




# Improving 3D Medical Image Segmentation at Boundary Regions using Local Self-attention and Global Volume Mixing

Daniya Najiha Abdul Kareem , Mustansar Fiaz , Noa Novershtern, Jacob Hanna , and Hisham Cholakkal

**Abstract**—Volumetric medical image segmentation is a fundamental problem in medical image analysis where the objective is to accurately classify a given 3D volumetric medical image with voxel-level precision. In this work, we propose a novel hierarchical encoder-decoder-based framework that strives to explicitly capture the local and global dependencies for volumetric 3D medical image segmentation. The proposed framework exploits local volume-based self-attention to encode the local dependencies at high resolution and introduces a novel volumetric MLP-mixer to capture the global dependencies at low-resolution feature representations, respectively. The proposed volumetric MLP-mixer learns better associations among volumetric feature representations. These explicit local and global feature representations contribute to better learning of the shape-boundary characteristics of the organs. Extensive experiments on three different datasets reveal that the proposed method achieves favorable performance compared to state-of-the-art approaches. On the challenging Synapse Multi-organ dataset, the proposed method achieves an absolute 3.82% gain over the state-of-the-art approaches in terms of HD95 evaluation metrics while a similar improvement pattern is exhibited in MSD Liver and Pancreas tumor datasets. We also provide a detailed comparison between recent architectural design choices in the 2D computer vision literature by adapting them for the problem of 3D medical image segmentation. Finally, our experiments on the ZebraFish 3D cell membrane dataset having limited training data demonstrate the superior transfer learning capabilities of the proposed vMixer model on the challenging 3D cell instance segmentation task, where accurate boundary prediction plays a vital role in distinguishing individual cell instances. Our source code is publicly available at <https://github.com/Daniyanaj/vMixer>.

**Impact Statement**—Automatic medical image segmentation is a crucial step in the healthcare systems to perform accurate diagnoses, and pixel-wise estimation of cancerous tissues and organs from diverse medical images such as CT, MRI, and more. In 3D volumetric segmentation, transformers and hybrid approaches are exposed to self-attention, struggle to learn the inherent complex boundaries of the tissue and exhibit quadratic complexity with the number of tokens. The proposed vMixer framework exploits explicit local and global volumetric features to better learn the shape-boundary details of the organs. We also provide an extensive study on the selection of architectural design that is adapted for 3D medical segmentation from 2D vision literature for better boundary localization. Finally, we exploit the transfer learning capabilities of the proposed vMixer where training data is limited. We hope that our contributions can facilitate the research community to make use of artificial intelligence in designing a 3D image segmentation framework.

**Index Terms**—Attention, Medical Image Segmentation, MLP-

mixer, Transfer Learning

## I. INTRODUCTION

In clinical diagnosis, volumetric segmentation is a fundamental task that has shown promising potential in a wide range of applications including organ localization [1], [2] and tumor identification [3], [4]. UNet [5] is a breakthrough volumetric medical image segmentation approach that utilizes a CNN-based encoder and decoder architecture, where the encoder generates hierarchical low-dimensional features and decoder maps learned features into a voxel-wise segmentation. However, UNet [5] and its variants including V-Net [6] and ESPNet [7], often struggle to capture the long-range feature dependencies due to the limited receptive field of the convolution operation. This is particularly problematic in the case of multi-organ segmentation having large variations in the shapes and scales of organs.

Recently, transformers-based [8] approaches have been explored to capture global feature dependencies, which are further improved by hybrid approaches [9]–[12] that leverage the benefits of self-attention along with CNN components. Although self-attention learns the global pair-wise dependencies from volumetric 3D medical data, it struggles to capture the underlying complex boundaries of the tissues. Moreover, standard self-attention operates on pairwise patches which have quadratic complexity with respect to the number of tokens. In this work, our objective is to provide a hybrid architecture that helps to learn local as well as global representation to capture better shape-boundary information of the complex organs (shape of the organs). Motivated by the success of Multi-Layer Perceptron-mixers (MLP-mixers) [13] for image classification in natural images, we propose a novel volumetric medical image segmentation approach that strives to reduce the segmentation errors by introducing MLP-mixer and window-based self-attention blocks to explicitly capture global and local dependencies of volumetric feature representations, respectively.

Generally, the performance of volumetric segmentation approaches is evaluated using dice similarity coefficient (DSC), Hausdorff distance (HD95), and Normalised surface distance (NSD) metrics. DSC measures the overlap index between the ground-truth and predicted segmentation masks and NSD measures the overlap between the segmented boundaries. On the other-hand, Hausdorff distance (HD95) is measured as a distance between the boundaries of predicted and ground-truth segmentations [14], [15], hence, higher values for the

Daniya Najiha Abdul Kareem, Mustansar Fiaz and Hisham Cholakkal are with the Mohamed bin Zayed University of Artificial Intelligence, UAE.

Noa Novershtern and Jacob Hanna are with the Weizmann Institute of Science, Israel.

DSC score and NSD score, and a lower HD95 score indicates better model performance. Among these metrics, HD95 is the most informative and useful criterion as it indicates the largest segmentation error, especially where organs have varying sizes and shapes. Moreover, in applications where segmentation is an initial step in a complex multi-step process, the largest segmentation error evaluated in the HD95 score is a good indicator of the usefulness of segmentation for the given application [14]. Despite this, state-of-the-art volumetric segmentation approaches are sub-optimal in their HD95 evaluation metric.

