Jing-Jing Li, et al. Deeply-Supervised Networks With Threshold Loss for Cancer Detection in Automated Breast Ultrasound. *IEEE Transactions on Medical Imaging*, 39(4):866–876, 2020. 1

[73] Chen Wei, Hqaoqi Fan, et al. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14668–14678, 2022. 3

[74] Yuxin Wei, Hu Hu, et al. Contrastive learning rivals masked image modeling in fine-tuning via feature distillation, 2022. 2

[75] Lequan Wu, Jiacheng Zhuang, and Huai Chen. Voco: A simple-yet-effective volume contrastive learning framework for 3d medical image analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22873–22882, 2024. 2

[76] Enze Xie, Wenhai Wang, et al. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34: 12077–12090, 2021. 6

[77] Zhenda Xie, Zheng Zhang, et al. SimMIM: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9653–9663, 2022. 3, 6, 1

[78] Yiming Xu, Bowen Zheng, Xiaohong Liu, Tao Wu, Jinxiu Ju, Shijie Wang, Yufan Lian, Hongjun Zhang, Tong Liang, Ye Sang, et al. Improving artificial intelligence pipeline for liver malignancy diagnosis using ultrasound images and video frames. *Briefings in Bioinformatics*, 24(1):bbac569, 2023. 2

[79] Zixun Xu, Ka-Wai Lee, et al. SSL-CPCD: Self-Supervised Learning With Composite Pretext-Class Discrimination for

Improved Generalisability in Endoscopic Image Analysis. *IEEE Transactions on Medical Imaging*, 43(12):4105–4119, 2024. 2

[80] Moi Hoon Yap, Gerard Pons, Joan Marti, Sergi Ganau, Melcior Sentis, Reyer Zwiggelaar, Adrian K Davison, and Robert Marti. Automated breast ultrasound lesions detection using convolutional neural networks. *IEEE journal of biomedical and health informatics*, 22(4):1218–1226, 2017. 2

[81] Lu Yuan, Dongdong Chen, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. 2

[82] Chaoyu Zhang, Haichuan Zheng, and Yafeng Gu. Dive into the details of self-supervised learning for medical image analysis. *Medical Image Analysis*, 89:102879, 2023. 2

[83] H. Zhang, Y. Wu, et al. Echocare: A fully open and generalizable foundation model for ultrasound clinical applications. *arXiv preprint arXiv:2509.11752*, 2025. 2

[84] Shaoting Zhang and Dimitris Metaxas. On the challenges and perspectives of foundation models for medical image analysis. *Medical image analysis*, 91:102996, 2024. 2

[85] Qi Zhao, Shuchang Lyu, Wenpei Bai, Linghan Cai, Binghao Liu, Guangliang Cheng, Meijing Wu, Xiubo Sang, Min Yang, and Lijiang Chen. Mmotu: A multi-modality ovarian tumor ultrasound image dataset for unsupervised cross-domain semantic segmentation. *arXiv preprint arXiv:2207.06799*, 2022. 2

[86] Xiaoyu Zheng, Xu Chen, et al. Xfmamba: Cross-fusion mamba for multi-view medical image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 672–682. Springer, 2025. 2

[87] Jinghao Zhou, Chen Wei, et al. ibot: Image bert pre-training with online tokenizer. In *International Conference on Learning Representations*, pages 1–12, 2022. 3, 4, 6

[88] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE Transactions on Medical Imaging*, 39(6):1856–1867, 2020. 6

[89] Zhemin Zhuang, Nan Li, Alex Noel Joseph Raj, Vijayalakshmi GV Mahesh, and Shunmin Qiu. An rdau-net model for lesion segmentation in breast ultrasound images. *PloS one*, 14(8):e0221535, 2019. 2

# *OpenUS*: A Fully Open-Source Foundation Model for Ultrasound Image Analysis via Self-Adaptive Masked Contrastive Learning

## Supplementary Material

This supplementary material begins by detailing the pre-training datasets in Sec. 6 and providing model implementation details for reproducibility in Sec. 7. Sec. 8 reports additional results on label efficiency, followed by further visualizations and qualitative analysis of segmentation results in Sec. 9. Finally, we include additional ablation experiments in Sec. 10 and a discussion of the study's limitations, future work, and broader impacts in Sec. 11.

