# MedSegDiff-V2: Diffusion based Medical Image Segmentation with Transformer

**Junde Wu**[1,6,7,8], **Wei Ji**[2], **Huazhu Fu**[3], **Min Xu**[*, 4,6] **Yueming Jin**[1], **Yanwu Xu**[*5],

[1]National University of Singapore, [2]University of Alberta, [3]A*STAR, [4]Carnegie Mellon University, [5]Singapore Eye Research Institute, [6]Mohamed bin Zayed University of Artificial Intelligence, [7]University of Oxford, [8]Kids with Tokens
jundewu@ieee.org, ywxu@ieee.org, xumin100@gmail.com

## Abstract

The Diffusion Probabilistic Model (DPM) has recently gained popularity in the field of computer vision, thanks to its image generation applications, such as Imagen, Latent Diffusion Models, and Stable Diffusion, which have demonstrated impressive capabilities and sparked much discussion within the community. Recent investigations have further unveiled the utility of DPM in the domain of medical image analysis, as underscored by the commendable performance exhibited by the medical image segmentation model across various tasks. Although these models were originally underpinned by a UNet architecture, there exists a potential avenue for enhancing their performance through the integration of vision transformer mechanisms. However, we discovered that simply combining these two models resulted in subpar performance. To effectively integrate these two cutting-edge techniques for the Medical image segmentation, we propose a novel Transformer-based Diffusion framework, called MedSegDiff-V2. We verify its effectiveness on 20 medical image segmentation tasks with different image modalities. Through comprehensive evaluation, our approach demonstrates superiority over prior state-of-the-art (SOTA) methodologies. Code is released at https://github.com/KidsWithTokens/MedSegDiff

## Introduction

Medical image segmentation is to divide a medical image into distinct regions of interest. It is a crucial step in many medical applications, such as diagnosis and image-guided surgery. In recent years, there has been a growing interest in automated segmentation methods, as they have the potential to improve the consistency and accuracy of results. With the advancement of deep learning techniques, several studies have successfully applied neural network-based models, including classical convolutional neural networks (CNNs) (Ji et al. 2021; Wu et al. 2022b) and the recently popular vision transformers (ViTs)(Chen et al. 2021; Wang et al. 2021b), to medical image segmentation tasks.

Very recently, the Diffusion Probabilistic Model (DPM)(Ho, Jain, and Abbeel 2020) has gained popularity as a powerful class of generative models, capable of generating high-quality and diverse images(Ramesh et al. 2022; Saharia

et al. 2022; Rombach et al. 2022). Inspired by its success, many researches have applied DPM in the field of medical image segmentation(Wu et al. 2022c; Wolleb et al. 2021; Kim, Oh, and Ye 2022; Guo et al. 2022; Rahman et al. 2023). Many of them reported new SOTA on several benchmarks by using the DPM. The remarkable performance of this model stems from its inherent stochastic sampling process(Wu et al. 2022c; Rahman et al. 2023). DPM has the capability to generate different segmentation predictions by running multiple times. The diversity among these samples directly captures the uncertainty associated with targets in medical images, where organs or lesions commonly have ambiguous boundaries. However, it is worth noting that all these methods rely on classical UNet backbones. In comparison to the increasingly popular vision transformers, classical UNet models compromise on segmentation quality, which can lead to the generation of divergent yet incorrect masks in ensemble, ultimately introducing noise that permanently hampers the performance.

A natural next step is to combine the transformer-based UNet, such as TransUNet(Chen et al. 2021), with DPM. However, we found that implementing it in a straightforward manner resulted in subpar performance. One issue is that the transformer-abstracted conditional feature is not compatible with the feature of the diffusion backbone. The transformer is able to learn deep semantic features from the raw image, whereas the diffusion backbone abstracts features from a corrupted and noisy mask, making feature fusion more challenging. Additionally, the dynamic and global nature of the transformer makes it more sensitive than CNNs (Naseer et al. 2021). Thus, the adaptive condition strategy used in previous diffusion-based methods(Wu et al. 2022c) will cause large variance in the transformer setting. This leads more ensemble and converge difficulties.

To overcome the aforementioned challenges, we have designed a novel Transformer-based Diffusion framework for the Medical image segmentation, called MedSegDiff-V2. The main idea is to employ two different conditioning techniques over the backbone with the raw image in the diffusion process. One is the Anchor Condition, which integrates the conditional segmentation features into the diffusion model encoder to reduce the diffusion variance. We design a novel $\mathcal{U}$ncertain Spatial Attention ($\mathcal{U}$-SA) mechanism for the integration, which relaxes the conditional segmentation feature with more uncer-
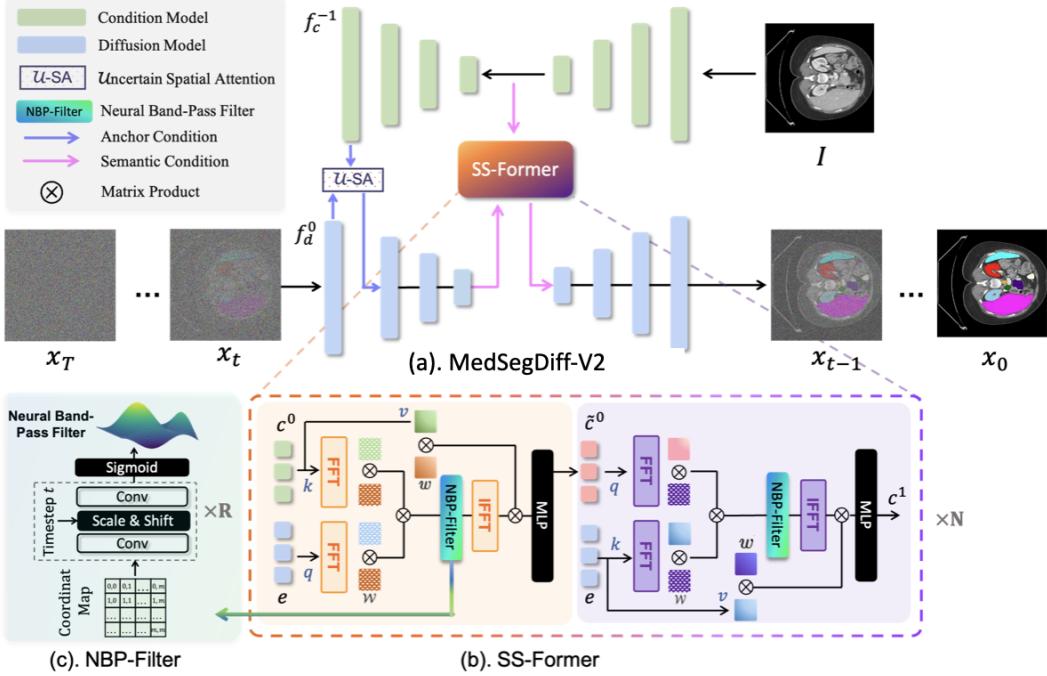
---