Transfer learning has shown great progress in various areas including computer vision [16], natural language processing [17], medical image segmentation [18], and remote sensing [19]. Transfer learning provides a powerful technique to benefit from the existing knowledge and to improve the efficiency and effectiveness of the target application and domain. Therefore, in this paper, we also explore different convolutional-based, transformer-based, mixer-based, and hybrid-based architectural design choices to get leverage from transfer learning where the training data is limited. Our experimental study reveals that the proposed hybrid architectural design exhibits better transfer learning compared to other architectural designs.

#### A. Contributions

(i) In this work, we propose a hybrid hierarchical encoder-decoder framework, termed vMixer, that strives to capture both local and global information for accurate boundary prediction during volumetric medical image segmentation. Based on our comprehensive studies, we propose a novel MLP mixer-based framework for medical image segmentation that utilizes volume-based self-attention (Swin attention) to capture local dependencies at the high-resolution stage and introduces a novel Global Volume Mixer (GVM) block to encode the global dependencies at lower-resolution stages. This explicit utilization of global and local representation leads to better learning of organ boundary regions.

(ii) We perform a comprehensive comparison between different architectural design choices available in the literature by adapting them for 3D medical image segmentation. To the best of our knowledge, we are the first to have such a comprehensive study, and we hope that our study will support future research in this direction.

(iii) Our comprehensive experiments on three different datasets (Synapse Multi-organ, MSD-Pancreas Tumour, and MSD-Liver Tumour) across HD95, Dice, and NSD evaluation metrics show the merits of the proposed approach. In addition, our experiments on ZebraFish 3D cell membrane dataset (with limited training data) show that the proposed vMixer exhibits superior transfer learning abilities for challenging 3D cell instance segmentation task where accurate boundary prediction is crucial for delineating different cell instances.

## II. RELATED WORK

3D volumetric segmentation helps in identifying the regions of interest within a 3D volume and aids in diagnostic imaging,

treatment evaluation, and surgical planning by accurate delineation of anatomical structures or abnormalities. Numerous studies have been evidenced by the literature which encompasses a wide range of approaches among which CNN-based models, transformer-based models, and hybrid models remain the most prominent ones. U-Net proposed by Long et al. [20] has gained popularity over the last decade and is widely utilized in image segmentation methods. These models, built upon the U-Net framework, continue to be at the forefront of cutting-edge design among the image segmentation algorithms incorporating both CNNs and transformers.

#### A. Convolutional Neural Networks Based Models

Convolutional Neural Networks (CNNs) have gained significant popularity in 3D medical segmentation due to their ability to process volumetric data to effectively capture and learn hierarchical features. Multiple convolutional layers in CNNs help to exploit local spatial dependencies within the 3D data. Milletari et al. [21] introduced a 3D U-Net architecture that incorporates residual blocks instead of cascaded CNNs, and pooling layers were replaced by strided convolutions. Additionally, a loss function based on the dice score was employed to address the issue of class imbalance among voxel groups. By leveraging these adaptations, the 3D U-Net architecture improved the accuracy and effectiveness of 3D medical segmentation. Rehman et al. [22] extended the MaxVit [23] into the UNet architecture and used it for 2D cell segmentation tasks similar to [24]–[26]. Peng et al. [27] introduced a variant of the 3D U-Net architecture which involved utilizing multiple U-Net modules for the extraction of long-range spatial details at different scales. Xception blocks were employed instead of conventional convolutions within the U-Net blocks to enhance the feature extraction. Moreover, Peng et al. utilized 3D convolutions with depthwise separable convolutions aiming to reduce the computational complexity. The nn-UNet framework proposed in [1] is another noteworthy contribution to the field of 3D medical image segmentation. This framework is essentially built upon three simple U-Net models and has introduced many optimal strategies for pre-processing, training, inference computations, and post-processing that greatly contribute to effective network implementation. Since its introduction, the nn-UNet framework has been widely used in the latest image segmentation methods.

The resolution of 3D images is a major cause of concern which often demand substantial computational resources while training a 3D model. Recently, patch-based segmentation methods have been developed to address the issue. Kamnitsas et al. [28] proposed a dual-network architecture for brain-lesion segmentation that captures information from two distinct receptive fields and simultaneously learns features at different scales. In this model, dense fully connected layers are utilized to classify voxels into different groups from patchified input images. With this approach, the accuracy and efficiency of brain lesion segmentation in 3D medical images have shown remarkable improvement compared to other SOTA methods.

