# U-Netmer: U-Net meets Transformer for medical image segmentation

Sheng He, Rina Bao, P. Ellen Grant, Yangming Ou

*Abstract*—The combination of the U-Net based deep learning models and Transformer is a new trend for medical image segmentation. U-Net can extract the detailed local semantic and texture information and Transformer can learn the long-rang dependencies among pixels in the input image. However, directly adapting the Transformer for segmentation has "token-flatten" problem (flattens the local patches into 1D tokens which losses the interaction among pixels within local patches) and "scale-sensitivity" problem (uses a fixed scale to split the input image into local patches). Compared to directly combining U-Net and Transformer, we propose a new global-local fashion combination of U-Net and Transformer, named U-Netmer, to solve the two problems. The proposed U-Netmer splits an input image into local patches. The global-context information among local patches is learnt by the self-attention mechanism in Transformer and U-Net segments each local patch instead of flattening into tokens to solve the 'token-flatten' problem. The U-Netmer can segment the input image with different patch sizes with the identical structure and the same parameter. Thus, the U-Netmer can be trained with different patch sizes to solve the "scale-sensitivity" problem. We conduct extensive experiments in 7 public datasets on 7 organs (brain, heart, breast, lung, polyp, pancreas and prostate) and 4 imaging modalities (MRI, CT, ultrasound, and endoscopy) to show that the proposed U-Netmer can be generally applied to improve accuracy of medical image segmentation. These experimental results show that U-Netmer provides state-of-the-art performance compared to baselines and other models. In addition, the discrepancy among the outputs of U-Netmer with different scales is linearly correlated to the segmentation accuracy which can be considered as a confidence score to rank test images by difficulty without ground-truth. The code will be available on GitHub.

*Index Terms*—Medical image segmentation, U-Net, Transformer, Image ranking without ground-truth, Deep learning, Confidence score
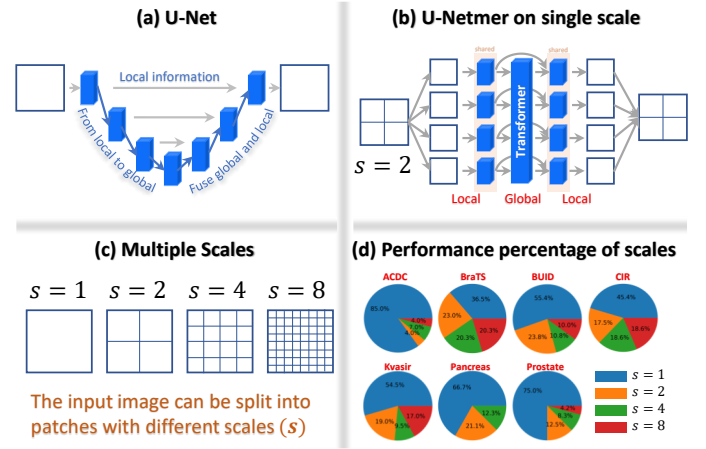


Fig. 1. Sketch of the proposed algorithm. (a) U-Net based models implicitly integrate local information by down-sampling the features. (b) The proposed U-Netmer model on single scale ($s = 2$) explicitly splits the input image into 4 equal-size local patches and uses U-Net to segment each local patch and uses Transformer to fuse the local information among local patches. (c) The input images can be split into patches with different scales $s = 1, 2, 4, 8$. When $s = 1$, the input is the whole image and $s = 2, 4, 8$ means splitting the input image into $s^2$ equal-size and no-overlap patches. (d) The percentage of the best performance achieved with different scales $s$ on test samples of 7 public datasets.

## I. INTRODUCTION

Medical image segmentation aims to use machine learning models (e.g., Convolutional Neural Networks or CNNs for short) to automatically segment the target regions (organs or lesions) from the input medical images with different modalities [1], [2], [3], [4]. One popular backbone of the deep learning model for segmentation is U-Net [5], which is a general CNN model with an encoder and decoder structure (Fig. 1(a)). The encoder path decomposes the input image from local to global deep features where the spatial size is gradually reduced by using the max-pooling operation. As the layer

goes to deep, it extracts the high-level contextual information by discarding the detailed information on each local pixel to remove noise and irrelevant information [6]. To recover the lost detailed spatial information, the decoder path hierarchically fuses global features from the output of the encoder and local features from the intermediate output of the encoder [5] for computing the final segmentation map. In summary, as shown in Fig. 1(a), the U-Net model implicitly extracts features from local information towards global contextual information and fuses these features gradually for the final segmentation output which has the same spatial size as the input image.

Based on the basic structure of the U-Net, many variations have been proposed [4], [7]. For example, U-Net++ [1] uses multiple decoder paths on different scales to fuse the global and local information. Inspired by the attention mechanism which can perform feature recalibration in deep neural networks [8], an attention module has been introduced in U-Net for medical image segmentation [2]. Inspired by the success of vision transformer [9], there is a new trend to integrate Transformer into U-Net to fuse information from different scales or resources to boost the performance of the U-Net [10], [11], [12]. For example, MedT [12] uses a two-

(*Sheng He and Yangming Ou are the corresponding authors.*)
S. He, R. Bao, P. Grant and Y. Ou are with the Boston Children's Hospital and Harvard Medical School, Harvard University, 300 Longwood Ave., Boston, MA, USA. E-mail: heshengxgd@gmail.com; rina.bao@childrens.harvard.edu, ellen.grant@childrens.harvard.edu, yangming.ou@childrens.harvard.edu

branch structure including a global branch to learn global information by CNN-based neural networks and a local branch to learn local information using Transformer. Trans U-Net [10] applies the Transformer on the last layer of the encoder from the U-Net to boost the encoder of the U-Net. UCTransNet [11] uses the Transformer to fuse the intermediate features of U-Net in different scales obtained after max-pooling operations.

Most models of combing the U-Net and Transformer follow the same structure of U-Net and consider the Transformer as a sub-module which can learn the long-range dependencies on deep features with different sizes to boost the performance of segmentation [7]. Few studies use the strategy of cutting the input image into local patches (which is also known as "patchification" [13]) for medical image segmentation (as shown in Fig. 1(c)) which which raises at least two issues existed [12]: (1) **"token-flatten issue"**: the vision Transformer flattens the local patches into 1D tokens, but the 1D functionality losses the interaction of the tokenized information on the local patches [11] and (2) **"scale-sensitivity issue"**: the vision Transformer usually uses a fixed scale to split the input image into patches and the performance of medical image segmentation is sensitive to the scale $s$ when cutting the input image into patches with different sizes (as shown in Fig. 1(c)). Fig. 1(d) shows the percentage of different scales which achieve the best performance on test images of 7 different datasets (the description of datasets can be see in Section III-A and the experiment is described in Section IV-A1). It shows that not all test samples have the highest segmentation accuracy with scale $s = 1$ and some test images achieve the best accuracy with other scales $s = 2, 3, 4$. For example, on BraTS, 36.5% of the test images have the best segmentation accuracy with scale $s = 1$, and 23.0%, 20.3%, 20.3% of the test images achieve the best accuracy with scale $s = 2, 3, 4$, respectively. Thus, directly cutting the input image into patches with a fixed scale and feeding them into Transformer is not necessarily optimal for segmentation.

In this paper, we propose a simple and efficient neural network to optimally combine the U-Net and Transformer for segmentation, named U-Netmer. Similar to vision Transformer, it explicitly splits the input image into local patches and uses Transformer to integrate local information among these patches (shown in Fig. 1(b)). The "patchification" can provide many potential applications for segmentation inspired by the successful application for classification, such as information fusion from different modalities [14], patch dropping and reconstruction for self-supervised learning [15] and few-shot learning for dense prediction [16].

To solve the **"token-flatten issue"**, U-Netmer uses a backbone of a standard segmentation neural network (such as U-Net) to perform the segmentation on local patches instead of flatten each patch into 1D tokens. A standard neural network usually contains *an encoder* and *a decoder* (as shown in Fig. 1(a)). The output deep feature of the encoder is reshaped into 1D tokens and all tokens from the local patches segmented from the input image are concatenated as a sequence of tokens as the input of a standard Transformer [17] (as shown in Fig. 1(b)). The Transformer uses a self-attention mechanism to learn the global-contextual information among local patches

to enhance the segmentation for each local patch.

To solve the **"scale-sensitivity issue"**, an identical U-Netmer with the same parameter is trained on local patches segmented with different scales $s = 1, 2, 4, 8$. Thanks to the flexible structure of the U-Netmer, it can be used on arbitrary patch sizes without any changes of the network structure (as shown in Fig. 3). Multi-scale patches are designed to reduce segmentation's sensitivity to patches at single scale [13].

The main contributions of the work are summarized below:

- We propose U-Netmer which consists of a backbone to extract deep features on local patches and a Transformer block to learn global-context information among local patches. The backbone can be any encoder and decoder structure for segmenting on local patches and we have evaluated three backbones of U-Net [5], Attention U-Net [2] and U-Net++ [1], yielding three variations of U-Netmer.
- U-Netmer is a flexible model which can segment the input image with different patch sizes with identical structure and the same parameters. Jointly training the U-Netmer with different patch sizes can solve the scale-sensitivity problem. Such crafted design and astutely devised training strategies of U-Netmer allow the network to seamlessly imbibe and incorporate bountiful multi-scale contextual knowledge in learning procedures. Therefore, U-Netmer consistently provide better results on 7 public datasets compared to baselines and state-of-the-art models for segmentation.
- U-Netmer can also output the segmentation maps with different scales and the discrepancy of these outputs is linearly correlated to the segmentation accuracy, which can be considered as a confidence score indicating the confidence of the segmentation map and ranks the test images by the difficulty.