Figure 1: An illustration of MedSegDiff-V2, which starts from (a) an overview of the pipeline, and continues with zoomed-in diagrams of individual Models, including (b) SS-Former, and (c) NBP-Filter.

tainty, thus providing the diffusion process more flexibility to further calibrate the predictions. The other is the Semantic Condition that integrates the conditional embedding into the diffusion embedding. To effectively bridge the gap between these two embedding, we propose a novel transformer mechanism called the Spectrum-Space Transformer (SS-Former) for the embedding integration. SS-Former is a cross-attention chain in frequency domain, with a timestep-adaptive Neural Band-pass Filter (NBP-Filter) to align the noise and semantic features each time.

In brief, the contributions of this paper are:

- We are the first to integrate transformer into a diffusion-based model for general medical image segmentation.
- We propose an Anchor Condition with $\mathcal{U}$-SA to mitigate the diffusion variance.
- We propose Semantic Condition with SS-Former to model the segmentation noise and semantic feature interaction.
- We achieve SOTA performance on 20 organ segmentation tasks including 5 image modalities.

## Related Work
### Medical Segmentation with Transformers

Previous studies have highlighted the potential of transformer-based models to achieve SOTA results in medical image segmentation. One notable example is TransUNet(Chen et al. 2021), which combined the transformer with UNet as a bottleneck feature encoder. Since then, several works have proposed incorporating cutting-edge transformer techniques into

the backbone of medical image segmentation models, including Swin-UNet(Cao et al. 2022), Swin-UNetr(Tang et al. 2022), and DS-TransUNet(Lin et al. 2022). As recently UNet based diffusion-based segmentation models have recently emerged as achieving new SOTA in medical image segmentation, it is worthwhile to explore ways to integrate recognized transformer architectures into this powerful new backbone.

### Diffusion Model for Medical Segmentation

Diffusion models have recently demonstrated significant potential in various segmentation tasks, including medical images (Armato III et al. 2011; Caron et al. 2021; Cao et al. 2022; Chen, Ma, and Zheng 2019). In fact, these models leverage a stochastic sampling process to generate an implicit ensemble of segmentations, leading to enhanced segmentation performance (Zhai et al. 2022). However, without effective control of diversity, the ensemble often struggles to converge Therefore, it is crucial to improve the sample accuracy with each sampling iteration.

## Method
### Diffusion process of MedSegDiff-V2

We have designed our model based on the diffusion model mentioned in (Ho, Jain, and Abbeel 2020). Diffusion models are generative models that consist of two stages: a forward diffusion stage and a reverse diffusion stage. In the forward process, Gaussian noise is gradually added to the segmentation label $x_0$ through a series of steps $T$. In the reverse process, a neural network is trained to recover the original

data by reversing the noise addition process. This can be mathematically represented as follows:

$$p_\theta(x_{0:T-1}|x_T) = \prod_{t=1}^{T} p_\theta(x_{t-1}|x_t), \quad (1)$$

where $\theta$ represents the parameters of the reverse process. Starting from a Gaussian noise distribution, $p_\theta(x_T) = \mathcal{N}(x_T; 0, I_{n \times n})$, where $I$ is the raw image, the reverse process transforms the latent variable distribution $p_\theta(x_T)$ to the data distribution $p_\theta(x_0)$. To maintain symmetry with the forward process, the reverse process recovers the noisy image step by step, ultimately obtaining the final clear segmentation.

Following the standard implementation of DPM, we utilize an encoder-decoder network for the learning. To achieve segmentation, we condition the step estimation function $\epsilon$ on the prior information from the raw image. This conditioning can be expressed as:

$$\epsilon_\theta(x_t, I, t) = D(TransF(E_t^I, E_t^x), t), \quad (2)$$

Here, $TransF$ denotes the transformer based attention mechanism. $E_t^I$ represents the conditional feature embedding, which, in our case, corresponds to the embedding of the raw image. $E_t^x$ represents the feature embedding of the segmentation map for the current step. These two components are incorporated together by transformer and passed through a UNet decoder $D$ for reconstruction. The step index $t$ is integrated with the combined embedding and decoder features, and each step index is embedded using a shared learned look-up table, following the approach described in (Ho, Jain, and Abbeel 2020).

## Overall architecture

The overall flow of MedSegDiff-V2 is shown in 1. To introduce the process, consider a single step $t$ of the diffusion process. The noisy mask $x_t$ is first inputted to a UNet, called the Diffusion Model. Diffusion Model is conditioned by the segmentation features extracted from the raw images through another standard UNet, called the Condition Model. Two different conditioning manners are applied to the Diffusion Model: Anchor Condition and Semantic Condition. Following the flow of the input, the Anchor Condition is first imposed on the encoder of the Diffusion Model. It integrates the anchor segmentation features, which are the decoded segmentation features of the Condition Model, into the encoded features of the Diffusion Model. This allows the diffusion model to be initialized by a rough but static reference, which helps to reduce the diffusion variances. The Semantic Condition is then imposed on the embedding of the Diffusion Model, which integrates the semantic segmentation embedding of the Condition Model into the embedding of the Diffusion Model. This conditional integration is implemented by SS-Former, which bridges the gap between the noise and semantic embedding, and abstracts a stronger representation with the advantage of the global and dynamic nature of transformer.

MedSegDiff-V2 is trained using a standard noise prediction loss $\mathcal{L}_n$ following DPM(Ho, Jain, and Abbeel 2020) and an anchor loss $\mathcal{L}_{anc}$ supervising the Condition Model. $\mathcal{L}_{anc}$ is a combination of soft dice loss $\mathcal{L}_{dice}$ and cross-entropy

loss $\mathcal{L}_{ce}$. Specifically, the total loss function is represented as:

$$\mathcal{L}_{total}^t = \mathcal{L}_n^t + (t \equiv 0 \pmod{\alpha})(\mathcal{L}_{dice} + \beta \mathcal{L}_{ce}) \quad (3)$$

where $t \equiv 0 \pmod{\alpha}$ control the times of supervision over Condition Model through hyper-parameter $\alpha$, cross-entropy loss is weighted by hyper-parameter $\beta$, which are set as 5 and 10 respectively.