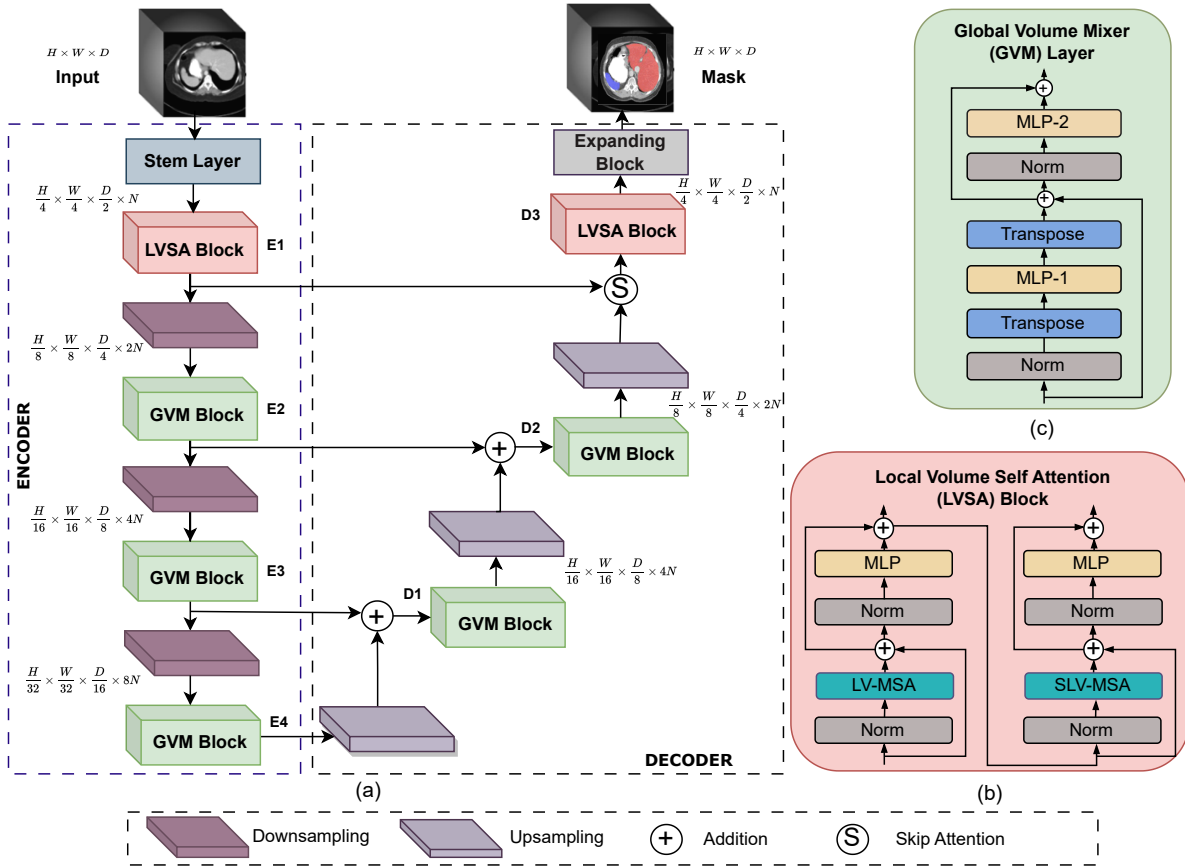


Fig. 1. (a) Overview of the proposed vMixer framework with hierarchical encoder-decoder architecture. The focus of our design is to explicitly capture the local and global feature dependencies for accurate segmentation. Our framework takes 3D images as input and employs local volume self-attention (LVSA) block to explicitly learn the local dependencies at high resolution ( $E1$ ,  $D3$ ). The  $E1$  features are downsampled and passed to the proposed global mixer block to explicitly learn the global dependencies. In the decoder, the features are first upsampled and then fused with the encoder features through a skip connection. We employ global volume mixer blocks at the first two decoder stages ( $D1$  and  $D2$ ) and a LVSA block at the last stage of the decoder ( $D3$ ). The final decoder features are fed to an expanding layer for producing the final segmentation mask. (b) Presents the LVSA block which comprises of local volume-based multi-head self-attention (LV-MSA) layer followed by a shifted local volume-based multi-head self-attention (SLV-MSA) layer. (c) Shows the structure of the volumetric MLP-mixer layer used in the GVM block. Each GVM block comprises two MLP-mixer layers. The volumetric MLP-mixer layer performs token mixing and channel mixing operations on the input volumetric tokens.

## B. Transformer Based Models

Transformers have demonstrated exceptional performance in the analysis of volumetric medical images by effectively modeling sequential data and encoding long-range dependencies. Transformers excel in capturing global interactions and contextual information across the entire 3D volume, unlike traditional CNN-based approaches that mainly focus on local spatial relationships [29], [30]. Karimi *et al.* [31] introduced a transformer-based model manipulating the self-attention mechanism between consecutive linear embeddings of image patches, enabling the model to effectively capture spatial relationships and dependencies. Additionally, an effective pre-training method was also adopted for the model, which proved valuable in scenarios where only limited annotated data is available. Experiments conducted on three different medical 3D datasets, including the hippocampus, pancreas, and brain cortical plate demonstrated the model's efficiency.

Swin Transformers introduced by Liu *et al.* [32] reduces the complexity linearly with the number of tokens compared to the quadratic computational complexity in ViTs [33] by utilizing window-based multi-scale attention. In addition, it allows a variable token size which enables them to handle objects of variable scales commonly found in medical images. Swin Transformers capture both local and global information effectively and also incorporate hierarchical feature extraction through patch merging layers and repeated Swin Transformer blocks at different image scales. Cao *et al.* [8] introduced Swin-UNet, an architecture that integrates Swin transformer blocks into the traditional U-Net framework in an encoder-decoder design. The Swin-UNet architectural design achieved promising segmentation accuracies for Synapse and ACDC datasets.