## II. METHOD

### A. Basic structure of U-Netmer

Fig. 2 shows the framework of the U-Netmer (with single scale $s = 2$ as an example). U-Netmer can be denoted as $\mathcal{M}_s = (\mathcal{P}_s, \mathcal{E}, \mathcal{T}, \mathcal{D})$, which consists of "patchification" $\mathcal{P}_s$, encoder $\mathcal{E}$, Transformer $\mathcal{T}$, and decoder $\mathcal{D}$, where $s \in [1, 2, 4, 8]$ is the scale (see Fig. 1(c)). Given the input image $x$ with the size of $h \times w$ ($h$ is height and $w$ is width, 2D image as an example), the segmentation output $y$ can be computed by: $y = \mathcal{M}_s(x) = \mathcal{D}(\mathcal{T}(\mathcal{E}(\mathcal{P}_s(x))))$. Each operation is described in the following sections.

*1) Patchification $\mathcal{P}_s$:* Patchification cuts the input image into (typically equal-sized and non-overlapping) patches which is an important step in vision Transformer [9]. Let $s$ be the number of the patches on one side of the input image and the output of the $\mathcal{P}_s$ is a set of $s \times s$ patches (for 2D input image as an example): $\mathbf{p} = \mathcal{P}_s(x)$. The size of each local patch is $h/s \times w/s$. Fig. 1(c) shows the examples of patchification with different scales $s = 1, 2, 4, 8$.

*2) Encoder $\mathcal{E}$:* In vision Transformer [9], the $i$th patch $p_i \in \mathbf{p}, i = 1, 2, ..., s \times s$ is flatten into 1D feature vector. However, for medical image segmentation, the aim is to make
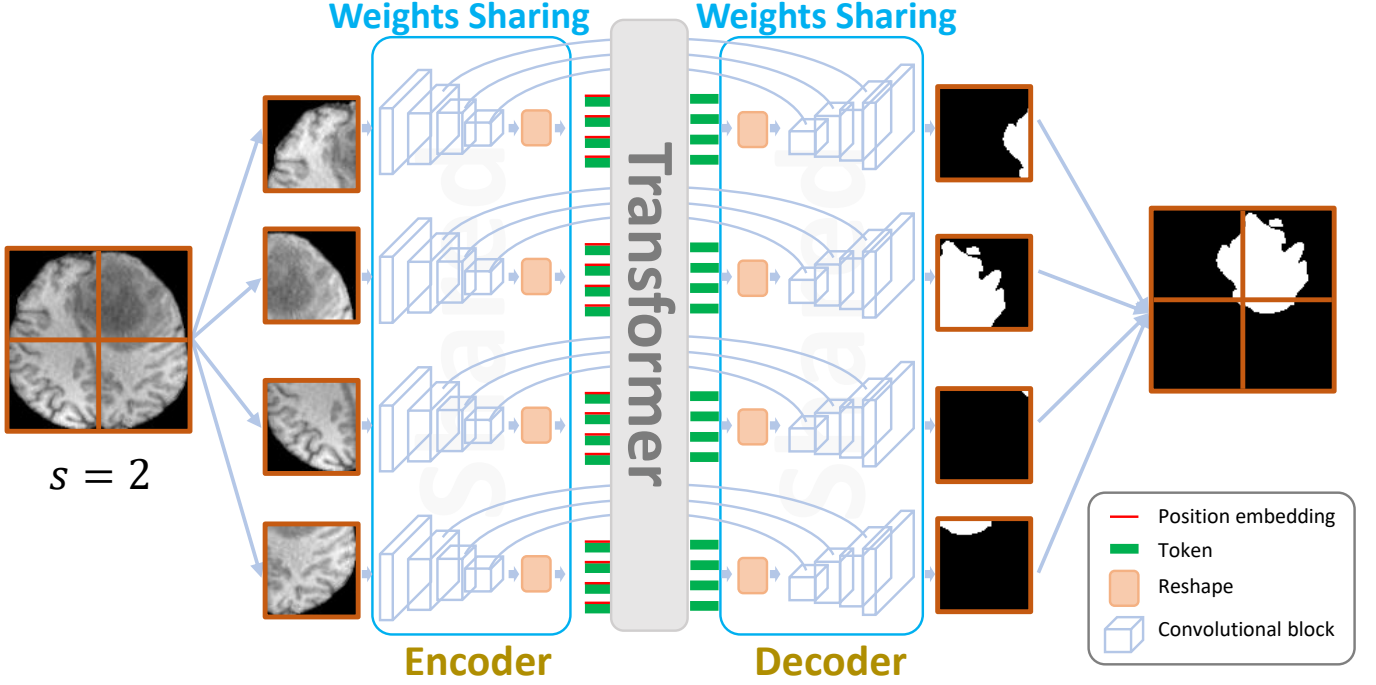
Fig. 2. Framework of the U-Netmer with an example of the scale $s = 2$, indicating 2 patches on each side. The input image is first split into 4 local patches and each local patch is encoded into tokens by an encoder. Tokens of all local patches are fed into Transformer for learning the global context among patches. The global-context enhanced tokens are then decoded by a decoder with the integrated information from the encoder for output prediction on each local patch. The weights on the encoder and decoder are shared among all local patches.

predictions on each pixel within local patches. Thus, we use a segmentation backbone to convert the patches into deep features instead of directly converting the patches into 1D tokens. The aim of using segmentation backbone is to learn the rich information among pixels within local patches for the pixel-level prediction. Given the local patch $p_i \in \mathbf{p}$, the encoder outputs $\mathbf{f}_i, \tau_i = \mathcal{E}(p_i)$, where $\mathbf{f}$ is a set of intermediate deep features while $\tau_i$ is the deep feature from the last layer which contains the deep abstract and contextual information of the input local patch. Any encoder block can be applied here to extract the deep features $\mathbf{f}$ and $\tau$, such as the encoder part of the U-Net [5] and U-Net++ [1], which usually consists of several convolutional layers, followed by Rectified Linear Unit (ReLU), Batch Normalization and Max-pooling layers. The size of $\tau_i$ is $h/(2^n s) \times w/(2^n s)$ where $n$ is the number of max-pooling layer in the encoder $\mathcal{E}$. As shown in Fig. 2, the encoder is shared for all local patches, which can be efficiently computed in parallel.

*3) Transformer $\mathcal{T}$:* Although it is efficient to segment small local patches $\mathbf{p}$, the global-context information of the input image is missed when splitting the image into local patches. To solve this problem, Transformer [17] is used to learn the global-context information among the local patches to enhance the segmentation on each local patch and further improve the accuracy of the final segmentation. We first reshape the deep feature $\tau_i$ obtained from the encoder into a sequence of $h/(2^n s) \times w/(2^n s)$. Note that the number of tokens does not vary with the scale $s$. The tokens from all local patches are concatenated as 1D sequences $\boldsymbol{\tau} = [\tau_1, \tau_2, ..., \tau_{s \times s}]$ with

the number of $s \times s \times h/(2^n s) \times w/(2^n s) = (hw)/2^n$ tokens, indicating the number of tokens does not related to the scale $s$. To keep the position information of local patches, a learnable position embedding vector $\nu$ is added to tokens: $\boldsymbol{\tau} = \boldsymbol{\tau} + \nu$, which is fed into the standard Transformer block [17] with a multi-head self-attention (MSA) $\hat{\boldsymbol{\tau}}_l = \text{MSA}(\hat{\boldsymbol{\tau}}_{l-1}) + \boldsymbol{\tau}_{l-1}$ and a feed-forward network (MLP): $\hat{\boldsymbol{\tau}}_l = \text{MLP}(\hat{\boldsymbol{\tau}}_l) + \hat{\boldsymbol{\tau}}_l$ where $l$ is the number of Transformer block and $\hat{\boldsymbol{\tau}}_0 = \boldsymbol{\tau}$. The detailed information on the multi-head self-attention (MSA) and feed-forward networks (MLP) can be found in studies [17], [9]. The output of Transformer $\hat{\boldsymbol{\tau}} = \mathcal{T}(\boldsymbol{\tau})$ contains the global-context information among all local patches learned by the self-attention mechanism.

*4) Decoder $\mathcal{D}$:* Similar to the encoder $\mathcal{E}$, the decoder aims to fuse the intermediate feature $f_i$ and global-context embedded feature $\hat{\tau}_i \in \hat{\boldsymbol{\tau}}$ to segment each pixel on the local patch $p_i$. The output of the decoder $\mathcal{D}$ is the segmentation result $o_i = \mathcal{D}(\hat{\tau}_i, f_i)$. All the outputs of local patches are stitched together as the final segmentation map $B$ of the input image $x$. The detailed structure of the decoder $\mathcal{D}$ is related to the encoder $\mathcal{E}$, which also consists several convolutional layers, followed by Rectified Linear Unit (ReLU), Batch Normalization and Up-pooling layers. Any decoder block of the typical segmentation neural networks can be applied, such as the encoder part of the U-Net [5], attention U-Net [2], and U-Net++ [1]. As shown in Fig. 2, the decoder is also shared among all local patches which can be efficiently computed in parallel.
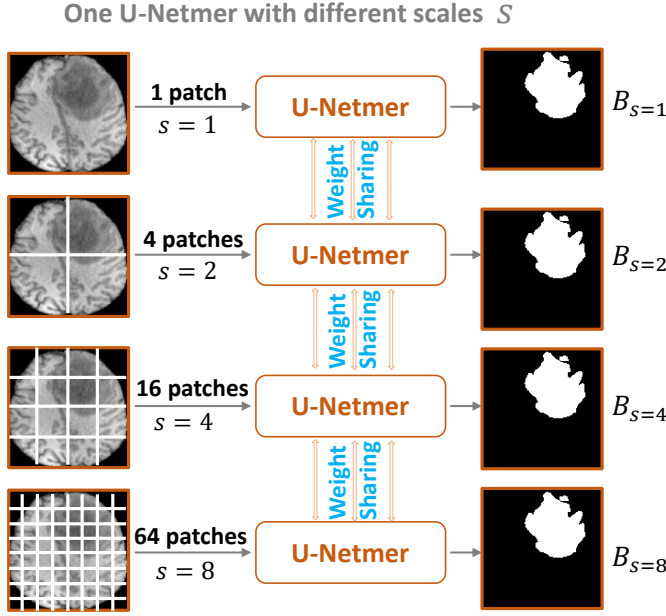
Fig. 3. An illustration of U-Netmer at different scales. The identical U-Netmer (with the same structure and parameters) which can be trained and tested with different scales $s$. $B_{s=i}$ is the segmentation map on scale $s=i$ (where $i \in [1,2,4,8]$).