## 6. Details of Pre-training Datasets.

Due to its significant domain gap from other medical imaging modalities, developing a foundation model for ultrasound (US) Image Analysis requires a large-scale, specialized dataset. To build the fully open-access dataset, we first followed US-43d [53], then we sourced US data from various online platforms: Kaggle, Mendeley dataset, Zenodo, Google Scholar, GitHub and ResearchGate. The resulting collection comprises 38 datasets for pre-training, encompassing 10 different organs, multiple clinical centers, and various device manufacturers, as shown in Tab. 6.

## 7. Implementation Details.

### 7.1. Pre-training Settings.

The pre-training dataset is augmented using techniques such as random horizontal flipping, color jitter, Gaussian blur, and exposure adjustment. The decoder consists of a lightweight, one-layer head [77], and the model is trained with an $\ell_2$ loss. The proposed model contains 59M parameters. Pre-training was performed on four NVIDIA GH200 GPUs, with detailed training configurations provided in Tab. 7.

### 7.2. Evaluation Methodology.

For downstream fine-tuning, the procedure is as follows: (1) **Classification**: Our analysis was performed on two US datasets: the BUSI dataset [1], with 1,560 images for 3-class classification (benign, malignant, and normal), and the Fetal Planes dataset [9], with 12,400 images for 6-class fetal plane classification (Abdomen, Brain, Femur, Thorax, maternal cervix, and other). The BUSI and Fetal Planes datasets were partitioned into training, validation, and testing sets, comprising 1090/230/240 and 7440/2480/2480 US images, respectively. All input US images were resized to $224 \times 224$. The network architecture consists of a pre-trained backbone and a new, randomly initialized linear classification head, which was subsequently trained for 100
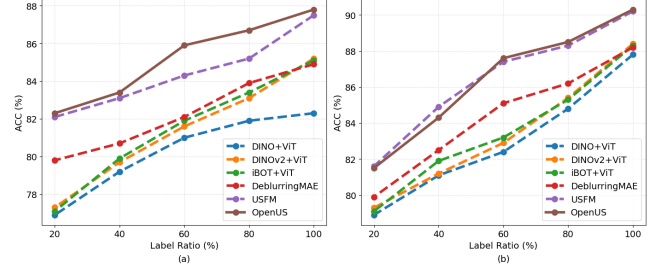


Figure 6. Label efficiency experiments of the downstream classifications tasks: (a) BUSI and (b) Fetal Planes. The ACC values (%) of the model trained at different label ratios are reported.
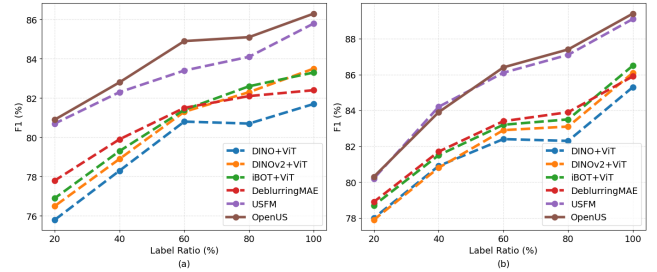


Figure 7. Label efficiency experiments of the downstream classifications tasks: (a) BUSI and (b) Fetal Planes. The F1 scores (%) of the model trained at different label ratios are reported.

epochs. For optimization, we employed the Stochastic Gradient Descent (SGD) with a batch size of 16, a learning rate of 1e-3, and a momentum value of 0.9. Additionally, the primary evaluation metrics for classification performance were Accuracy (ACC) and the F1-score. (2) **Segmentation**: This evaluation utilizes two US datasets: the BUSBRA [21] (breast lesions) and TN3K [22] (thyroid nodule) datasets. The BUSBRA dataset is partitioned into 1,200 training, 299 validation, and 376 testing images, while the TN3K dataset consists of 2,303 training, 576 validation, and 614 testing images. Each US image is annotated with a ground truth mask that delineates the boundaries of the polyp region. We employ the channel-aware Mamba decoder as the segmentation decoder. We optimized the network using the AdamW optimizer, setting the learning rate to 1e-4, weight decay to 5e-2, and the batch size to 16. All input US images are resized to $224 \times 224$ and fine-tuned for 100 epochs. Furthermore, the primary evaluation metrics for segmentation performance were Dice Similarity Coefficient (DSC) and Intersection over Union (IoU).