## Anchor Condition with $\mathcal{U}$-SA

Without the inductive bias of convolution layer, transformer blocks have stronger representation capability but are also more sensitive to the input variance when training data is limited(Naseer et al. 2021). Directly adding the transformer block to the Diffusion Model will cause the large variance on each time outputs. To overcome this negative effect, we adapt the structure of MedSegDiff(Wu et al. 2022c) and introduce the Anchor Condition operation to the Diffusion Model.

Anchor Condition provides a rough anchor feature from the Condition Model and integrates it into the Diffusion Model. This provides the Diffusion Model with a correct range for predictions while also allowing it to further refine the results. Specifically, we integrate the decoded segmentation features of the Condition Model into the encoder features of the Diffusion Model. We propose $\mathcal{U}$-SA mechanism for the feature fusion to represent the uncertainty nature of the given conditional features. Formally, consider we integrate the last conditional feature $f_c^{-1}$ into the first diffusion feature $f_d^0$. $\mathcal{U}$-SA can be expressed as:

$$f_{anc} = Max(f_c^{-1} * k_{Gauss}, f_c^{-1}), \quad (4)$$

$$f_d^{'0} = Sigmoid(f_{anc} * k_{Conv_{1 \times 1}}) \cdot f_d^0 + f_d^0, \quad (5)$$

where $*$ denotes slide-window kernel manipulation, $\cdot$ denotes general element-wise manipulation. In the equation, we first apply a learnable Gaussian kernel $k_G$ over $f_c^{-1}$ to smooth the activation, as $f_c^{-1}$ serves as an anchor but may not be completely accurate. We then select the maximum value between the smoothed map and the original feature map to preserve the most relevant information, resulting in a smoothed anchor feature $f_{anc}$. Then we integrate $f_{anc}$ into $f_d^0$ to obtain an enhanced feature $f_d^{'0}$. Specifically, we first apply a $1 \times 1$ convolution $k_{1 \times 1conv}$ to reduce the anchor feature channels to 1 and multiply it with $f_d^0$ after the Sigmoid activation, then add it to each channel of $f_d^0$, similar to the implementation of spatial attention(Woo et al. 2018).

## Semantic Condition with SS-Former

The Diffusion Model predicts redundant noise from a noisy mask input, leading to a domain gap between its embedding and the conditional segmentation semantic embedding. This divergence compromises performance when using matrix manipulations, such as in a stranded transformer. To address this challenge, we propose a novel Spectrum-Space Transformer (SS-Former). Our key idea is to learn the interaction of condition semantic feature and diffusion noise feature in the frequency domain. We use a filter, called the Neural Band-pass Filter (NBP-Filter) to align them to a unified range of

frequencies, i.e., spectrum. NBP-Filter learns to pass a specific spectrum while constraining the others. We learn this spectrum in a self-adaptive way to the diffusion time steps, as the noise-level (frequency range) is specific for each step.

A bird-eye view of SS-Former is shown in the 1 (b), which is composed of $N$ blocks that share the same architecture. We set $N = 4$ in the paper. Each block consists of two cross-attention-like modules. The first encodes the diffusion noise embedding into the condition semantic embedding, and the next symmetric module encodes the last semantic embedding into the diffusion noise embedding. This allows the model to learn the interaction between noise and semantic features and achieve a stronger representation. Formally, consider $c^0$ is the deepest feature embedding of Condition Model and $e$ is that of Diffusion Model. We first transfer $c^0$ and $e$ to the Fourier space, denoted as $F(c^0)$ and $F(e)$, respectively. Note that the feature maps are all patchlized and liner projected in accordance with the standard vision transformer method. Then we compute an affinity weight map over Fourier space taking $e$ as the $query$ and $c^0$ as the $key$, which can be represented by $\mathcal{M} = (F(c^0)\mathcal{W}^q)(F(e)\mathcal{W}^k)^T$, where $\mathcal{W}^q$ and $\mathcal{W}^k$ are the learnable $query$ and $key$ weights in Fourier space.

We then apply a NBP-Filter to align the representation of frequency. We note that each point in $\mathcal{M}$ now represents a particular frequency, and since we need to control a continuous range of frequencies, it is intuitive to establish a smooth projection from the feature map position to the frequency magnitude. To accomplish this, we use a neural network to learn a weight map from a coordinate map. By doing so, inductive bias of the network will facilitate the learning of a smooth projection, as similar inputs will naturally produce similar outputs(Sitzmann et al. 2020; Wu and Fu 2019). This idea is widely used in 3D vision tasks and is known as Neural Radiance Fields (NeRF)(Mildenhall et al. 2020). But different from the original NeRF, we further condition it with time-step information. Specifically, the network takes a coordinate map as input and produces an attention map to serve as the filter, both of which have the same size $\mathcal{M}$. We implement it using a simple stack of convolutional blocks with intermediate layer normalization. To condition the network with timestep information, we scale and shift the normalized features with the timestep embedding of the diffusion model. We use two MLP layers to project the current timestep embedding to two values representing the mean and variance, which are used for scaling and shifting, respectively. We stack a total of $R = 6$ such blocks and a Sigmoid function to produce the final filter. Finally, the filter is element-wise multiplied with the affinity map $\mathcal{M}$ in the pipeline. NBP-Filter is trained in an end-to-end manner with the whole pipeline.

The filtered affinity map $\mathcal{M}'$ is then transferred back to Euclidean space using inverse fast Fourier transform (IFFT) and applied to condition features in $value$: $f = F^{-1}(M')(c^0 w^v)$, where $W^v$ is the learnable value weights. We also use a MLP to further refine the attention result, obtaining the final feature $\tilde{c}^0$. The following attention module is symmetric to the first one, but using the combined feature $\tilde{c}^0$ as the query and noise embedding $e$ as the key and value, in order to transform the segmentation features to the noise domain. The transformed feature $c^1$ will serve as the condition

embedding for the next block.

# Experiments

## Dataset

We conduct the experiments on total five different medical image segmentation datasets. Two datasets are used to verify the general segmentation performance, which are public AMOS2022(Ji et al. 2022) dataset with sixteen anatomies and public BTCV(Fang and Yan 2020) dataset with twelve anatomies annotated for abdominal multi-organ segmentation. The other four public datasets REFUGE-2 (Fang et al. 2022), BraTs-2021 dataset (Baid et al. 2021), ISIC 2018 dataset(Milton 2019) and TNMIX dataset (Pedraza et al. 2015) are used to verify the model performance on multi-modal images, which are the optic-cup segmentation from fundus images, the brain tumor segmentation from MRI images, and the thyroid nodule segmentation from ultrasound images. More details about datasets are shown in the appendix.