### C. Hybrid Architectures

In recent studies, the dominant approach utilizes U-shaped architectures incorporating transformers and CNNs through various strategies incorporating multi-scale feature extraction techniques and self-attention layers in the network [34]. The excellence of CNNs in capturing local spatial features and extracting hierarchical representations makes them well-suited for encoding fine-grained details and local patterns within the 3D volume. Transformers, on the other hand, have the ability to model global dependencies and encompass long-range contextual information. As identification of the complex relationships between different regions and structures within the volume is important in 3D medical image segmentation, hybrid architectures play a significant role.

The TransUNet was proposed by Chen et al. [35] which utilizes a transformer for global feature encoding and CNN for high-resolution feature extraction. Transformer features are added with skip connections from the encoder network for precise localization. This network achieved better performance in organ segmentation compared to ViTs. In [36], UNETR was designed by Hatamizadeh et al. based on ViTs [33] that includes a transformer encoder and a CNN decoder which is connected using skip connections. This network achieves good performance on MSD and BTCV segmentation datasets. In another approach, Hatamizadeh et al. [37] introduced Swin UNETR, an architecture based on Swin UNet, incorporating Swin Transformers in the encoder and utilizing a CNN decoder. This model stands out as one of the top-performing approaches in the BraTs2021 MRI dataset and Synapse multi-organ CT dataset. With an efficient window partitioning scheme and attention mechanism, Swin UNETR showcases its capabilities in achieving accurate segmentation results.

Many of the aforementioned studies focused on encoding global context by incorporating transformers along with CNN as supplementary modules but they did not fully optimize the integration of convolution and self-attention operations. To overcome this limitation, Zhou et al. introduced nnFormer [10] within the nn-UNet [1] framework, presenting an interleaved architecture combining both convolution and self-attention mechanisms. nnFormer generates hierarchical features and expands the receptive fields by effectively leveraging both local volume-based and global volume-based self-attention mechanisms. Compared to existing methods, we propose a hybrid architectural design that benefits in the extraction of local details using self-attention at high-resolution and global dependencies using MLP-Mixer at low-resolution features for volumetric 3D medical image segmentation. This combination handles the complexity associated with 3D image segmentation tasks and helps in achieving an optimal design choice. We also leverage the transfer learning capabilities of our design over a 3D cell membrane dataset.

## III. THE PROPOSED METHOD

**Motivation:** As discussed earlier, state-of-the-art transformer-based and hybrid approaches generally employ self-attention operations to obtain high-quality segmentation results. However, these approaches often struggle to predict accurate organ

boundaries. Here, we argue that it is desired to learn the boundary regions of the organs occurring in the local and larger spatial context. To this end, it is desired to analyze different architectural choices in the 2D computer vision literature by adapting them for 3D medical image segmentation. In this work, we explicitly learn local dependencies at high-resolution features while capturing global dependencies at lower-resolution features. This approach benefits from learning better associations among the volumetric feature representation, which leads to better prediction of organ boundary regions.

### A. Overall Architecture

Fig. 1-(a) shows our overall hierarchical encoder-decoder framework. The proposed framework has four encoder and three decoder stages. As discussed earlier, accurate organ boundary segmentation is a complex task that requires local contextual information for the precise delineation of boundary pixels from the background, while it simultaneously requires global contextual information to prevent erroneous predictions. Therefore, to learn shape-boundary information about the varying shapes of organs, we propose to explicitly capture local and global dependencies. To be specific, we capture local feature dependencies at the first encoder and last decoder stages having the highest feature resolution, and global feature dependencies are captured at the remaining encoder and decoder stages having relatively lower feature resolutions. The first stage of the encoder consists of the stem layer followed by the local volumetric self-attention (LVSA) block which comprises local volume-based multi-head self-attention and shifted local volume-based multi-head self-attention layers. While the latter three encoder stages are composed of a downsampling layer followed by a global volumetric mixer (GVM) block (which has two global volumetric MLP-mixer layers). GVM block at low-resolution features exhibits a holistic approach which helps in better extraction of global detail that can possibly learn complex organ shapes in the volumetric context by performing token mixing. Similar to the encoder, the decoder also follows a multi-stage hierarchical architecture. Each decoder stage utilizes upsampling layer to increase the feature resolution followed by a GVM block. Finally, in the last stage of the decoder, we follow a structure similar to the first encoder stage i.e., we employ a LVSA block followed by a feature-expanding layer that predicts the final masks.

### B. Local Volume Self Attention

Our encoder takes 3D input images as input to the stem layers. These stem features and last-stage decoder features have high-resolution features. Hence, applying self-attention on uniformly sampled dense patches from these high-resolution feature maps leads to a quadratic complexity with respect to the number of tokens. To learn the explicit local dependencies, we adopt a local volume-based self-attention block instead of 2D local windows as in Swin Transformer [38], which has reduced computational complexity compared to standard self-attention. Similar to the Swin Transformer block, we endorse



TABLE I

COMPARISON WITH OTHER STATE-OF-THE-ART METHODS OVER SYNAPSE MULTI-ORGAN DATASET. THE BEST RESULTS ARE IN BOLD.