Table 6. Overview of pre-training ultrasound datasets.

| Dataset | Organ | Number |
|---|---|---|
| 105US [28] | Liver | 105 |
| AbdomenUS [70] | Abdomen | 617 |
| ACOUSLIC [20] | Abdomen | 6620 |
| ASUS [67] | Abdomen | 2865 |
| AUL [78] | Liver | 735 |
| brachial plexus [66] | Nerve | 40788 |
| BrEaST [59] | Breast | 252 |
| BUID [2] | Breast | 232 |
| BUS_UC [34] | Breast | 810 |
| BUS_UCML [68] | Breast | 264 |
| BUS [80] | Breast | 163 |
| CAMUS [41] | Cardiac | 19232 |
| CCAUI [55] | Echocardiogram | 1100 |
| DDTI [60] | Thyroid | 637 |
| EchoNet-Dynamic [58] | Echocardiogram | 20048 |
| EchoNet-Pediatric [61] | Echocardiogram | 15450 |
| FALLMUD [54] | Muscle | 813 |
| FASS [14] | Fetal abdomen | 1588 |
| Fast-U-Net [3] | Abdominal Circumference and Head Circumference | 1411 |
| GIST514-DB [32] | Gastrointestinal Stromal Tumour | 43656 |
| HC [63] | Fetal head | 1334 |
| kidneyUS | Kidney | 487 |
| LUSS_phantom [52] | Lung | 564 |
| MicroSeg [36] | Prostate | 1931 |
| MMOTU-2D [85] | Ovarian Tumor | 1469 |
| MMOTU-3D [85] | Ovarian Tumor | 170 |
| regPro [27] | Prostate | 4706 |
| S1 [26] | Breast | 201 |
| Segthy [40] | Thyroid | 12737 |
| STMUS_NDA [51] | Transverse musculoskeletal | 4355 |
| STU-Hospital [89] | Breast | 42 |
| TG3K [60] | Thyroid | 3585 |
| Thyroid US Cineclip [64] | Thyroid | 17412 |
| UPBD [16] | Brachial Plexus | 955 |
| US nerve Segmentation [56] | Nerve | 11134 |
| USAnotAI [19] | Anatomical regions | 366 |
| Cactus [18] | Cardiac | 37736 |
| IUGC25 [6] | Fetal head | 33466 |

## 8. Additional Label Efficiency Analysis.

The label efficiency results for the downstream classification tasks are presented in Fig. 6 and Fig. 7. In the BUSI classification task, the pre-trained on ImageNet or thyroid nodule ultrasound SSL methods, i.e., DINO, DINOv2, iBOT, and DeblurringMAE, showed strong dependence on the amount of annotated data. At a 20% label ratio, these four methods failed to achieve a satisfactory classification performance, with ACC values and F1 scores below 80%. In the Fetal Planes classification task, all methods besides *OpenUS* and USFM demonstrated suboptimal performance when using a low ratio of labeled data. Specifically,

at a 20% label ratio, their ACC values and F1 scores both fell below 80%.

Compared to *OpenUS*, USFM pre-trained on 2 million private US images shows improved performance over our *OpenUS* for Fetal Planes classification task, particularly with the lower label ratio (20% and 40%). Conversely, at higher label ratios (60%, 80%, and 100%), *OpenUS* surpasses the performance of USFM, achieving ACC of 90.3%, 88.5% and 87.6%, and F1 of 89.4%, 87.5% and 86.3%. This result indicates that *OpenUS* effectively learns features from US images across various organs and possesses strong generalizability, even though it was pre-trained on a much smaller US dataset than USFM.

Table 7. Pre-training settings.