## Implementation Details

All experiments were conducted using the PyTorch platform and trained/tested on 4 NVIDIA A100 GPUs. All images were uniformly resized to a resolution of $256 \times 256$ pixels. The networks were trained in an end-to-end manner using the AdamW(Loshchilov and Hutter 2017) optimizer with a batch size of 32. The initial learning rate was set to $1 \times 10^{-4}$. We employed 100 diffusion steps for the inference. We run the model 10 times for the ensemble, which is much fewer than the 25 times in MedSegDiff(Wu et al. 2022c). Then we use STAPLE algorithm(Warfield, Zou, and Wells 2004) to fuse the different samples. We evaluate the segmentation performance by Dice score, IoU, HD95 metrics.

## Main Results

**Comparing with SOTA on Abdominal Multi-organ Segmentation** To verify the general medical image segmentation performance, we compare MedSegDiff-V2 with SOTA segmentation methods on multi-organ segmentation dataset AMOS and BTCV. The quantitative results of Dice score are shown in 2 and 3 respectively. In the table, we compare with the segmentation methods which are widely-used and well-recognized in the community, including the CNN-based method nnUNet(Isensee et al. 2021), the transformer-based methods TransUNet(Chen et al. 2021), UNetr(Hatamizadeh et al. 2022), Swin-UNetr(Jiang et al. 2022) and the diffusion based method EnsDiff (Wolleb et al. 2021), SegDiff(Amit et al. 2021), MedSegDiff (Wu et al. 2022c). We also compare with a simple combination of diffusion and transformer model. We replace the UNet model in MedSegDiff to TransUNet and denoted it as 'MedSegDiff + TransUNet' in the table.

As seen in 2 and 3, advanced network architectures and sophisticated designs are crucial for achieving good performance. Considering the architecture, transformer-based models such as Swin-UNetr outperform the carefully designed CNN-based model, nnUNet. The diffusion-based model MedSegDiff again outperforms the transformer-based models on
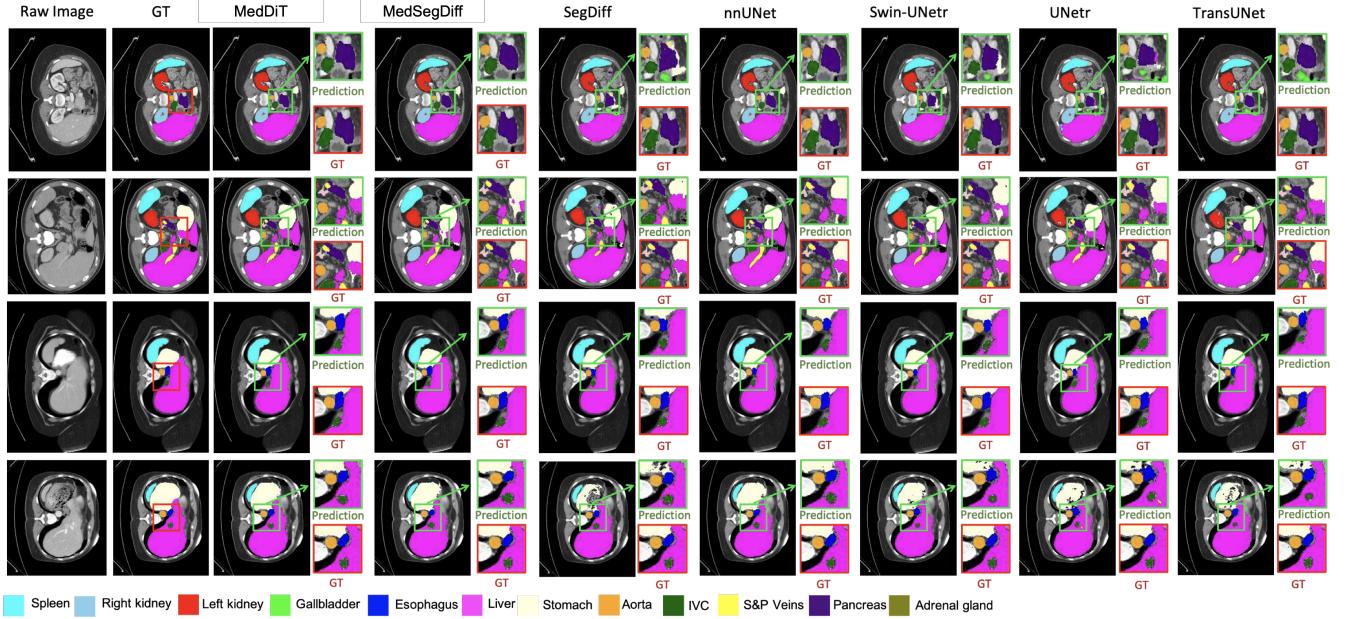
Figure 2: The visual comparison with SOTA segmentation models on BTCV.

most of the organs. However, network architecture alone is not the determining factor for performance. For example, the well-designed CNN-based model nnUNet considerably outperforms the transformer-based model TransUNet and UNetr in the table. This is also true for diffusion-based models. We can see that a straightforward adoption of the diffusion model for medical image segmentation, i.e., EnsDiff and SegDiff, perform worse than UNetr and Swin-UNetr. A simple combination of transformer and diffusion model, i.e., MedSegDiff + TransUNet, obtains even worse performance than the standard MedSegDiff. By introducing Anchor Condition and SS-Former in the diffusion + transformer model, MedSegDiff-V2 overcomes these challenges and shows superior performance. We also present a qualitative comparison in 2. It can be observed that MedSegDiff-V2 predicts segmentation maps with more precise details, even in low-contrast or ambiguous areas.

**Comparing with SOTA on Multi-modality Images** We also compare MedSegDiff-V2 to SOTA segmentation methods proposed for three specific tasks with different image modalities. The results are presented in 1. In the table, ResUnet(Yu et al. 2019) and BEAL(Wang et al. 2019) are proposed for optic cup segmentation, TransBTS(Wang et al. 2021b) and SwinBTS(Wang et al. 2021b) are proposed for brain tumor segmentation, MTSeg(Gong et al. 2021) and UltraUNet(Chu, Zheng, and Zhou 2021) are proposed for thyroid nodule segmentation, and FAT-Net(Wu et al. 2022a) and BAT(Wang et al. 2021a) are proposed for skin lesion segmentation.