## B. Variations of U-Netmer

As discussed before, the encoder $\mathcal{E}$ and decoder $\mathcal{D}$ can be obtained from any segmentation backbone. In this paper, we use encoders and decoders from three typical segmentation neural networks: U-Net [5], attention U-Net [2], and U-Net++ [1], yielding to the U-Netmer, attention U-Netmer, and U-Netmer++, respectively. Any other Transformer can be used as Transformer $\mathcal{T}$ block to learn the global-context information among local patches. In this paper, we use the standard Transformer block for simplification and generalization.

## C. Joint training U-Netmer with multiple scales

Segmentation accuracy is sensitive to patch scale $s$ (as show in Fig. 1(d)). To overcome this limitation, we train the U-Netmer with different scale $s$. The reasons are that (1) the encoder $\mathcal{E}$ and decoder $\mathcal{D}$ can compute the deep features on any size of local patches (with a minimal size of $2^n$ due to the $n$ number of max-pooling layers) and (2) the number of tokens $(hw)/2^n$ in Transformer $\mathcal{T}$ does not rely on the scale $s$, indicating that cutting the input image into local patches with different scales $s$ results in the same number of tokens. As shown in Fig. 3, the same model can be trained with different scales $s = 1, 2, 4, 8$ with no added change and cost. Thus, the U-Netmer can learn the information across different patch sizes to boost the segmentation accuracy with an identical setup. The patch size $(h/s \times w/s)$ is small when the scale size is large and we only consider the scale values $s = 1, 2, 4, 8$ for easy computation.

In the following sections, we use $s = \langle i|j|...\rangle$ to denote the U-Netmer which is trained on all local patches segmented with multiple scales $i$, $j$ and others where $i < j$ and

$i, j \in [1, 2, 4, 8]$. For example, U-Netmer++$_{s=\langle 1|2|4\rangle}$ is the U-Netmer++ trained with all local patches split with scales of $s = 1, 2, 4$ from the input image. U-Netmer$_{s=\langle 2\rangle}$ means the U-Netmer is only trained on local patches split with single scale $s = \langle 2 \rangle$.

## III. EXPERIMENTS

### A. Datasets

To evaluate the accuracy of U-Netmer, we conduct experiments on 7 publicly available datasets for medical image segmentation. The datasets used in the experiments include (1) **ACDC** is from the Automated Cardiac Diagnosis Challenge [18] with the purpose of cardiac MRI (CMR) assessment. It consists of 150 cardiac magnetic resonance images with 100 for training and the rest of 50 for testing. (2) **BraTS** is from 2020 Multimodal Brain Tumor Segmentation Challenge [19], [20], [21] with the purpose of segmenting brain tumor (including the peritumoral edema and tumor core) segmentation. 369 scans with four modalities (T1, T1GT, T2, FLAIR) have been split into 295 ($\approx$80%) for training and 74 ($\approx$20%) for testing. (3) **BUID** is from Ultrasound & Breast Ultrasound Images Dataset [22] with the purpose of breast cancer segmentation. There are 780 images which are randomly split into training (80%) and testing (20%) samples. (4) **CIR** [23] consists of 956 CT images on segmented lung nodules from two public datasets which are randomly split into training (80%) and testing (20%) samples. (5) **Kvasir** is from Kvasir-Seg [24] which consists of 1000 polyp images (800 for training and 200 for testing). (6) **Pancreas** [25], [3] consists of 285 CT scans with the purpose of pancreatic parenchyma and mass segmentation. The dataset is randomly split into 228 (80%) scans for training and 57 (20%) scans for testing. (7) **Prostate** is from a Multi-site Dataset [26] which consists 116 prostate T2-weighted MRI from three different sites. The dataset is separated into 80% and 20% for training and testing. For 3D images, we extract 2D slices for training which are stitched into 3D for evaluation. For CT scans, the intensity values are truncated to the range of 5% and 95% percentile to remove the irrelevant details [3]. All images are normalized with zero mean and one standard deviation. Fig. 4 shows examples of images for the 7 datasets.

### B. Neural network training

All models are trained with the PyTorch package, with the Adam optimizer of an initial learning rate 0.0001 which is decayed to half after every 20 epochs. We totally trained 100 epochs with a batch size of 16. The cross-entropy is used as the loss function in the training and Dice score is used as the evaluation metric in the testing. To evaluate the performance of the model itself, no data augmentation or post-processing is applied for all models. All segmentation models, including U-Netmer and other state-of-the-art models, are trained with the same dataset and same training configuration for a fair comparison.
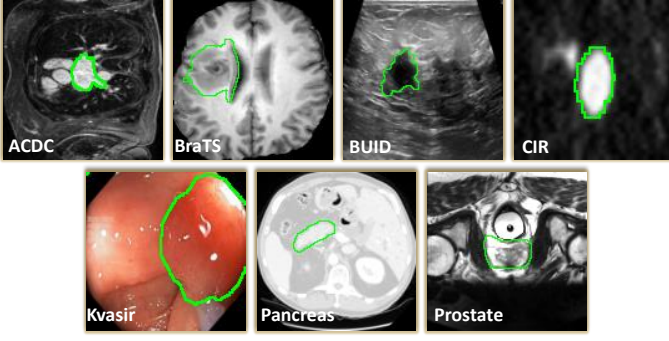
Fig. 4. Examples of images in the 7 datasets used in the experiments. The green contours are the boundary of ground-truth. The 7 public datasets contain images from 7 organs (brain, heart, breast, lung, poly, pancreas, and prostate) and 4 imaging modalities (MRI, CT, Ultrasound, and endoscopy).

## IV. RESULTS

In this section, we present the ablation studies of the U-Netmer with the comparison to state-of-the-art models and its potential application for ranking the test images by difficulty without ground-truth.

### A. Accuracy of U-Netmer

*1) Transformer supplements U-Net:* To evaluate the importance of Transformer $\mathcal{T}$ on U-Netmer with different encoders and decoders, we train models of U-Netmer$_{s=\langle i \rangle}$, attention U-Netmer$_{s=\langle i \rangle}$ and U-Netmer++$_{s=\langle i \rangle}$ with and without Transformer on local patches segmented from input image with a single scale $s = i$. Without Transformer $\mathcal{T}$, the U-Netmer$_{s=\langle i \rangle}$, attention U-Netmer$_{s=\langle i \rangle}$ and U-Netmer++$_{s=\langle i \rangle}$ are similar to original U-Net, attention U-Net, and U-Net++ which are applied on local patches segmented from the input image with scale $s = i$.

Fig. 5 shows the accuracy of these models on 7 datasets. Several observations can be obtained: (1) Models trained with Transformer $\mathcal{T}$ have a higher accuracy than models trained without Transformer $\mathcal{T}$, especially on local patches segmented from scale $i = 2, 4, 8$. The results show that Transformer $\mathcal{T}$ can learn the global-context information among these patches. The results are consistent of three different backbones (U-Net, attention U-Net and U-Net++) over 7 datasets. (2) Unlike the vision Transformer [9], splitting the input image into local patches with a single scale does not improve the accuracy for medical image segmentation on the 7 datasets and the accuracy decreases when patch sizes decrease (the scale $s$ increases) for all models with and without Transformer $\mathcal{T}$. We also plot the percentage of the best performance of U-Netmer with Transformer $\mathcal{T}$ among different scales $s$ on test images (shown in Fig. 1(d)). Results show that most test images achieve the highest accuracy on $s = 1$. For example, 85.0%, 36.5%, 55.4%, 45.4%, 54.5%, 66.7%, and 75.0% of test samples achieve the best accuracy with scale $s = 1$ on ACDC, BraTS, BUID, CIR, Kvasir, Pancreas, and Prostate datasets, respectively, which are larger than other scales $s = 2, 4, 8$. Thus, the average accuracy decreases when the scale $s$ increases in single-scale split of input images.

*2) Joint training U-Netmer with multi-scales improves the accuracy compared with single-scale split:* This section presents the results of the U-Netmer$_{s=\langle i|j|... \rangle}$ trained with local patches segmented with multi-scales. If the U-Netmer is trained only on the single scale $s = \langle 1 \rangle$, the structure of the U-Netmer is the same to the Trans U-Net [10] with U-Net as the backbone. Thus, U-Netmer$_{s=1}$ is the one baseline for comparison. The advantage of the U-Netmer is that it can be also jointly trained with local patches segmented from multi-scales. For example, U-Netmer$_{s=\langle 1|2 \rangle}$ indicates that the U-Netmer is trained on all local patches segmented with both scales $s = 1$ and scale $s = 2$. We train the three variations of U-Netmer$_{s=\langle 1 \rangle}$ (baseline), U-Netmer$_{s=\langle 1|2 \rangle}$, U-Netmer$_{s=\langle 1|2|4 \rangle}$, and U-Netmer$_{s=\langle 1|2|4|8 \rangle}$. When applying the trained U-Netmer with multi-scale patches, different segmentation outputs $B_{s=i}$ can also be obtained on the corresponding scale $s = i$ (see Fig. 3). We evaluate the accuracy of each output $B_{s=i}$ and the results are shown in Fig. 6. During testing, the accuracy of the output $B_{s=i}$ slightly decreases when scale $s = i$ increases. For three models of U-Netmer$_{s=\langle 1|2 \rangle}$, U-Netmer$_{s=\langle 1|2|4 \rangle}$ and U-Netmer$_{s=\langle 1|2|4|8 \rangle}$, the best results is achieved by the output of $B_{s=1}$ of U-Netmer jointly trained with multi-scale split, which are reported in the following sections.