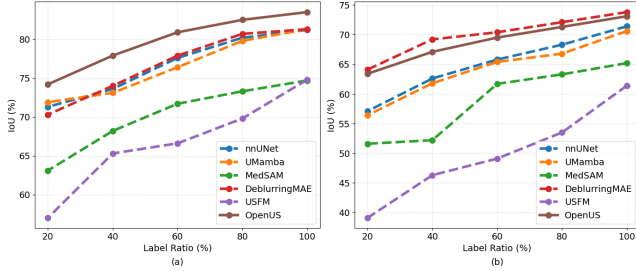| Config | Value |
|---|---|
| optimizer | AdamW [48] |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ |
| weight decay | $4e{-}2$ |
| base learning rate | $5e{-}4$ |
| learning rate schedule | cosine [47] |
| warmup epochs | 30 |
| pretraining epochs | 150 |
| batch size | 128 |
| temperature parameters | $\tau_t, \tau_s = 0.04, 0.07$ |
| masking rate | $r = 0.8$ |
| momentum coefficient | $\lambda = 0.996$ |
| $ALP$ coefficient | $\alpha \in [0.1, 0.9]$ |
| $ALP$ schedule | cosine schedule |



Figure 8. Label efficiency experiments of the downstream segmentation tasks: (a) BUSBRA and (b) TN3K. The IoU (%) of the model trained at different label ratios are reported.

The additional label efficiency results for the downstream segmentation tasks are demonstrated in Fig. 8. *OpenUS* exhibits remarkable label efficiency, consistently outperforming both supervised and the majority of semi-supervised learning (SSL) methods on the BUSBRA and TN3K datasets. Even with a highly limited training set of only 20% of labels, the model attained impressive IoU scores of 74.2% in breast tumor segmentation and 63.4% in thyroid nodule segmentation. The performance of OpenUS scales effectively with the amount of labeled data. Increasing the label ratio to 40% boosts the IoU scores to 77.9% for breast tumor and 67.1% for thyroid nodule segmentation. Remarkably, at just a 80% label ratio, the model's performance becomes comparable to that of a fully supervised model using 100% of the labels. Although DeblurringMAE exhibits slightly better performance in thyroid nodule segmentation, achieving a 0.7% to 2.1% higher IoU across all label ratios—this advantage is likely attributable to its pretraining on the same imaging modality as the thyroid nodule task. Conversely, for the breast tumor segmentation task, *OpenUS* demonstrated superior performance to DeblurringMAE at all levels of supervision.

## 9. Additional Qualitative Evalutation.

Fig. 9 demonstrates the segmentation results of our method and other superivised and self-supervised pre-training methods on the TN3K and BUS-BRA datasets. The results indicate that *OpenUS* outperforms other methods in the accurate recognition of indistinct tissue borders and demonstrates greater robustness to speckle noise in US images. Specifically, for thyroid nodule segmentation tasks involving noisy, low-quality US images ($1^{st}$ and $3^{rd}$ row in Fig. 9), nnUNet, UMamba, MedSAM, Deblurring-MAE, and USFM inaccurately segment artefacts, whereas *OpenUS* remains robust. While some small thyroid nodules occupy very few pixels and are prone to omission ($1^{st}$ row in Fig. 9), *OpenUS* successfully segments these challenging structures. In the breast cancer segmentation task ($2^{nd}$ row, Fig. 9), our method produces segmentation masks with smoother and more continuous edges. In contrast, other methods often yield results corrupted by speckle noise or containing inaccurate lesion boundaries.

## 10. Additional Ablation Studies

### 10.1. Masking Ratio.

The impact of different masking ratios is illustrated in Tab. 8. An increase in the masking ratio from 60% to 80% yields a notable improvement in performance across two downstream tasks. When the masking ratio in US images reaches 80%, the task becomes challenging as the model must learn feature correspondences from a very limited set of visible patches. This forces the model to develop more robust and meaningful representations, which in turn boosts its overall learning capacity. In contrast, the model's performance degrades significantly with a 95% masking ratio, which is the lowest-performing of the four ratios. This indicates that at such a high level of occlusion, the model becomes prohibitively difficult. The extreme scarcity of visible patches prevents the model from effectively learning feature correspondences, leading to a decline in the quality of its learned representations.

Table 8. Comparison of classification and segmentation performance for different self-adaptive mask ratios.