Method	DSC	HD95	NSD
UNet [5]	76.85	-	-
ViT [39]+CUP [35]	67.86	36.11	-
R50-ViT [39] + CUP [35]	71.29	32.87	-
TransUNet [35]	77.48	31.69	-
SwinUNet [8]	79.13	21.55	-
MissFormer [9]	81.96	18.20	-
UNETR [36]	79.56	22.97	85.34
Swin UNETR [37]	80.24	17.65	85.75
nnFormer [10]	<b>86.57</b>	10.63	92.04
<b>vMixer (Ours)</b>	86.53	<b>6.78</b>	<b>92.96</b>

local feature encoding with the help of a local volume-based multi-head self-attention layer followed by a locally shifted volume-based multi-head self-attention layer as shown in Fig. 1-(b).

### C. Global Volume Mixing

As mentioned earlier, to learn the complex shapes of the organs in the volumetric context, we propose a global volume mixer block to explicitly capture the global dependencies from low-resolution stages. Standard self-attention [39] operates over dense patches that have quadratic complexity with respect to the total number of tokens. MLP-Mixer [13] can possibly learn complex relationships across the entire input, which makes them effective at capturing global information by performing token mixing followed by a pointwise feature refinement. In contrast to other context aggregators [33], [40], MLP-mixer is more dense, static, and does not require parameter sharing [41]. The core operation of MLP-mixer is the dense transposed affinity matrix on a single feature group. Therefore, we introduce a global volume mixer (GVM) block that has two volumetric MLP-mixer layers to capture the global dependencies for the underlying feature representation to better learn the complex boundaries of the tissues. The volumetric MLP-mixer layer is composed of layer norm, transposed token-mixing, and a MLP with fully-connected layers with a GELU nonlinearity, as shown in Figure 1-(c).

Suppose  $\mathcal{F} \in \mathcal{R}^{H \times W \times D \times N}$ , which is reshaped to  $\mathcal{F} \in \mathcal{R}^{M \times N}$  and taken as input to volumetric global MLP-mixer layer, where  $M = (H \times W \times D)$  represents the size of the 3D input (volume) and  $N$  denotes the number of channels. The volumetric global MLP-mixer operations can be summarized as:

$$\begin{aligned} \hat{\mathcal{F}} &= (W_{mlp-1}(\text{Norm}(\mathcal{F})^T))^T + \mathcal{F}, \\ \bar{\mathcal{F}} &= W_{mlp-2}(\text{Norm}(\hat{\mathcal{F}}) + \hat{\mathcal{F}}), \end{aligned} \quad (1)$$

where  $W_{mlp-1}$  and  $W_{mlp-2}$  denote the learnable multi-perceptron layer weights and  $\hat{\mathcal{F}}$  and  $\bar{\mathcal{F}}$  represent the intermediate and final volumetric global mixing features, respectively.

### D. Additional Layers

**Stem Layer:** The stem layer is responsible to generate high dimensional tensor  $\mathcal{F}_s \in \mathcal{R}^{\frac{H}{4} \times \frac{W}{4} \times \frac{D}{2}}$  for the input image

$\mathcal{X} \in \mathcal{R}^{H \times W \times D}$ , where  $(H, W, D)$  is the shape of 3D input (volume). We apply four successive convolutional layers with kernel size 3 for tokenization. This reduces the computational complexity compared to the usage of large convolutional kernels as in ViT [39] and also helps in encoding pixel-level spatial details. Each convolution is followed by GELU activation and layer normalization operations.

**Down Sampling:** We perform the downsampling operation after each stage except the last stage utilizing convolution with a kernel size of 3 with a stride of 2. This helps in modeling objects at different scales as hierarchical details are obtained by convolution downsampling.

**Up-Sampling and Patch Expanding:** In contrast to down-sampling in the encoder, up-sampling is performed at each stage of the decoder with the help of convolutional upsampling. A 3D transposed convolution layer is used with a kernel size of 2 and a stride 2 that helps in upsampling the low-resolution feature maps to a higher resolution. In the decoder, these up-sampled feature maps are then added from the encoder to encode the fine-grained details along with the semantic details. At the last stage of the decoder, patch expanding is performed to obtain the final prediction masks using the deconvolutional operation.

## IV. EXPERIMENTS

### A. Dataset and Evaluation Metrics

We select four 3D medical segmentation datasets for evaluation, with diverse task objectives and versatile granularities. Synapse multi-organ dataset [42] aids for the segmentation of multiple organs while MSD Liver [43] and MSD Pancreas [43] datasets provide data for segmenting respective organs and tumors growths over it. On the other hand, ZebraFish Cell Membrane Dataset [44] is a cell instance segmentation dataset that has images with numerous cell targets present in it.

**Synapse Multi-organ dataset:** is a multi-organ dataset [42] having 30 abdominal CT scans has eight organs including the liver, right kidney, left kidney, pancreas, gall bladder, stomach, spleen, and aorta, and we use a train-test split of 18-12 scans with a resolution of  $512 \times 512 \times 160$ . The abdomen CT scans were acquired from a chemotherapy trial for colorectal cancer and ventral hernia study under the supervision of the Institutional Review Board (IRB).