Fig. 7 shows the accuracy of $B_{s=1}$ obtained from U-Netmer models trained with different number of scales: U-Netmer$_{s=\langle 1 \rangle}$, U-Netmer$_{s=\langle 1|2 \rangle}$, U-Netmer$_{s=\langle 1|2|4 \rangle}$ and U-Netmer$_{s=\langle 1|2|4|8 \rangle}$. Three variations of the U-Netmer have consistent accuracy on datasets ACDC, BraTS, Pancreas and Prostate. For ACDC, the best accuracy is achieved by U-Netmer$_{s=\langle 1|2 \rangle}$ and for BraTS, Pancreas and Prostate, the best accuracy is obtained on U-Netmer$_{s=\langle 1|2|4 \rangle}$. For Kvasir, U-Netmer$_{s=\langle 1|2|4 \rangle}$ provides the best performance with U-Net as the backbone while the highest accuracies are achieved by attention U-Netmer$_{s=\langle 1|2|4|8 \rangle}$ and U-Netmer++$_{s=\langle 1|2|4|8 \rangle}$ with attention U-Net and U-Net++ as the backbone, respectively. A similar trend is found on BUID and CIR where the best accuracy is achieved on different scales for different backbones. In general, the results on Fig. 7 show that training U-Netmer with multi-scales with different backbones can improve the performance on all 7 datasets, providing higher accuracies than training the U-Netmer$_{s=\langle 1 \rangle}$ which is the baseline model. Table I shows the accuracy of the U-Netmer with single-scale split and multi-scale split. The results show that training the U-Netmer with multi-scale split improves the accuracy. Similar results have also found on attention U-Netmer and U-Netmer++.

*3) U-Netmer trained with multi-scales outperforms state-of-the-art models:* We first compare three variations of U-Netmer with their corresponding baselines: U-Net [5], Attention U-Net [2], U-Net++ [1] and Trans U-Net [10]. Table II shows the accuracy measured by Jaccard index, Dice coefficient, pixelwise accuracy, sensitivity and and specificity on 7 datasets. Results show that U-Netmer outperforms its baseline models in most cases. We also conduct the comparison study between the U-Netmer and other state-of-the-art models including other pure U-Net based models (such as BiOnet [27], ConvUNeXt [28], ResUnet [29]) and U-Net with Transformer models (such as UNext [30], UCTransNet [11], MedT [12]).
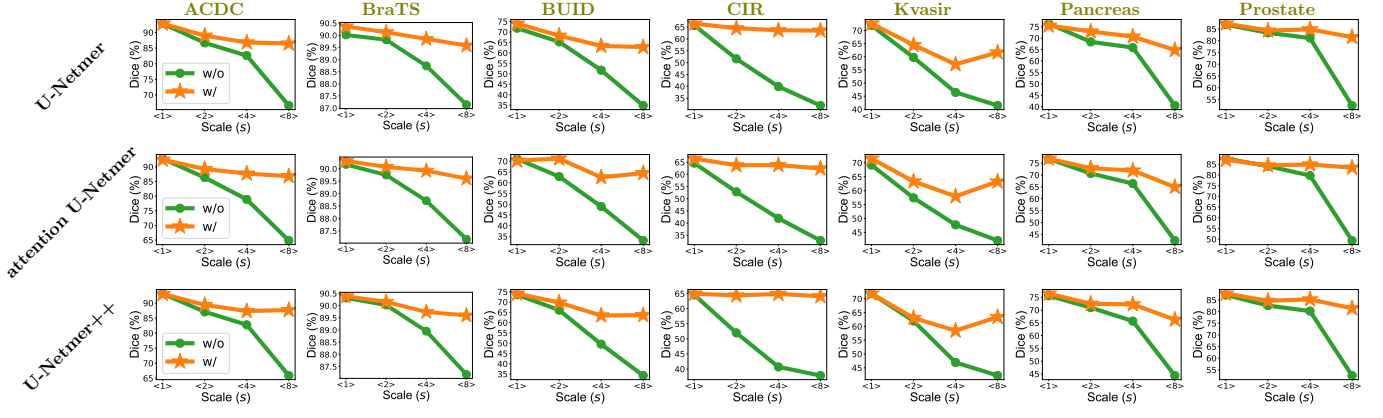
Fig. 5. Dice accuracy comparison of different variations of U-Netmer trained with a single scale $s = \langle i \rangle$ with Transformer $\mathcal{T}$ (w/, orange lines) and without Tranformer (w/o, green lines) on 7 datasets. The accuracy significantly drops for models without Transformer when the scale $s = i$ increases, indicating that the global-context information learned by Transformer is important for segmentation, especially with a high scale $s = \langle i \rangle$.
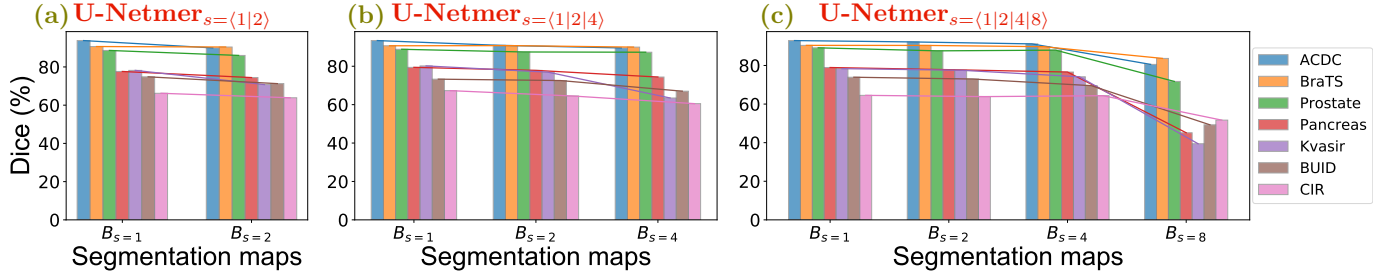


Fig. 6. The segmentation accuracy of different outputs $B_{s=i}$ from the joint training of (a) U-Netmer$_{s=\langle 1|2 \rangle}$, (b) U-Netmer$_{s=\langle 1|2|4 \rangle}$, and (c) U-Netmer$_{s=\langle 1|2|4|8 \rangle}$.
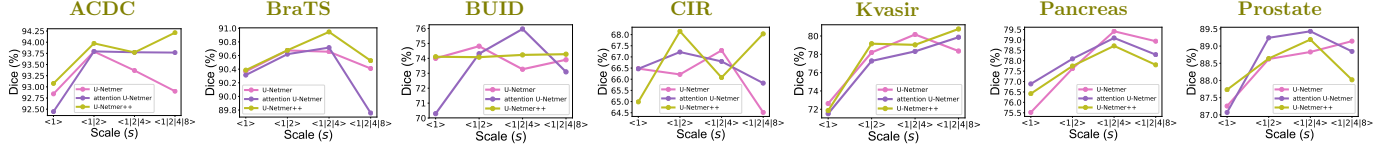


Fig. 7. Accuracy of U-Netmer trained with multiple scales on 7 datasets. $s = \langle 1|2 \rangle$ indicates that models trained on local patches segmented with scales $s = 1, 2$. The same definition to $s = \langle 1 \rangle$ (baseline), $s = \langle 1|2|4 \rangle$ and $s = \langle 1|2|4|8 \rangle$.

TABLE I
THE DICE PERFORMANCE OF U-NETMER WITH SINGLE SCALE $s = \langle i \rangle$, $i = 1, 2, 4, 8$ AND MULTI-SCALE $s = \langle 1|2 \rangle$, $s = \langle 1|2|4 \rangle$, AND $s = \langle 1|2|4|8 \rangle$ ON 7 DATASETS.

| | Scale | ACDC | BUID | BraTS | CIR | Kvasir | Pancreas | Prostate |
|---|---|---|---|---|---|---|---|---|
| Single-Scale | $s = \langle 1 \rangle$ | $92.84_{\pm 4.11}$ | $74.00_{\pm 26.47}$ | $90.36_{\pm 5.73}$ | $66.48_{\pm 23.00}$ | $72.63_{\pm 26.56}$ | $75.52_{\pm 9.18}$ | $87.26_{\pm 3.97}$ |
| | $s = \langle 2 \rangle$ | $88.98_{\pm 7.69}$ | $68.22_{\pm 28.59}$ | $90.13_{\pm 5.54}$ | $64.51_{\pm 22.08}$ | $64.49_{\pm 25.85}$ | $72.92_{\pm 10.09}$ | $84.29_{\pm 5.68}$ |
| | $s = \langle 4 \rangle$ | $86.81_{\pm 11.64}$ | $63.33_{\pm 27.89}$ | $89.85_{\pm 5.88}$ | $63.66_{\pm 21.71}$ | $57.14_{\pm 27.05}$ | $70.69_{\pm 10.98}$ | $84.82_{\pm 5.26}$ |
| | $s = \langle 8 \rangle$ | $86.46_{\pm 11.52}$ | $62.78_{\pm 28.89}$ | $89.58_{\pm 6.61}$ | $63.52_{\pm 21.06}$ | $61.63_{\pm 25.85}$ | $64.84_{\pm 13.10}$ | $81.52_{\pm 5.28}$ |
| Multi-Scale | $s = \langle 1|2 \rangle$ | $\mathbf{93.79}_{\pm 3.30}$ | $\mathbf{74.82}_{\pm 26.40}$ | $\mathbf{90.67}_{\pm 5.25}$ | $66.22_{\pm 21.62}$ | $78.20_{\pm 22.58}$ | $77.62_{\pm 8.67}$ | $88.62_{\pm 3.52}$ |
| | $s = \langle 1|2|4 \rangle$ | $93.37_{\pm 4.83}$ | $73.27_{\pm 28.61}$ | $90.66_{\pm 5.48}$ | $\mathbf{67.29}_{\pm 21.65}$ | $\mathbf{80.16}_{\pm 20.94}$ | $\mathbf{79.42}_{\pm 7.59}$ | $88.83_{\pm 3.30}$ |
| | $s = \langle 1|2|4|8 \rangle$ | $92.90_{\pm 5.04}$ | $73.91_{\pm 25.71}$ | $90.41_{\pm 5.27}$ | $64.51_{\pm 23.14}$ | $78.36_{\pm 21.22}$ | $78.94_{\pm 7.90}$ | $\mathbf{89.14}_{\pm 3.46}$ |