From the table, we can see that MedSegDiff-V2 surpasses all other methods in five different tasks, highlighting its remarkable generalization capability across various medical segmentation tasks and image modalities. In comparison to

the UNet-based MedSegDiff, MedSegDiff-V2 exhibits improvements of 2.0% on Optic-Cup, 1.9% on Brain-Tumor, and 3.9% on Thyroid Nodule in terms of the Dice score, underscoring the effectiveness of its transformer-based backbone. Furthermore, when compared to MedSegDiff plus TransUNet, MedSegDiff-V2 outperforms it by an even larger margin, clearly demonstrating the efficacy of the proposed $\mathcal{U}$-SA and SS-Former in enhancing performance.

**Ablation Study** We conducted a comprehensive ablation study to verify the effectiveness of the proposed modules. The results are shown in 4, where Anc.Cond. and Sem.Cond. denote Anchor Condition and Semantic Condition, respectively. As shown in the table, Anc.Cond. significantly improves the vanilla diffusion model, with the proposed $\mathcal{U}$-SA outperforming the previous Spatial Attention on all datasets. In Sem.Cond., using SS-Former alone provides only marginal improvement, but combining it with the NBP-Filter results in a significant improvement, demonstrating the effectiveness of the proposed SS-Former design.

## Analysis and Discussion

**Implicit Ensemble Effect** As confirmed by numerous previous studies (Wolleb et al. 2021; Wu et al. 2022c; Amit et al. 2021), the implicit ensemble of multiple sampling runs plays a crucial role in diffusion-based methods. In diffusion model context, implicit ensemble refers to combining predictions from multiple samplings of a single diffusion model, rather than fusing predictions from different models.

In this study, we evaluate the ensemble performance of various diffusion-based medical segmentation models, as shown in 3. The evaluation is based on the average Dice Score calculated on the AMOS dataset. Each configuration is run 20 times, and the average Dice Score is used as the performance

Table 1: The comparison of MedSegDiff-V2 with SOTA segmentation methods on different image modalities. The grey background denotes the methods are proposed for that/these particular tasks.

| | | REFUGE2-Disc | | REFUGE2-Cup | | BraTs | | | TNMIX | | ISIC | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Dice | IoU | Dice | IoU | Dice | IoU | HD95 | Dice | IoU | Dice | IoU |
| Optic Disc/Cup | ResUNet | 92.9 | 85.5 | 80.1 | 72.3 | 78.4 | 71.3 | 18.71 | 78.3 | 70.7 | 87.1 | 78.2 |
| | BEAL | 93.7 | 86.1 | 83.5 | 74.1 | 78.8 | 71.7 | 18.53 | 78.6 | 71.6 | 86.6 | 78.0 |
| Brain Tumor | TransBTS | 94.1 | 87.2 | 85.4 | 75.7 | 87.6 | 78.44 | 12.44 | 83.8 | 75.5 | 88.1 | 80.6 |
| | SwinBTS | 95.2 | 87.7 | 85.7 | 75.9 | 88.7 | 81.2 | 10.03 | 84.5 | 76.1 | 89.8 | 82.4 |
| Thyroid Nodule | MTSeg | 90.3 | 83.6 | 82.3 | 73.1 | 82.2 | 74.5 | 15.74 | 82.3 | 75.2 | 87.5 | 79.7 |
| | UltraUNet | 91.5 | 82.8 | 83.1 | 73.78 | 84.5 | 76.3 | 14.03 | 84.5 | 76.2 | 89.0 | 81.8 |
| Skin Lesion | FAT-Net | 91.8 | 84.8 | 80.9 | 71.5 | 79.2 | 72.8 | 17.35 | 80.8 | 73.4 | 90.7 | 83.9 |
| | BAT | 92.3 | 85.8 | 82.0 | 73.2 | 79.6 | 73.5 | 15.49 | 81.7 | 74.2 | 91.2 | 84.3 |
| General Med Seg | nnUNet | 94.7 | 87.3 | 84.9 | 75.1 | 88.5 | 80.6 | 11.20 | 84.2 | 76.2 | 90.8 | 83.6 |
| | TransUNet | 95.0 | 87.7 | 85.6 | 75.9 | 86.6 | 79.0 | 13.74 | 83.5 | 75.1 | 89.4 | 82.2 |
| | UNetr | 94.9 | 87.5 | 83.2 | 73.3 | 87.3 | 80.6 | 12.81 | 81.7 | 73.5 | 89.7 | 82.8 |
| | Swin-UNetr | 95.3 | 87.9 | 84.3 | 74.5 | 88.4 | 81.8 | 11.36 | 83.5 | 74.8 | 90.2 | 83.1 |
| Diffusion Based | EnsemDiff | 94.3 | 87.8 | 84.2 | 74.4 | 88.7 | 80.9 | 10.85 | 83.9 | 75.3 | 88.2 | 80.7 |
| | SegDiff | 92.6 | 85.2 | 82.5 | 71.9 | 85.7 | 77.0 | 14.31 | 81.9 | 74.8 | 87.3 | 79.4 |
| | MedsegDiff | 95.1 | 87.6 | 85.9 | 76.2 | 88.9 | 81.2 | 10.41 | 84.8 | 76.4 | 91.3 | 84.1 |
| | MedsegDiff+TransUNet | 91.8 | 84.5 | 82.1 | 72.6 | 86.1 | 78.0 | 13.88 | 79.2 | 71.4 | 84.6 | 75.5 |
| Proposed | MedSegDiff-V2 | **96.7** | **88.9** | **87.9** | **80.3** | **90.8** | **83.4** | **7.53** | **88.7** | **81.5** | **93.2** | **85.3** |

Table 2: The comparison of MedSegDiff-V2 with SOTA segmentation methods over AMOS dataset evaluated by Dice Score. Best results are denoted as **bold**.