**MSD Liver Tumour Dataset:** The Liver Tumour dataset was introduced as a part of the Medical Segmentation Decathlon Challenge [43]. It contains 131 3D contrast-enhanced CT scan images from patients with primary cancers and metastatic liver disease, as a consequence of colorectal, breast, and lung primary cancers with a resolution of  $512 \times 512 \times 482$ . Regions of interest include the liver and tumor regions inside the liver region. We selected this dataset due to its challenging aspect where a major label unbalance is present between the ROIs, i.e. between a larger liver region and a smaller tumor region. The data was collected from IRCAD Hopitaux Universitaires, Strasbourg which contained random samples from the 2017 Liver Tumor Segmentation (LiTS) challenge [43].

**MSD Pancreas Tumour Dataset:** This dataset consists of 282 3D CT scan (with a resolution of  $512 \times 512 \times 96$ ) volumes

TABLE II  
ORGAN-WISE SEGMENTATION COMPARISON BETWEEN UNETR, nnFORMER, AND OUR VMIXER OVER SYNAPSE MULTI-ORGAN DATASET. THE BEST RESULTS ARE IN BOLD.

Methods	Average		Aorta		Gall Bladder		Kidney(L)		Kidney(R)		Liver		Pancreas		Spleen		Stomach	
	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95	DSC	HD95
UNETR [36]	79.5	22.97	89.9	5.48	60.55	28.69	85.66	17.76	84.80	22.44	94.45	30.40	59.24	15.82	87.8	47.12	73.99	16.05
nnFormer [10]	<b>86.57</b>	10.63	<b>92.04</b>	11.38	70.17	11.55	86.57	18.09	86.25	12.76	<b>96.84</b>	<b>2.00</b>	<b>83.35</b>	<b>3.72</b>	<b>90.51</b>	16.92	86.83	8.58
Ours	86.53	<b>6.78</b>	90.63	<b>6.13</b>	<b>70.33</b>	<b>9.04</b>	<b>88.74</b>	<b>5.60</b>	<b>87.38</b>	<b>7.25</b>	96.74	2.43	80.34	4.53	89.71	<b>12.82</b>	<b>87.10</b>	<b>8.45</b>

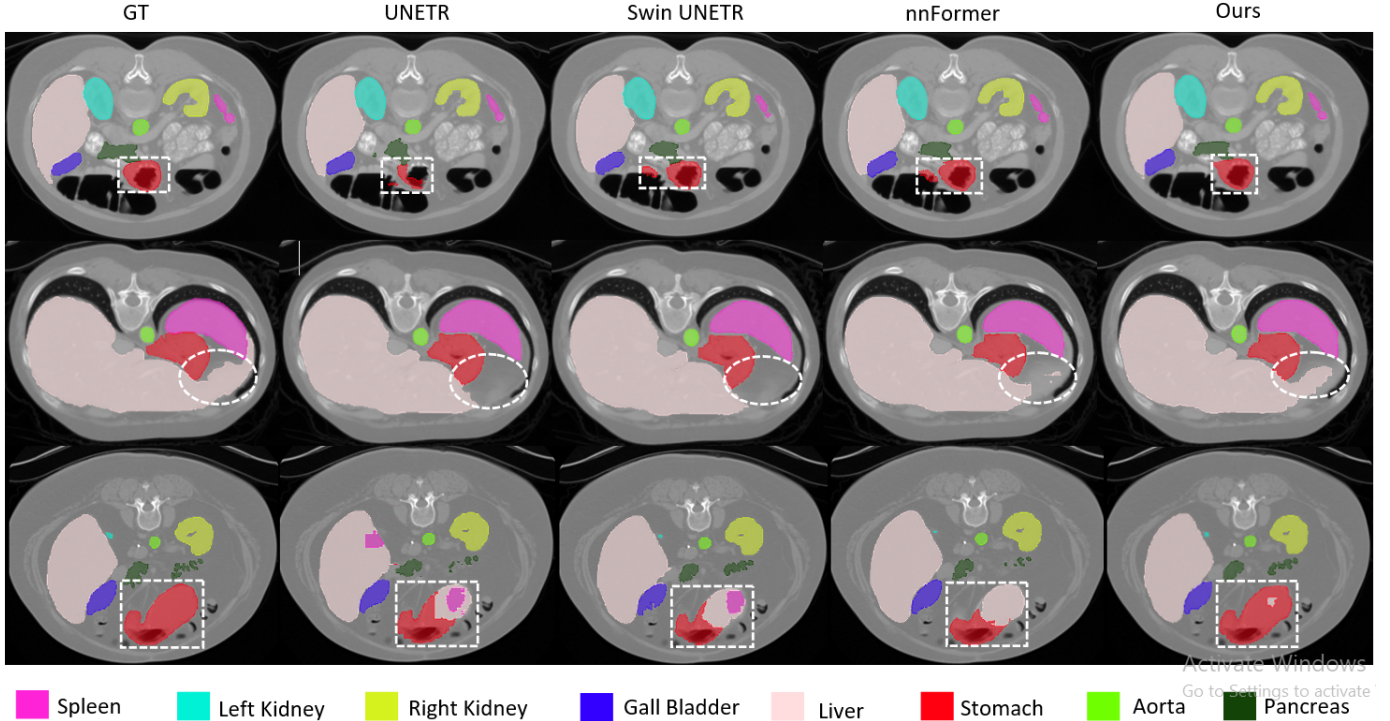


Fig. 2. Qualitative comparison on the Synapse multi-organ dataset. Our method provides improved segmentation by accurately detecting the boundaries of the organ.