All of these models are trained with the same training setup for a fair comparison. Table III shows the accuracies on the 7 datasets which shows that U-Netmer based methods (U-Netmer, Attention U-Netmer and U-Netmer++) provide higher

accuracy than other models. U-Netmer provides the highest accuracy on the Pancreas dataset, Attention U-Netmer provides the highest accuracy on the BUID, Prostate datasets and U-Netmer++ provides the highest accuracy on ACDC, BraTS,

TABLE II
ACCURACY COMPARISON IN TERMS OF JACCARD INDEX, DICE,
PIXEL-WISE ACCURACY, SENSITIVITY AND SPECIFICITY BETWEEN THE
U-NETMER AND CORRESPONDING BASELINES ON THE 8 DATASETS.

| ACDC | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
|---|---|---|---|---|---|
| U-Net [5] | $86.88_{\pm6.42}$ | $92.84_{\pm4.06}$ | $99.38_{\pm0.31}$ | $92.04_{\pm6.54}$ | $99.74_{\pm0.16}$ |
| Attention U-Net [2] | $86.78_{\pm5.88}$ | $92.81_{\pm3.58}$ | $99.37_{\pm0.31}$ | $92.15_{\pm5.89}$ | $99.72_{\pm0.16}$ |
| U-Net++ [1] | $87.41_{\pm5.96}$ | $93.16_{\pm3.68}$ | $99.41_{\pm0.27}$ | $92.70_{\pm5.69}$ | $99.73_{\pm0.14}$ |
| Trans U-Net [10] | $86.89_{\pm6.52}$ | $92.84_{\pm4.11}$ | $99.38_{\pm0.30}$ | $92.15_{\pm6.65}$ | $99.73_{\pm0.16}$ |
| U-Netmer | $88.48_{\pm5.52}$ | $93.79_{\pm3.30}$ | $99.46_{\pm0.27}$ | $\mathbf{93.12}_{\pm5.66}$ | $99.76_{\pm0.13}$ |
| Attention U-Netmer | $88.55_{\pm6.40}$ | $93.79_{\pm4.01}$ | $99.46_{\pm0.30}$ | $93.06_{\pm6.44}$ | $99.77_{\pm0.11}$ |
| U-Netmer++ | $\mathbf{89.23}_{\pm5.49}$ | $\mathbf{94.21}_{\pm3.26}$ | $\mathbf{99.49}_{\pm0.34}$ | $92.89_{\pm4.73}$ | $\mathbf{99.80}_{\pm0.20}$ |
| BraTS | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $82.31_{\pm8.66}$ | $90.01_{\pm6.00}$ | $98.41_{\pm0.79}$ | $88.33_{\pm9.18}$ | $\mathbf{99.32}_{\pm0.48}$ |
| Attention U-Net [2] | $82.51_{\pm8.12}$ | $90.17_{\pm5.49}$ | $98.43_{\pm0.72}$ | $88.74_{\pm8.24}$ | $99.30_{\pm0.48}$ |
| U-Net++ [1] | $82.72_{\pm7.96}$ | $90.31_{\pm5.36}$ | $98.44_{\pm0.73}$ | $89.12_{\pm8.26}$ | $99.28_{\pm0.49}$ |
| Trans U-Net [10] | $82.86_{\pm8.37}$ | $90.36_{\pm5.73}$ | $98.46_{\pm0.76}$ | $89.09_{\pm8.62}$ | $99.30_{\pm0.47}$ |
| U-Netmer | $83.31_{\pm7.82}$ | $90.67_{\pm5.25}$ | $98.50_{\pm0.75}$ | $89.38_{\pm8.23}$ | $\mathbf{99.32}_{\pm0.47}$ |
| Attention U-Netmer | $83.38_{\pm7.82}$ | $90.72_{\pm5.24}$ | $98.49_{\pm0.75}$ | $90.13_{\pm7.78}$ | $99.26_{\pm0.48}$ |
| U-Netmer++ | $\mathbf{83.76}_{\pm7.72}$ | $\mathbf{90.95}_{\pm5.17}$ | $\mathbf{98.52}_{\pm0.73}$ | $90.37_{\pm7.57}$ | $99.29_{\pm0.46}$ |
| BUID | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $62.65_{\pm29.13}$ | $71.81_{\pm29.57}$ | $95.93_{\pm5.07}$ | $72.84_{\pm31.58}$ | $98.63_{\pm1.77}$ |
| Attention U-Net [2] | $61.70_{\pm29.36}$ | $71.09_{\pm29.20}$ | $95.87_{\pm5.20}$ | $70.45_{\pm31.22}$ | $98.88_{\pm1.48}$ |
| U-Net++ [1] | $63.73_{\pm26.85}$ | $73.58_{\pm26.51}$ | $96.00_{\pm4.74}$ | $74.62_{\pm28.23}$ | $98.61_{\pm1.83}$ |
| Trans U-Net [10] | $64.33_{\pm27.11}$ | $74.00_{\pm26.47}$ | $\mathbf{96.08}_{\pm4.93}$ | $74.16_{\pm28.44}$ | $98.70_{\pm1.68}$ |
| U-Netmer | $65.41_{\pm27.17}$ | $74.82_{\pm26.40}$ | $95.96_{\pm5.16}$ | $76.47_{\pm28.12}$ | $98.60_{\pm1.78}$ |
| Attention U-Netmer | $\mathbf{66.35}_{\pm25.81}$ | $\mathbf{75.97}_{\pm24.95}$ | $95.95_{\pm5.17}$ | $\mathbf{76.60}_{\pm26.85}$ | $98.61_{\pm1.88}$ |
| U-Netmer++ | $64.64_{\pm27.10}$ | $74.28_{\pm26.24}$ | $95.90_{\pm5.05}$ | $72.75_{\pm28.65}$ | $\mathbf{98.90}_{\pm1.46}$ |
| CIR | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $52.80_{\pm22.71}$ | $65.71_{\pm23.24}$ | $97.57_{\pm2.86}$ | $67.59_{\pm25.51}$ | $98.98_{\pm1.20}$ |
| Attention U-Net [2] | $52.09_{\pm24.09}$ | $64.62_{\pm24.89}$ | $97.57_{\pm2.90}$ | $66.38_{\pm27.57}$ | $99.00_{\pm1.21}$ |
| U-Net++ [1] | $52.25_{\pm24.53}$ | $64.57_{\pm25.64}$ | $97.57_{\pm2.93}$ | $65.58_{\pm28.20}$ | $99.03_{\pm1.24}$ |
| Trans U-Net [10] | $53.61_{\pm22.58}$ | $66.48_{\pm23.00}$ | $\mathbf{97.68}_{\pm2.68}$ | $67.78_{\pm25.63}$ | $99.04_{\pm1.24}$ |
| U-Netmer | $54.21_{\pm21.86}$ | $67.29_{\pm21.65}$ | $97.60_{\pm2.93}$ | $\mathbf{70.67}_{\pm24.32}$ | $98.96_{\pm1.25}$ |
| Attention U-Netmer | $54.13_{\pm21.92}$ | $67.22_{\pm21.61}$ | $97.58_{\pm2.82}$ | $69.90_{\pm24.05}$ | $98.99_{\pm1.14}$ |
| U-Netmer++ | $\mathbf{55.02}_{\pm21.35}$ | $\mathbf{68.15}_{\pm20.96}$ | $97.65_{\pm3.02}$ | $70.31_{\pm23.86}$ | $\mathbf{99.05}_{\pm1.11}$ |
| Kvasir | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $62.01_{\pm28.11}$ | $71.88_{\pm27.34}$ | $92.31_{\pm10.00}$ | $71.91_{\pm30.10}$ | $97.87_{\pm3.69}$ |
| Attention U-Net [2] | $59.11_{\pm28.89}$ | $69.18_{\pm28.65}$ | $91.92_{\pm10.27}$ | $67.73_{\pm31.57}$ | $\mathbf{98.24}_{\pm2.87}$ |
| U-Net++ [1] | $61.42_{\pm26.88}$ | $71.87_{\pm25.87}$ | $92.12_{\pm9.93}$ | $71.62_{\pm28.50}$ | $97.81_{\pm3.38}$ |
| Trans U-Net [10] | $62.65_{\pm27.48}$ | $72.63_{\pm26.56}$ | $92.53_{\pm9.80}$ | $71.85_{\pm29.05}$ | $97.95_{\pm3.64}$ |
| U-Netmer | $70.95_{\pm23.64}$ | $80.16_{\pm20.94}$ | $93.87_{\pm8.40}$ | $82.80_{\pm22.71}$ | $97.63_{\pm4.26}$ |
| Attention U-Netmer | $70.57_{\pm23.62}$ | $79.86_{\pm21.17}$ | $93.91_{\pm8.16}$ | $83.19_{\pm23.31}$ | $97.42_{\pm3.85}$ |
| U-Netmer++ | $\mathbf{71.66}_{\pm23.34}$ | $\mathbf{80.76}_{\pm20.40}$ | $\mathbf{93.98}_{\pm8.10}$ | $\mathbf{84.53}_{\pm21.48}$ | $97.35_{\pm3.96}$ |
| Pancreas | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $62.51_{\pm10.63}$ | $76.38_{\pm8.53}$ | $99.39_{\pm0.26}$ | $75.64_{\pm12.05}$ | $99.74_{\pm0.16}$ |
| Attention U-Net [2] | $62.95_{\pm10.08}$ | $76.76_{\pm8.03}$ | $99.41_{\pm0.24}$ | $75.18_{\pm12.82}$ | $99.77_{\pm0.13}$ |
| U-Net++ [1] | $61.52_{\pm10.56}$ | $75.61_{\pm8.64}$ | $99.37_{\pm0.26}$ | $74.88_{\pm13.07}$ | $99.73_{\pm0.16}$ |
| Trans U-Net [10] | $61.47_{\pm10.98}$ | $75.52_{\pm9.18}$ | $99.38_{\pm0.25}$ | $73.79_{\pm13.71}$ | $99.75_{\pm0.15}$ |
| U-Netmer | $\mathbf{66.46}_{\pm9.63}$ | $\mathbf{79.42}_{\pm7.59}$ | $\mathbf{99.47}_{\pm0.23}$ | $\mathbf{78.49}_{\pm11.77}$ | $99.78_{\pm0.11}$ |
| Attention U-Netmer | $66.02_{\pm9.58}$ | $79.09_{\pm7.63}$ | $99.46_{\pm0.23}$ | $77.67_{\pm11.63}$ | $99.79_{\pm0.13}$ |
| U-Netmer++ | $65.50_{\pm9.69}$ | $78.71_{\pm7.56}$ | $99.46_{\pm0.21}$ | $76.62_{\pm12.87}$ | $\mathbf{99.80}_{\pm0.11}$ |
| Prostate | Jaccard index | Dice | Accuracy | Sensitivity | Specificity |
| U-Net [5] | $76.91_{\pm5.59}$ | $86.83_{\pm3.74}$ | $98.99_{\pm0.50}$ | $86.62_{\pm6.78}$ | $99.58_{\pm0.23}$ |
| Attention U-Net [2] | $78.17_{\pm5.57}$ | $87.63_{\pm3.66}$ | $99.05_{\pm0.48}$ | $87.64_{\pm6.36}$ | $99.59_{\pm0.21}$ |
| U-Net++ [1] | $77.38_{\pm5.30}$ | $87.14_{\pm3.58}$ | $99.01_{\pm0.48}$ | $86.60_{\pm6.55}$ | $99.60_{\pm0.22}$ |
| Trans U-Net [10] | $77.60_{\pm5.82}$ | $87.26_{\pm3.97}$ | $99.03_{\pm0.47}$ | $86.65_{\pm6.53}$ | $99.62_{\pm0.22}$ |
| U-Netmer | $80.58_{\pm5.39}$ | $89.14_{\pm3.46}$ | $99.15_{\pm0.56}$ | $87.01_{\pm6.91}$ | $\mathbf{99.73}_{\pm0.10}$ |
| Attention U-Netmer | $\mathbf{81.02}_{\pm4.85}$ | $\mathbf{89.43}_{\pm3.13}$ | $\mathbf{99.19}_{\pm0.41}$ | $\mathbf{88.71}_{\pm5.27}$ | $99.68_{\pm0.14}$ |
| U-Netmer++ | $80.63_{\pm4.94}$ | $89.19_{\pm3.13}$ | $99.16_{\pm0.49}$ | $87.23_{\pm6.11}$ | $\mathbf{99.73}_{\pm0.13}$ |