| Methods | Spleen | R.Kid | L.Kid | Gall. | Eso. | Liver | Stom. | Aorta | IVC | Panc. | RAG | LAG | Duo. | Blad. | Pros. | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TransUNet | 0.881 | 0.928 | 0.919 | 0.813 | 0.740 | 0.973 | 0.832 | 0.919 | 0.841 | 0.713 | 0.638 | 0.565 | 0.685 | 0.748 | 0.692 | 0.792 |
| UNetr | 0.926 | 0.936 | 0.918 | 0.785 | 0.702 | 0.969 | 0.788 | 0.893 | 0.828 | 0.732 | 0.717 | 0.554 | 0.658 | 0.683 | 0.722 | 0.762 |
| Swin-UNetr | 0.959 | 0.960 | 0.949 | 0.894 | 0.827 | **0.979** | 0.899 | 0.944 | 0.899 | 0.828 | 0.791 | 0.745 | 0.817 | **0.875** | 0.841 | 0.880 |
| nnUNet | 0.965 | 0.959 | 0.951 | 0.889 | 0.820 | 0.980 | 0.890 | 0.948 | 0.901 | 0.821 | 0.785 | 0.739 | 0.806 | 0.869 | 0.839 | 0.878 |
| EnsDiff | 0.905 | 0.918 | 0.904 | 0.732 | 0.723 | 0.947 | 0.838 | 0.915 | 0.838 | 0.704 | 0.677 | 0.618 | 0.715 | 0.673 | 0.680 | 0.786 |
| SegDiff | 0.885 | 0.872 | 0.891 | 0.703 | 0.654 | 0.852 | 0.702 | 0.874 | 0.819 | 0.715 | 0.654 | 0.632 | 0.697 | 0.652 | 0.695 | 0.753 |
| MedSegDiff | 0.963 | 0.965 | 0.953 | 0.917 | 0.846 | 0.971 | 0.906 | 0.952 | 0.918 | 0.854 | 0.803 | 0.751 | 0.819 | 0.868 | 0.855 | 0.889 |
| MedSegDiff + TransUNet | 0.941 | 0.932 | 0.921 | 0.934 | 0.813 | 0.946 | 0.867 | 0.921 | 0.880 | 0.821 | 0.793 | 0.528 | 0.788 | 0.813 | 0.837 | 0.849 |
| Anchor | 0.872 | 0.901 | 0.892 | 0.784 | 0.802 | 0.910 | 0.835 | 0.908 | 0.810 | 0.735 | 0.682 | 0.651 | 0.583 | 0.631 | 0.728 | 0.781 |
| MedSegDiff-V2 | **0.971** | **0.969** | **0.964** | **0.932** | **0.864** | 0.976 | **0.934** | **0.968** | **0.925** | **0.871** | **0.815** | **0.762** | **0.827** | 0.873 | **0.871** | **0.901** |

metric. In the figure, we denote "MedSegDiff+TransUNet" setting as "MSD-Trans". Our findings indicate a common trend, where the model performance improves rapidly in the initial 50 ensembles and then stabilizes. Typically, the best performance is achieved after approximately 50 ensembles.

When comparing MedSegDiff-V2 variant with other diffusion methods, we observe that it requires fewer ensembles to converge. It starts with a significantly better performance, surpassing MedSegDiff by 5%, and consistently maintains a lead of over 2% throughout. This highlights the efficiency of MedSegDiff-V2, as it achieves satisfactory results even with fewer ensemble iterations. Moreover, it suggests that a superior starting point and more stable predictions can lead to a higher performance ceiling. This aligns with our assumption that low-quality samples can consistently degrade the model's performance by introducing noise. This again demonstrates the importance of introducing $\mathcal{U}$-SA for divergence control and utilizing SS-Former to attain a better starting point.

**Analysis of Uncertainty** In 5, we compare the sample diversity on REFUGE2-Cup dataset. We compare the previous DPM-based methods, backbone with individual proposed modules, and final MedSegDiff-V2 together. We evaluated the variance among the samples using the Generalized En-

ergy Distance (GED) and Confidence Interval (CI). GED is a commonly used metric to measures the agreement between predictions and the ground truth distribution of segmentation by comparing their distributions(Kohl et al. 2018). A lower energy value indicates better agreement.

From the table, we can see that the proposed $\mathcal{U}$-SA achieves lower CI and higher GED compared to previous methods, indicating a larger sample diversity. However, it is also observed that the proposed model reaches a higher or comparable performance, suggesting that its generated samples mostly fall within the uncertainty region of the targets. When using SS-Former alone, without $\mathcal{U}$-SA, the model achieves the best agreement with the highest CI and lowest GED. Although SS-Former gets a fine performance with largest confidence, it fails to fully use the diversity ensemble capability of the diffusion model. By combining $\mathcal{U}$-SA and SS-Former as MedSegDiff-V2, the performance is significantly improved with still high confidence. It suggests that SS-Former helps mitigate the noise generated in $\mathcal{U}$-SA, while $\mathcal{U}$-SA provides more diversity to the model, resulting in mutual improvement.

**Model Efficiency and Complexity** In 5, we also present a comparison of model complexity and Gflops with other

Table 3: The comparison of MedSegDiff-V2 with SOTA segmentation methods over BTCV dataset evaluated by Dice Score. Best results are denoted as **bold**.

| Model | Spleen | R.Kid | L.Kid | Gall. | Eso. | Liver | Stom. | Aorta | IVC | Veins | Panc. | AG | Ave |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TransUNet | 0.952 | 0.927 | 0.929 | 0.662 | 0.757 | 0.969 | 0.889 | 0.920 | 0.833 | 0.791 | 0.775 | 0.637 | 0.838 |
| UNetr | 0.968 | 0.924 | 0.941 | 0.750 | 0.766 | 0.971 | 0.913 | 0.890 | 0.847 | 0.788 | 0.767 | 0.741 | 0.856 |
| Swin-UNetr | 0.971 | 0.936 | 0.943 | 0.794 | 0.773 | 0.975 | 0.921 | 0.892 | 0.853 | 0.812 | 0.794 | 0.765 | 0.869 |
| nnUNet | 0.942 | 0.894 | 0.910 | 0.704 | 0.723 | 0.948 | 0.824 | 0.877 | 0.782 | 0.720 | 0.680 | 0.616 | 0.802 |
| EnsDiff | 0.938 | 0.931 | 0.924 | 0.772 | 0.771 | 0.967 | 0.910 | 0.869 | 0.851 | 0.802 | 0.771 | 0.745 | 0.854 |
| SegDiff | 0.954 | 0.932 | 0.926 | 0.738 | 0.763 | 0.953 | 0.927 | 0.846 | 0.833 | 0.796 | 0.782 | 0.723 | 0.847 |
| MedSegDiff | 0.973 | 0.930 | 0.955 | 0.812 | 0.815 | 0.973 | 0.924 | 0.907 | 0.868 | 0.825 | 0.788 | 0.779 | 0.879 |
| MedSegDiff +TransUNet | 0.912 | 0.876 | 0.846 | 0.645 | 0.718 | 0.947 | 0.824 | 0.876 | 0.715 | 0.775 | 0.672 | 0.618 | 0.785 |
| Anchor | 0.928 | 0.882 | 0.873 | 0.652 | 0.750 | 0.951 | 0.829 | 0.855 | 0.731 | 0.714 | 0.683 | 0.602 | 0.787 |
| MedSegDiff-V2 | **0.978** | **0.941** | **0.963** | **0.848** | **0.818** | **0.985** | **0.940** | **0.928** | **0.869** | **0.823** | **0.831** | **0.817** | **0.895** |

Table 4: An ablation study on Anchor Conditioning and SS-Former. SA denotes Spatial Attention.