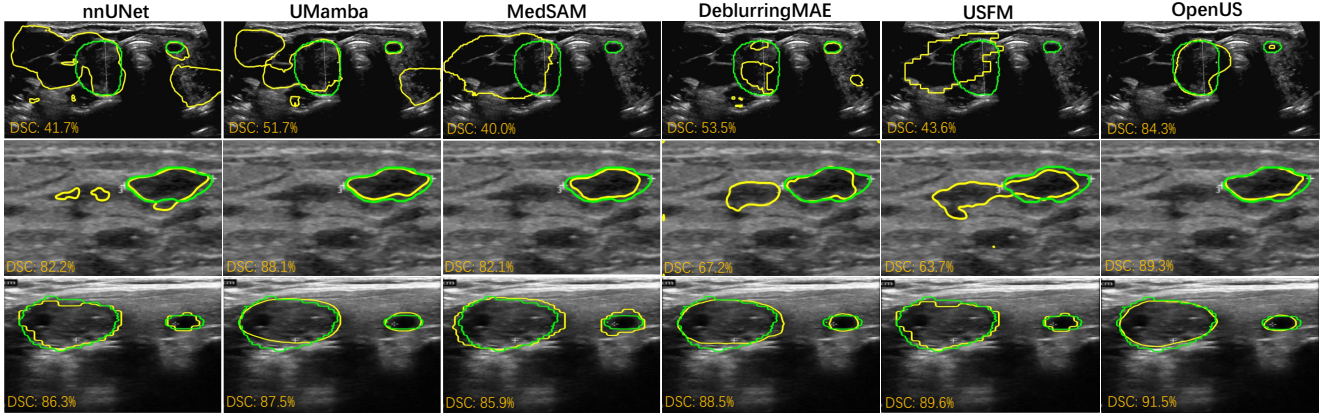| Self-Adaptive | Cla. | Seg. |
|---|---|---|
| Mask Ratio | ACC(%) | DSC(%) |
| 60% | $89.6 \pm 0.2$ | $81.8 \pm 1.4$ |
| 70% | $89.5 \pm 0.3$ | $81.9 \pm 1.1$ |
| 80% | $\mathbf{90.3 \pm 0.4}$ | $\mathbf{82.7 \pm 1.2}$ |
| 90% | $89.0 \pm 0.5$ | $80.3 \pm 1.3$ |

Figure 9. Visualization of US segmentation results on TN3K and BUS-BRA. The ground truth is depicted in green, and the prediction is shown in yellow.

## 11. Limitations, Future Work, and Broader Implications.

### 11.1. Limitations.

We introduce *OpenUS*, a fully open-access foundation model for US image analysis, designed for a novel self-adaptive masked contrastive learning framework. Our *OpenUS* foundation model targets fully open-access, reproducibility, and wide applicability. However, we built *OpenUS* on all public datasets for pre-training and downstream task validation. A key challenge to the reproducibility of *OpenUS* is its reliance on public datasets for pre-training and downstream task validation, as their future availability cannot be guaranteed. Additionally, the extensive pre-training required by *OpenUS* is computationally expensive, demanding significant energy consumption and specialized high-performance hardware, such as GPUs.

### 11.2. Future Work.

First, we will expand the pre-training corpus by integrating additional public US datasets. This will aim to enhance the robustness and generalization capabilities of *OpenUS*. Second, we plan to advance the model beyond its current unimodal (image-only) pre-training by incorporating multi-modal data, such as US videos and corresponding textual information. This multi-modal approach is expected to enrich the feature representations learned during self-supervision, leading to a more powerful model. Finally, we will demonstrate the utility and versatility of our model by fine-tuning the pre-trained *OpenUS* for a broader range of downstream tasks, such as medical image detection, enhancement, and generation.

### 11.3. Broader Implications.

Our approach underscores the potential of SSL for US image analysis. By pre-training on large-scale unlabeled US images, our method mitigates the dependence on costly annotated medical data and improves the label efficiency. The resulting models serve as a versatile foundation for diverse downstream applications, including classification and segmentation, thereby offering a scalable solution to enhance the quality and efficiency of clinical diagnostics.