CIR, and Kvasir datasets.

*4) Deep features learned by U-Netmer have higher segmentation ability than those from its counterpart (U-Net):* This section presents the segemtnation ability (SA) scores of each layer of the U-Netmer$_{s=\langle1|2|4\rangle}$ and the original U-Net computed by the ProtoSeg [6] which can compute the binary segmentation map on deep features based on prototypes of background and target regions. There are 18 layers in a typical U-Net and the U-Netmer with U-Net as the backbone. The ProtoSeg can also be applied on input image which is considered as a specifical feature map [6]. Fig. 8 shows the SA scores of U-Net and U-Netmer including input images (index 0), 18 deep features (index 1-19) and the final segmentation output (index 20). It can be seen from the figure that the segmentation ability scores of the U-Netmer trained with different patch sizes are higher than the corresponding SA scores of the U-Net trained with the whole input image, especially on ACDC, BUID, CIR, Kvasir, and Prostate datasets. The results show that U-Netmer trained with multi-scales can improve the segmentation ability on deep features as well as on the final output.

*5) Discussion:* Experimental results show that the vanilla U-Net [5] provides a moderate accuracy based on small input patches which contain limited information for separating the target region and background. Applying a Transformer block among the local patches segmented from the same input image can significantly improve the accuracy (see Fig. 5) because the local patches can learn the global-context information by the self-attention mechanism in Transformer [31]. This inspires the design of the U-Netmer which can make a prediction based on local patches with different patch sizes without changes to the model structure and parameters. However, training the U-Netmer with a single scale (a fixed patch size) does not improve the accuracy (see Fig. 5) due to the fact that segmentation is sensitive to the scales. To solve this problem, we train the U-Netmer with local patches segmented with multi-scales which can improve the accuracies for segmentation. The same results are also found on the three different variations (see Fig. 7). The U-Netmer trained by patches with different sizes (multi-scales) provides the highest accuracy compared to four different baselines and other U-Net based and U-Net combined with Transformer models (see Table III), indicating that training neural networks with multiscale information can improve the accuracy of segmentation. By computing the segmentation ability scores [6] on layers of the U-Netmer, we find that the training U-Netmer with different patch sizes can improve the segmentation ability on different layers and further improve the final segmentation. Overall, experimental results have shown that U-Netmer trained with multi-scales can improve the accuracy on the test 7 datasets with different modalities.

### B. U-Netmer can be used to rank test images by difficulty without ground-truth

Ranking the test images by difficulty without ground-truth can provide useful information for end users to automatically identify the most challenging examples for human experts to review [32], [6]. Most segmentation models only output the segmentation result without such a "confidence" evaluation. This limits the use of segmentation algorithms in clinical practice when certain acceptance criteria are required [33] or when we need to prioritize users' time on inspection and auditing [32]. There is an unmet need to evaluate segmentation accuracies and even to reject failed segmentations in the real-world applications when the ground-truth is absent in reality [34]. Most studies estimate the pixel/voxel level uncertainty [32], [35], [36] to highlight the challenging regions

TABLE III
THE DICE SCORES OF THE U-NETMER AND OTHER STATE-OF-THE-ART MODELS ON THE 7 DATASETS.

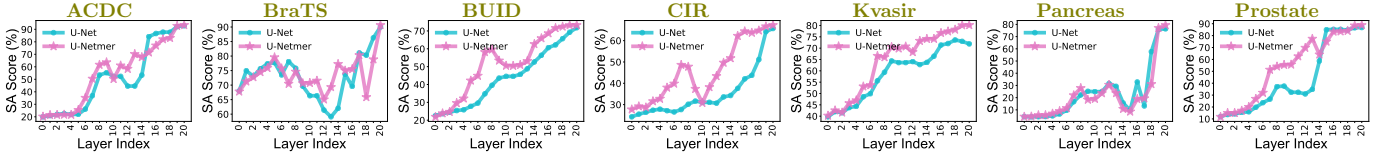| Models | ACDC | BUID | BraTS | CIR | Kvasir | Pancreas | Prostate |
|---|---|---|---|---|---|---|---|
| U-Net [5] | $92.84_{\pm 4.06}$ | $71.81_{\pm 29.57}$ | $90.01_{\pm 6.00}$ | $65.71_{\pm 23.24}$ | $71.88_{\pm 27.34}$ | $76.38_{\pm 8.53}$ | $86.83_{\pm 3.74}$ |
| Attention U-Net [2] | $92.81_{\pm 3.58}$ | $71.09_{\pm 29.20}$ | $90.17_{\pm 5.49}$ | $64.62_{\pm 24.89}$ | $69.18_{\pm 28.65}$ | $76.76_{\pm 8.03}$ | $87.63_{\pm 3.66}$ |
| U-Net++ [1] | $93.16_{\pm 3.68}$ | $73.58_{\pm 26.51}$ | $90.31_{\pm 5.36}$ | $64.57_{\pm 25.64}$ | $71.87_{\pm 25.87}$ | $75.61_{\pm 8.64}$ | $87.14_{\pm 3.58}$ |
| BiOnet [27] | $91.01_{\pm 4.93}$ | $67.91_{\pm 32.27}$ | $90.51_{\pm 4.91}$ | $65.77_{\pm 22.72}$ | $68.29_{\pm 29.08}$ | $77.59_{\pm 8.49}$ | $86.39_{\pm 4.71}$ |
| ConvUNeXt [28] | $85.28_{\pm 8.59}$ | $64.21_{\pm 27.79}$ | $88.23_{\pm 6.67}$ | $60.97_{\pm 25.43}$ | $49.30_{\pm 26.47}$ | $56.25_{\pm 12.60}$ | $80.46_{\pm 6.49}$ |
| ResUnet [29] | $91.94_{\pm 6.08}$ | $57.91_{\pm 29.84}$ | $90.29_{\pm 4.93}$ | $58.03_{\pm 28.39}$ | $67.66_{\pm 23.35}$ | $76.31_{\pm 8.78}$ | $87.01_{\pm 4.42}$ |
| UNext [30] | $86.32_{\pm 7.63}$ | $64.12_{\pm 30.62}$ | $89.23_{\pm 5.90}$ | $58.31_{\pm 28.61}$ | $57.26_{\pm 25.50}$ | $56.08_{\pm 11.73}$ | $81.86_{\pm 5.99}$ |
| UCTransNet [11] | $93.39_{\pm 4.19}$ | $72.39_{\pm 28.12}$ | $90.53_{\pm 5.77}$ | $64.79_{\pm 23.99}$ | $78.25_{\pm 24.05}$ | $75.31_{\pm 8.99}$ | $88.57_{\pm 3.31}$ |
| Trans U-Net [10] | $92.84_{\pm 4.11}$ | $74.00_{\pm 26.47}$ | $90.36_{\pm 5.73}$ | $66.48_{\pm 23.00}$ | $72.63_{\pm 26.56}$ | $75.52_{\pm 9.18}$ | $87.26_{\pm 3.97}$ |
| MedT [12] | $77.41_{\pm 12.03}$ | $48.67_{\pm 29.51}$ | $89.91_{\pm 6.29}$ | $59.73_{\pm 26.29}$ | $37.69_{\pm 29.27}$ | $43.27_{\pm 14.87}$ | $74.65_{\pm 5.93}$ |
| U-Netmer | $93.79_{\pm 3.30}$ | $74.82_{\pm 26.40}$ | $90.67_{\pm 5.25}$ | $67.29_{\pm 21.65}$ | $80.16_{\pm 20.94}$ | $\mathbf{79.42}_{\pm 7.59}$ | $89.14_{\pm 3.46}$ |
| Attention U-Netmer | $93.79_{\pm 4.01}$ | $\mathbf{75.97}_{\pm 24.95}$ | $90.72_{\pm 5.24}$ | $67.22_{\pm 21.61}$ | $79.86_{\pm 21.17}$ | $79.09_{\pm 7.63}$ | $\mathbf{89.43}_{\pm 3.13}$ |
| U-Netmer++ | $\mathbf{94.21}_{\pm 3.26}$ | $74.28_{\pm 26.24}$ | $\mathbf{90.95}_{\pm 5.17}$ | $\mathbf{68.15}_{\pm 20.96}$ | $\mathbf{80.76}_{\pm 20.40}$ | $78.71_{\pm 7.56}$ | $89.19_{\pm 3.13}$ |