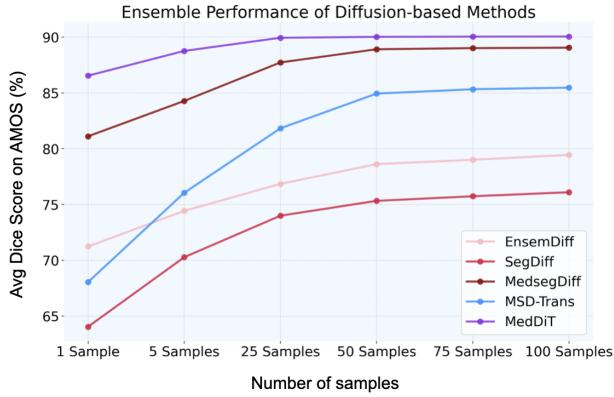| Anc.Cond. | | Sem.Cond. | | AMOS | BTCV | OpticCup | BrainTumor | ThyroidNodule |
|---|---|---|---|---|---|---|---|---|
| SA | $\mathcal{U}$-SA | SS-Former (w/o Filter) | NBP-Filter | Ave-Dice (%) | Ave-Dice (%) | Dice (%) | Dice (%) | Dice (%) |
| | | | | 78.6 | 85.4 | 84.6 | 88.2 | 84.1 |
| ✓ | | | | 83.5 | 85.8 | 85.2 | 88.7 | 84.6 |
| | ✓ | | | 86.7 | 86.6 | 85.7 | 89.4 | 86.5 |
| | ✓ | ✓ | | 87.8 | 87.1 | 86.5 | 89.8 | 86.8 |
| | ✓ | ✓ | ✓ | **90.1** | **89.5** | **87.9** | **90.8** | **88.7** |



Figure 3: The comparison of ensemble effect of DPM-based methods. We show their performance of average Dice Score on AMOS with increasing sampling times.

stability in fewer steps. In comparison, MedSegDiff-V2 consumes only half the Gflops of MedSegDiff while outperforming it in various segmentation tasks, as demonstrated above. This underscores the efficiency of MedSegDiff-V2 and its potential in real-world application.

Table 5: Comparison of model parameters, Gflops, and generated samples uncertainty

| Model | Params (M) | Gflops | CI | GED | Dice |
|---|---|---|---|---|---|
| EnsemDiff | **23** | 2203 | 76.3 | 28.9 | 84.2 |
| SegDiff | **23** | 2399 | 75.4 | 26.4 | 82.5 |
| MedSegDiff | 25 | 1770 | 77.5 | 27.9 | 85.9 |
| MSD-Trans | 118 | 2581 | 75.8 | 28.7 | 82.1 |
| bone + $\mathcal{U}$-SA | - | - | 73.2 | 34.6 | 85.7 |
| bone + SS-Former | - | - | **84.2** | **21.7** | 86.1 |
| MedSegDiff-V2 | 46 | **983** | 82.6 | 23.5 | **87.9** |

diffusion-based segmentation methods. The reported Gflops is the processing speed for a single $256 \times 256$ image until stability is reached in the implicit ensemble. We consider a variance of performance less than 0.1% across the last ten ensembles as an indicator of convergence. This metric is significant for the practical application of diffusion-based segmentation models, as users commonly run the diffusion model iteratively to obtain a stable result.

We can see from the table that, unlike traditional deep learning models, the amount of parameters in diffusion-based models is not directly correlated with Gflops, due to the presence of the implicit ensemble. For instance, even though MedSegDiff-V2 incorporates transformer blocks and occupies more parameters, it requires fewer Gflops as it achieves

## Conclusion

In this paper, we enhance the diffusion-based medical image segmentation framework by incorporating the transformer mechanism into the original UNet backbone, called MedSegDiff-V2. We propose a novel SS-Former architecture to learn the interaction between noise and semantic features. The comparative experiments show our model outperformed previous SOTA methods on 20 different medical image segmentation tasks with various image modalities. As the first transformer-based diffusion model for medical image segmentation, we believe MedSegDiff-V2 will serve as a benchmark for future research.

# References

Amit, T.; Nachmani, E.; Shaharbany, T.; and Wolf, L. 2021. Segdiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*.

Armato III, S. G.; McLennan, G.; Bidaut, L.; McNitt-Gray, M. F.; Meyer, C. R.; Reeves, A. P.; Zhao, B.; Aberle, D. R.; Henschke, C. I.; Hoffman, E. A.; et al. 2011. The lung image database consortium (LIDC) and image database resource initiative (IDRI): a completed reference database of lung nodules on CT scans. *Medical physics*, 38(2): 915–931.

Baid, U.; Ghodasara, S.; Mohan, S.; Bilello, M.; Calabrese, E.; Colak, E.; Farahani, K.; Kalpathy-Cramer, J.; Kitamura, F. C.; Pati, S.; et al. 2021. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*.

Cao, H.; Wang, Y.; Chen, J.; Jiang, D.; Zhang, X.; Tian, Q.; and Wang, M. 2022. Swin-unet: Unet-like pure transformer for medical image segmentation. In *European conference on computer vision*, 205–218. Springer.

Caron, M.; Touvron, H.; Misra, I.; Jégou, H.; Mairal, J.; Bojanowski, P.; and Joulin, A. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 9650–9660.

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

Chen, S.; Ma, K.; and Zheng, Y. 2019. Med3d: Transfer learning for 3d medical image analysis. *arXiv preprint arXiv:1904.00625*.

Chu, C.; Zheng, J.; and Zhou, Y. 2021. Ultrasonic thyroid nodule detection method based on U-Net network. *Computer Methods and Programs in Biomedicine*, 199: 105906.

Fang, H.; Li, F.; Fu, H.; Sun, X.; Cao, X.; Son, J.; Yu, S.; Zhang, M.; Yuan, C.; Bian, C.; et al. 2022. REFUGE2 Challenge: Treasure for Multi-Domain Learning in Glaucoma Assessment. *arXiv preprint arXiv:2202.08994*.

Fang, X.; and Yan, P. 2020. Multi-organ segmentation over partially labeled datasets with multi-scale feature abstraction. *IEEE Transactions on Medical Imaging*, 39(11): 3619–3629.

Gong, H.; Chen, G.; Wang, R.; Xie, X.; Mao, M.; Yu, Y.; Chen, F.; and Li, G. 2021. Multi-task learning for thyroid nodule segmentation with thyroid region prior. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, 257–261. IEEE.