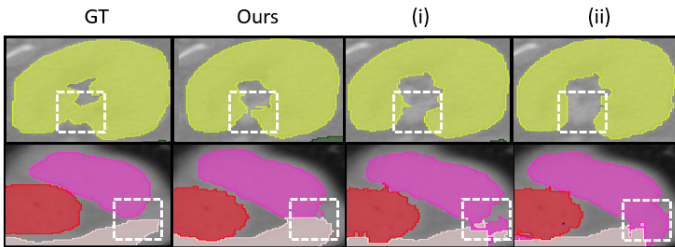


Fig. 3. **Qualitative comparison** of ablation experiments on Synapse multi-organ dataset. In a closer inspection, it can be seen that our method (LVSA (stage 1) and GVM (stage 2, stage 3, stage 4)) performs better than (i) GVM (Stages 1-4) and (ii) LVSA (Stages 1-4). Row 1 corresponds to the cross-section of the left kidney and row 2 shows a cross-section containing portions of the liver (light pink), spleen (magenta), and stomach (red). It can be clearly observed that our method has better shape preservation capabilities compared to other settings.

of patients undergoing resection of pancreatic masses which is also a subtask of [43]. The corresponding target ROIs were the pancreas organ and pancreatic mass (cyst or tumor). This dataset was also selected due to the label unbalance between small (tumor), medium (pancreas), and large (background) structures. The data was acquired at the Memorial Sloan Kettering Cancer Center, in New York, US. [43]

**ZebraFish Cell Membrane Dataset:** For transfer learning, we select 3D zebrafish cell dataset provided by the Department of Systems Biology at Harvard Medical School (HMS) that contains data for cell instance segmentation [44]. It comprises 36 images with a resolution of  $181 \times 331 \times 160$  which has 32 images for training and 4 images for the test. We utilized Dice (DSC) similarity, and Jaccard index (JI) as the evaluation metrics for the zebrafish cell dataset. We evaluate the performance using overall accuracy and cell count accuracy. Overall accuracy indicates the mean value of JI and DSC for all the cells, whereas cell count accuracy is represented as the fractions of cells whose JI or DSC is greater than 50% or 70%.

### B. Training Setup

The network is implemented using PyTorch 1.8.0 and trained using an NVIDIA GeForce RTX 3090 GPU. Following nnFormer [10], we adopt the pre-processing and augmentation strategies and set the batch size 2 and initial learning rate to 0.01. We utilize a poly decay strategy to adjust the learning rate (lr) as:

$$lr = \text{initial\_lr} \times \left(1 - \frac{\text{epoch\_number}}{\text{final\_epoch\_number}}\right)^{0.9}. \quad (2)$$

TABLE III  
COMPARISON OVER IN MSD LIVER TUMOUR DATASET. THE BEST RESULTS ARE IN BOLD.

Method	Liver			Tumor			Average		
	DSC	HD95	NSD	DSC	HD95	NSD	DSC	HD95	NSD
3D UNet [21]	94.37	-	-	53.94	-	-	74.15	-	-
UNETR [36]	86.69	20.79	83.06	50.94	44.68	50.14	68.85	32.73	66.64
Swin UNETR [37]	87.31	19.42	85.07	52.29	42.96	51.47	69.8	31.19	68.27
nnFormer [10]	94.87	12.77	93.00	55.78	30.77	60.42	76.75	21.77	76.72
Ours	<b>94.89</b>	<b>10.26</b>	<b>93.41</b>	<b>58.11</b>	<b>28.74</b>	<b>63.48</b>	<b>78.45</b>	<b>19.48</b>	<b>78.44</b>

TABLE IV  
COMPARISON OVER MSD PANCREAS TUMOUR DATASET. THE BEST RESULTS ARE IN BLUE

Method	Pancreas			Tumor			Average		
	DSC	HD95	NSD	DSC	HD95	NSD	DSC	HD95	NSD
3D-UNet [21]	69.20	-	-	35.64	-	-	52.42	-	-
UNETR [36]	71.91	13.97	86.47	35.93	26.41	52.92	53.92	20.19	69.68
Swin UNETR [37]	72.42	12.65	87.38	36.98	23.84	53.02	54.70	18.24	70.23
nnFormer [10]	78.80	6.04	94.91	44.36	12.83	62.33	61.53	9.42	78.62
Ours	<b>79.81</b>	<b>4.52</b>	<b>95.92</b>	<b>48.45</b>	<b>7.16</b>	<b>66.32</b>	<b>64.13</b>	<b>5.85</b>	<b>81.12</b>