Fig. 8. The segmentation ability (SA) scores of the U-Net and U-Netmer$_{s=<1|2|4>}$ on different layers (there are 18 layers on U-Net indexed from 1 to 19) measured by the ProtoSeg [6]. The input image is indexed as 0 while the segmentation output is indexed as 20.
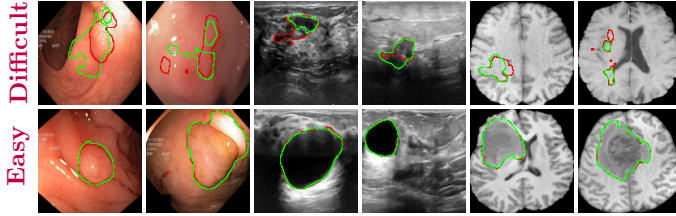


Fig. 9. Examples of difficult and easy examples for segmentation. The red contours are the segmentation results of U-Netmer with scale $s = 1$ and the green contours are the segmentation results with scale $s = 2$.
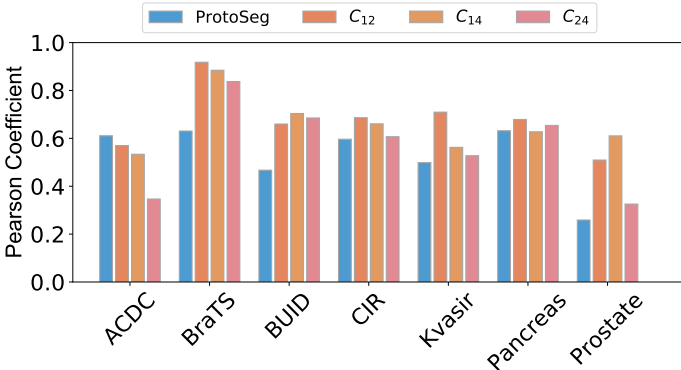


Fig. 10. Comparison of the Pearson coefficient between the segmentation accuracy and confidence scores computed by ProtoSeg [6] and $\mathcal{C}_{12}$, $\mathcal{C}_{14}$, ad $\mathcal{C}_{24}$ from U-Netmer$_{s=\langle 1|2|4\rangle}$.

within the target region or background on the single input image. These methods do not provide the image level confidence scores to automatically select a subset of challenging samples.

The U-Netmer can provide an estimation of the confidence score which can be used to rank the test image by difficulty without ground-truth. Given a test image, U-Netmer can make predictions with different scales. We use $B_{s=i}$ to denote the segmentation map computed with scale $s = i$. Fig. 9 shows testing examples of the segmentation maps from the scale $s = 1$ ($B_{s=1}$) and $s = 2$ ($B_{s=2}$) of the trained U-Netmer$_{s=\langle 1|2|4\rangle}$. One observation is that the segmentation maps of $B_{s=1}$ and $B_{s=2}$ are similar on easy samples while they are different on difficult samples. Part of the reason is that the target regions (e.g., lesions) on difficult samples have low contrast information compared to their background or have a small size, yielding to different segmentation results with different scales. Therefore, their discrepancy can be used to measure the difficulty of the test images or rank the test images by difficulty without ground-truth for segmentation.

To measure the difficulty of the test images without ground-truth, we compute the discrepancy between the estimation of $B_{s=i}$ and $B_{s=j}$ (where $i \neq j$) by: $\mathcal{C}_{ij} = \mathcal{D}(B_{s=i}, B_{s=j})$, where $\mathcal{D}$ is a distance metric to measure the difference between two estimated segmentation maps of $B_{s=i}$ and $B_{s=j}$. $\mathcal{C}$ has a different meanings given different distance metric $\mathcal{D}$. If the distance metric is defined as $\mathcal{D} = |B_{s=i} - B_{s=j}|$, the $\mathcal{C}$ indicates the uncertainty regions of the segmentation map [37]. In this paper, we use the Dice coefficient as the distance metric to measure the difficulty to segment each image, defined as: $\mathcal{C}_{ij} = 2|B_{s=i} \cap B_{s=j}|/|B_{s=i} + B_{s=j}|$ where $\mathcal{C}_{ij}$ measures the consistency between the segmentation maps between $B_{s=i}$ and $B_{s=j}$ obtained from U-Netmer with different scales $s = i$ and $s = j$. $\mathcal{C}_{ij}$ is considered as the confidence score to estimate the segmentation accuracy of the testing images without ground-
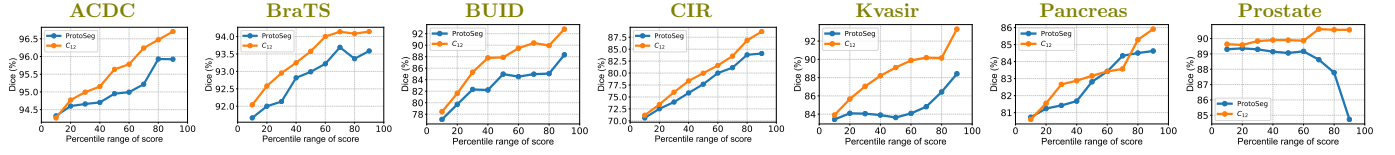
Fig. 11. The segmentation accuracy (y-axis) for the test images thresholded by the confidence score (x-axis) computed by ProtoSeg [6] and U-Netmer.
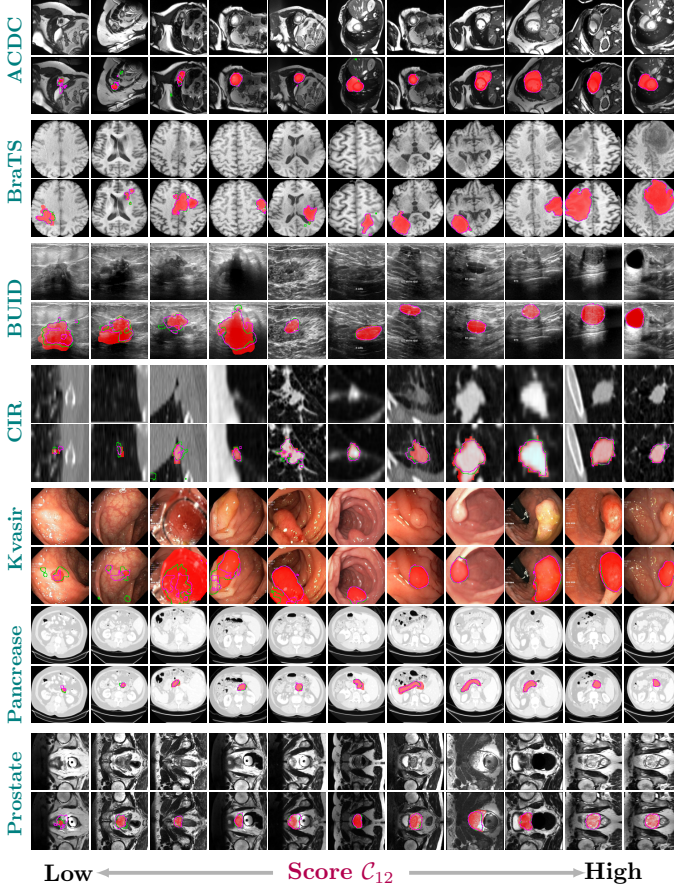


Fig. 12. Examples of the images (the first row on each data set) and their corresponding segmentation results (the second row on each dataset) ranked by the confidence score $\mathcal{C}_{12}$. The red masks denote the ground-truth and the green and pink contours denote the segmentation results of the $B_{s=1}$ and $B_{s=2}$, respectively.

truth for end users.

Fig. 10 shows the Pearson correlation between the Dice accuracy of the final segmentation and the confidence score $\mathcal{C}_{ij}$ obtained by U-Netmer$_{s=\langle 1|2|4\rangle}$. We also compare them with the baseline of the mean SA score computed by ProtoSeg [6] which is the mean segmentation ability score computed on the last two layers of the neural network for ranking the test images by difficulty. The results show that the Pearson correlation between $\mathcal{C}$ and the segmentation accuracy is higher than the one between ProtoSeg and the segmentation accuracy on 6 datasets except on ACDC dataset. In addition, $\mathcal{C}_{12}$ provides the highest correlation on BraTS, CIR, Kvasir, Pancreas while $\mathcal{C}_{14}$ provides the highest correlation on the BUID and Prostate

datasets.