Guo, X.; Yang, Y.; Ye, C.; Lu, S.; Xiang, Y.; and Ma, T. 2022. Accelerating Diffusion Models via Pre-segmentation Diffusion Sampling for Medical Image Segmentation. *arXiv preprint arXiv:2210.17408*.

Hatamizadeh, A.; Tang, Y.; Nath, V.; Yang, D.; Myronenko, A.; Landman, B.; Roth, H. R.; and Xu, D. 2022. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 574–584.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840–6851.

Isensee, F.; Jaeger, P. F.; Kohl, S. A.; Petersen, J.; and Maier-Hein, K. H. 2021. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2): 203–211.

Ji, W.; Yu, S.; Wu, J.; Ma, K.; Bian, C.; Bi, Q.; Li, J.; Liu, H.; Cheng, L.; and Zheng, Y. 2021. Learning calibrated medical image segmentation via multi-rater agreement modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12341–12351.

Ji, Y.; Bai, H.; Yang, J.; Ge, C.; Zhu, Y.; Zhang, R.; Li, Z.; Zhang, L.; Ma, W.; Wan, X.; et al. 2022. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *arXiv preprint arXiv:2206.08023*.

Jiang, Y.; Zhang, Y.; Lin, X.; Dong, J.; Cheng, T.; and Liang, J. 2022. SwinBTS: A method for 3D multimodal brain tumor segmentation using swin transformer. *Brain sciences*, 12(6): 797.

Kim, B.; Oh, Y.; and Ye, J. C. 2022. Diffusion adversarial representation learning for self-supervised vessel segmentation. *arXiv preprint arXiv:2209.14566*.

Kohl, S.; Romera-Paredes, B.; Meyer, C.; De Fauw, J.; Ledsam, J. R.; Maier-Hein, K.; Eslami, S.; Jimenez Rezende, D.; and Ronneberger, O. 2018. A probabilistic u-net for segmentation of ambiguous images. *Advances in neural information processing systems*, 31.

Lin, A.; Chen, B.; Xu, J.; Zhang, Z.; Lu, G.; and Zhang, D. 2022. Ds-transunet: Dual swin transformer u-net for medical image segmentation. *IEEE Transactions on Instrumentation and Measurement*, 71: 1–15.

Loshchilov, I.; and Hutter, F. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing scenes as neural radiance fields for view synthesis. In *The European Conference on Computer Vision (ECCV)*.

Milton, M. A. A. 2019. Automated skin lesion classification using ensemble of deep neural networks in isic 2018: Skin lesion analysis towards melanoma detection challenge. *arXiv preprint arXiv:1901.10802*.

Naseer, M. M.; Ranasinghe, K.; Khan, S. H.; Hayat, M.; Shahbaz Khan, F.; and Yang, M.-H. 2021. Intriguing properties of vision transformers. *Advances in Neural Information Processing Systems*, 34: 23296–23308.

Pedraza, L.; Vargas, C.; Narváez, F.; Durán, O.; Muñoz, E.; and Romero, E. 2015. An open access thyroid ultrasound image database. In *10th International symposium on medical information processing and analysis*, volume 9287, 188–193. SPIE.

Rahman, A.; Valanarasu, J. M. J.; Hacihaliloglu, I.; and Patel, V. M. 2023. Ambiguous medical image segmentation using diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11536–11546.

Ramesh, A.; Dhariwal, P.; Nichol, A.; Chu, C.; and Chen, M. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10684–10695.

Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E.; Ghasemipour, S. K. S.; Ayan, B. K.; Mahdavi, S. S.; Lopes, R. G.; et al. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *arXiv preprint arXiv:2205.11487*.

Sitzmann, V.; Martel, J.; Bergman, A.; Lindell, D.; and Wetzstein, G. 2020. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33: 7462–7473.

Tang, Y.; Yang, D.; Li, W.; Roth, H. R.; Landman, B.; Xu, D.; Nath, V.; and Hatamizadeh, A. 2022. Self-supervised pretraining of swin transformers for 3d medical image analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20730–20740.

Wang, J.; Wei, L.; Wang, L.; Zhou, Q.; Zhu, L.; and Qin, J. 2021a. Boundary-aware transformers for skin lesion segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24*, 206–216. Springer.

Wang, S.; Yu, L.; Li, K.; Yang, X.; Fu, C.-W.; and Heng, P.-A. 2019. Boundary and entropy-driven adversarial learning for fundus image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 102–110. Springer.

Wang, W.; Chen, C.; Ding, M.; Yu, H.; Zha, S.; and Li, J. 2021b. Transbts: Multimodal brain tumor segmentation using transformer. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 109–119. Springer.

Warfield, S. K.; Zou, K. H.; and Wells, W. M. 2004. Simultaneous truth and performance level estimation (STAPLE): an algorithm for the validation of image segmentation. *IEEE transactions on medical imaging*, 23(7): 903–921.

Wolleb, J.; Sandkühler, R.; Bieder, F.; Valmaggia, P.; and Cattin, P. C. 2021. Diffusion Models for Implicit Image Segmentation Ensembles. *arXiv preprint arXiv:2112.03145*.

Woo, S.; Park, J.; Lee, J.-Y.; and Kweon, I. S. 2018. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, 3–19.

Wu, H.; Chen, S.; Chen, G.; Wang, W.; Lei, B.; and Wen, Z. 2022a. FAT-Net: Feature adaptive transformers for automated skin lesion segmentation. *Medical image analysis*, 76: 102327.

Wu, J.; Fang, H.; Shang, F.; Yang, D.; Wang, Z.; Gao, J.; Yang, Y.; and Xu, Y. 2022b. SeATrans: Learning Segmentation-Assisted Diagnosis Model via Transformer. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part II*, 677–687. Springer.

Wu, J.; Fang, H.; Zhang, Y.; Yang, Y.; and Xu, Y. 2022c. MedSegDiff: Medical Image Segmentation with Diffusion Probabilistic Model. *arXiv preprint arXiv:2211.00611*.

Wu, J.; and Fu, R. 2019. Universal, transferable and targeted adversarial attacks. *arXiv preprint arXiv:1908.11332*.

Yu, S.; Xiao, D.; Frost, S.; and Kanagasingam, Y. 2019. Robust optic disc and cup segmentation with deep learning for glaucoma detection. *Computerized Medical Imaging and Graphics*, 74: 61–71.

Zhai, X.; Kolesnikov, A.; Houlsby, N.; and Beyer, L. 2022. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12104–12113.