To test whether the confidence score obtained by U-Netmer is discriminative between easy and difficult test samples, we plot the segmentation accuracy of test samples bucketed by the decile of confidence scores [32]. We first rank all testing samples based on the estimated confidence score and the testing samples within the $d\%$ percentile are included to compute the accuracy of the segmentation. This is similar to the coverage [38] which rejects the (100-$d$)% difficult samples for further attention. Fig. 11 shows the accuracy of segmentation with different percentile $d$. We show that examples at the highest percentiles on the rank often have high segmentation accuracy and the scores computed by U-Netmer provide higher accuracy than ProtoSeg [6] on 6 datasets except for Pancreas. The results also demonstrate that the confidence scores have a high correlation with the segmentation accuracies on the test images.

Fig. 12 visualizes test images and their corresponding segmentation results ranked by the confidence scores $\mathcal{C}$ on the 7 datasets. Images with low confidence scores tend to have small target regions, smooth boundaries and poor contrast between the target and background tissues. The accuracies of their segmentation results are usually low. Images with high confidence scores often have large target regions and clear boundaries, yielding more consistent segmentation maps obtained from different scales of U-Netmer. For test images without ground-truth, the confidence score can be used by end-users with a human-in-the-loop strategy: it can suggest the segmentation results of these images with low confidence scores to the human users for further review.

## V. CONCLUSION

In conclusion, we have presented the U-Netmer which is a combination of CNN-based neural network and Transformer for medical image segmentation. We have studied three variations of the U-Netmer where the backbones are from three typical segmentation neural networks: U-Netmer, Attention U-Netmer, and U-Netmer++. Experimental results on 7 datasets for medical image segmentation with different modalities have shown that the U-Netmer can provide competitive results compared to four baselines and six state-of-the-art models. The U-Nemter can also provide segmentation maps with different scales on test images. The discrepancy of these segmentations is linearly correlated to the segmentation accuracy, which can be considered as a confidence score to rank the test images by difficulty when the ground-truth is absent. This is important in real world applications, as it can highlight most difficult cases, or least accuracy cases, for users' further inspection and edit.

## REFERENCES

[1] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: Redesigning skip connections to exploit multiscale features in image segmentation," *IEEE transactions on medical imaging*, vol. 39, no. 6, pp. 1856–1867, 2019.

[2] J. Schlemper, O. Oktay, M. Schaap, M. Heinrich, B. Kainz, B. Glocker, and D. Rueckert, "Attention gated networks: Learning to leverage salient regions in medical images," *Medical image analysis*, vol. 53, pp. 197–207, 2019.

[3] M. Antonelli, A. Reinke, S. Bakas, K. Farahani, A. Kopp-Schneider, B. A. Landman, G. Litjens, B. Menze, O. Ronneberger, R. M. Summers *et al.*, "The medical segmentation decathlon," *Nature communications*, vol. 13, no. 1, p. 4128, 2022.

[4] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, and K. H. Maier-Hein, "nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation," *Nature methods*, vol. 18, no. 2, pp. 203–211, 2021.

[5] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[6] S. He, Y. Feng, P. E. Grant, and Y. Ou, "Segmentation ability map: Interpret deep features for medical image segmentation," *Medical Image Analysis*, vol. 84, p. 102726, 2023.

[7] R. Azad, E. K. Aghdam, A. Rauland, Y. Jia, A. H. Avval, A. Bozorgpour, S. Karimijafarbigloo, J. P. Cohen, E. Adeli, and D. Merhof, "Medical image segmentation review: The success of u-net," *arXiv preprint arXiv:2211.14830*, 2022.

[8] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.

[9] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2020.

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

[11] H. Wang, P. Cao, J. Wang, and O. R. Zaiane, "Uctransnet: rethinking the skip connections in u-net from a channel-wise perspective with transformer," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 36, no. 3, 2022, pp. 2441–2449.

[12] J. M. J. Valanarasu, P. Oza, I. Hacihaliloglu, and V. M. Patel, "Medical transformer: Gated axial-attention for medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 36–46.

[13] L. Beyer, P. Izmailov, A. Kolesnikov, M. Caron, S. Kornblith, X. Zhai, M. Minderer, M. Tschannen, I. Alabdulmohsin, and F. Pavetic, "Flexi-ViT: One model for all patch sizes," *arXiv preprint arXiv:2212.08013*, 2022.

[14] Q. Diao, Y. Jiang, B. Wen, J. Sun, and Z. Yuan, "Metaformer: A unified meta framework for fine-grained recognition," *arXiv preprint arXiv:2203.02751*, 2022.

[15] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 000–16 009.

[16] K. Donggyun, K. Jinwoo, C. Seongwoong, L. Chong, and H. Seunghoon, "Universal few-shot learning of dense prediction tasks with visual token matching," *https://openreview.net/forum?id=88nT0j5jAn*, 2023.

[17] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[18] O. Bernard, A. Lalande, C. Zotti, F. Cervenansky, X. Yang, P.-A. Heng, I. Cetin, K. Lekadir, O. Camara, M. A. G. Ballester *et al.*, "Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved?" *IEEE transactions on medical imaging*, vol. 37, no. 11, pp. 2514–2525, 2018.

[19] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, "The multimodal brain tumor image segmentation benchmark (BRATS)," *IEEE transactions on medical imaging*, vol. 34, no. 10, pp. 1993–2024, 2014.

[20] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, and C. Davatzikos, "Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features," *Scientific data*, vol. 4, no. 1, pp. 1–13, 2017.

[21] S. Bakas, M. Reyes, A. Jakab, S. Bauer, M. Rempfler, A. Crimi, R. T. Shinohara, C. Berger, S. M. Ha, M. Rozycki *et al.*, "Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge," *arXiv preprint arXiv:1811.02629*, 2018.

[22] W. Al-Dhabyani, M. Gomaa, H. Khaled, and A. Fahmy, "Dataset of breast ultrasound images," *Data in brief*, vol. 28, p. 104863, 2020.

[23] W. Choi, N. Dahiya, and S. Nadeem, "CIRDataset: A large-scale dataset for clinically-interpretable lung nodule radiomics and malignancy prediction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2022, pp. 13–22.

[24] D. Jha, P. H. Smedsrud, M. A. Riegler, P. Halvorsen, T. d. Lange, D. Johansen, and H. D. Johansen, "Kvasir-seg: A segmented polyp dataset," in *International Conference on Multimedia Modeling*. Springer, 2020, pp. 451–462.

[25] M. A. Attiyeh, J. Chakraborty, A. Doussot, L. Langdon-Embry, S. Mainarich, M. Gönen, V. P. Balachandran, M. I. D'Angelica, R. P. DeMatteo, W. R. Jarnagin *et al.*, "Survival prediction in pancreatic ductal adenocarcinoma by quantitative computed tomography image analysis," *Annals of surgical oncology*, vol. 25, no. 4, pp. 1034–1042, 2018.

[26] Q. Liu, Q. Dou, L. Yu, and P. A. Heng, "MS-Net: multi-site network for improving prostate segmentation with heterogeneous mri data," *IEEE transactions on medical imaging*, vol. 39, no. 9, pp. 2713–2724, 2020.

[27] T. Xiang, C. Zhang, D. Liu, Y. Song, H. Huang, and W. Cai, "BiO-Net: learning recurrent bi-directional connections for encoder-decoder architecture," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2020, pp. 74–84.

[28] Z. Han, M. Jian, and G.-G. Wang, "ConvUNeXt: An efficient convolution neural network for medical image segmentation," *Knowledge-Based Systems*, vol. 253, p. 109512, 2022.

[29] Z. Zhang, Q. Liu, and Y. Wang, "Road extraction by deep residual u-net," *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 5, pp. 749–753, 2018.

[30] J. M. J. Valanarasu and V. M. Patel, "UNeXt: Mlp-based rapid medical image segmentation network," *arXiv preprint arXiv:2203.04967*, 2022.

[31] S. He, P. E. Grant, and Y. Ou, "Global-local transformer for brain age estimation," *IEEE transactions on medical imaging*, vol. 41, no. 1, pp. 213–224, 2021.

[32] C. Agarwal, D. D'souza, and S. Hooker, "Estimating example difficulty using variance of gradients," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 368–10 378.

[33] S. Budd, E. C. Robinson, and B. Kainz, "A survey on active learning and human-in-the-loop deep learning for medical image analysis," *Medical Image Analysis*, vol. 71, p. 102062, 2021.

[34] V. V. Valindria, I. Lavdas, W. Bai, K. Kamnitsas, E. O. Aboagye, A. G. Rockall, D. Rueckert, and B. Glocker, "Reverse classification accuracy: predicting segmentation performance in the absence of ground truth," *IEEE transactions on medical imaging*, vol. 36, no. 8, pp. 1597–1606, 2017.

[35] Y. Shi, J. Zhang, T. Ling, J. Lu, Y. Zheng, Q. Yu, L. Qi, and Y. Gao, "Inconsistency-aware uncertainty estimation for semi-supervised medical image segmentation," *IEEE transactions on medical imaging*, vol. 41, no. 3, pp. 608–620, 2021.

[36] K. Wickstrøm, M. Kampffmeyer, and R. Jenssen, "Uncertainty and interpretability in convolutional neural networks for semantic segmentation of colorectal polyps," *Medical image analysis*, vol. 60, p. 101619, 2020.

[37] T. Nair, D. Precup, D. L. Arnold, and T. Arbel, "Exploring uncertainty measures in deep networks for multiple sclerosis lesion detection and segmentation," *Medical image analysis*, vol. 59, p. 101557, 2020.

[38] F. C. Ghesu, B. Georgescu, A. Mansoor, Y. Yoo, E. Gibson, R. Vishwanath, A. Balachandran, J. M. Balter, Y. Cao, R. Singh *et al.*, "Quantifying and leveraging predictive uncertainty for medical image assessment," *Medical Image Analysis*, vol. 68, p. 101855, 2021.