During the training, we use outputs from intermediate and final prediction maps and compute the combined soft dice and cross-entropy losses as in [10]. We set the momentum and weight decay as 0.99 and 3e-5 with SGD optimizer and trained for 1000 epochs with 250 iterations per epoch. The final loss is calculated as,

$$\mathcal{L}_{all} = \alpha_1 \mathcal{L}_{\{H, W, D\}} + \alpha_2 \mathcal{L}_{\{\frac{H}{4}, \frac{W}{4}, \frac{D}{2}\}} + \alpha_3 \mathcal{L}_{\{\frac{H}{8}, \frac{W}{8}, \frac{D}{4}\}}. \quad (3)$$

Here,  $\alpha_{\{1, 2, 3\}}$  refers to weights for losses, and their values  $\alpha_{\{1, 2, 3\}}$  are reduced by half with respect to reduction in resolution i.e.  $\alpha_3 = \frac{\alpha_1}{4}$  and  $\alpha_2 = \frac{\alpha_1}{2}$ . All weights are finally normalized to one.

We follow nnUnet [1] and nnFormer [10] for the pre-processing and augmentation strategies. After pre-processing, we obtained crops with a resolution of 128x128x64 for Synapse Multi-organ, 128x128x128 for MSD Liver Tumour, and 40x224x224 for MSD Pancreas Tumour dataset.

### C. State-of-the-art Comparison

**Synapse Multi-organ dataset:** In Table I, we compare our vMixer over Synapse Multi-organ dataset with the existing state-of-the-art methods. MISSFormer [9], UNETR [36], swin UNETR [37] and nnFormer [10] have more than 79% DSC scores and HD95 scores are 18.20, 22.97, 17.65 and 10.63, respectively. Although our approach obtains a comparable 86.53% DSC score, it achieves a better 6.78 HD95 score. This indicates that our method is more capable of capturing the shape-boundary characteristics of the organs compared to other SOTA methods. Furthermore, in table II, we conduct a detailed performance analysis between our method and nnFormer [10] and UNETR [36]. Our method shows significant improvement in terms of HD95 scores. For example, it can be observed that for the smaller organs with complex boundaries such as the left kidney, right kidney, and gall bladder, our method achieves 5.60, 7.25, and 9.04 HD95 scores. In addition, qualitative comparison with other SOTA methods in 2. It can be clearly seen that our approach provides superior segmentation results compared to other methods and preserves the boundaries for the organs. For example, in figure 2 (row 1), our method better preserves the shape of the gall bladder organ.

### MSD Liver Tumour and MSD Pancreas Tumour Datasets:

We also compare our method with nnFormer over MSD Liver Tumour and MSD Pancreas Tumour datasets and showed consistent improvement as shown in tables III and IV, respectively. The qualitative results in figure 4 show that our method preserves better boundary information.

In addition, we also perform the statistical analysis to validate the significance of our method, we report the mean, median, and standard deviation of 3-fold experimental results over MSD Liver and MSD Pancreas datasets in Table V. The statistical significance analysis shows that our method clearly outperforms nnFormer [10] in all the quantities analyzed.

### D. Transfer Learning: 3D Cell Instance segmentation

**ZebraFish 3D cell membrane dataset:** We perform experiments on ZebraFish 3D cell membrane dataset (also named as HMS dataset) to validate the transfer learning abilities of our vMixer on the challenging 3D cell instance segmentation task where accurate boundary prediction is required to delineate different cell instances. In this study, we leverage the transfer learning abilities of different architectural choices by taking their pre-trained models from the multi-organ Synapse dataset and fine-tuned the models for the 3D cell instance segmentation task on HMS (ZebraFish 3D cell membrane) dataset having limited training data. For a fair comparison, we set the same fine-tuning parameters for all architectures, across all our fine-tuning experiments. We set the batch size as 5, maximum epochs to 500, and the learning rate as 1e-4 using the Adam optimizer. We set the input size of  $64 \times 64 \times 64$  and utilize weighted Dice loss for loss backpropagation. To show the generalizability of our method, we present different transfer learning approaches across various architectural choices over the HMS dataset, as shown in Tab. IX. Please note that the models trained on the Synapse multi-organ dataset can not be directly utilized on the HMS (ZebraFish 3D cell membrane) dataset due to the following reasons: (i) the input resolution differences (the input resolution for the Synapse multi-organ dataset is  $128 \times 128 \times 64$ , whereas the input size for the HMS dataset is  $64 \times 64 \times 64$ ). The stem needs to be learned for the target HMS dataset so that the stem can generate appropriate features for the model. (ii) Furthermore, the Synapse Multi-organ dataset has 14 classes and the HMS dataset has 3 classes. Therefore, there is a need to have different expanding layers to handle classes for the HMS dataset. Hence, we take the whole Synapse multi-organ model and make these two minimum modifications to the architecture to adapt to the HMS dataset.

To adapt the multi-organ Synapse weights for the HMS dataset, and show the generalizability of our method, we present two different ways to show the transfer learning abilities of the models over the HMS dataset, as shown in Tab. IX. To compare the generalizability of our model, we first fine-tune only the stem and patch expanding layers by freezing all remaining layers. This minimal learning validates that our frozen model weights have a better generalization ability compared to baseline methods (see Tab. IX column 1), even without fine-tuning most layers of the multi-organ Synapse pre-trained weights. Next, we perform end-to-end



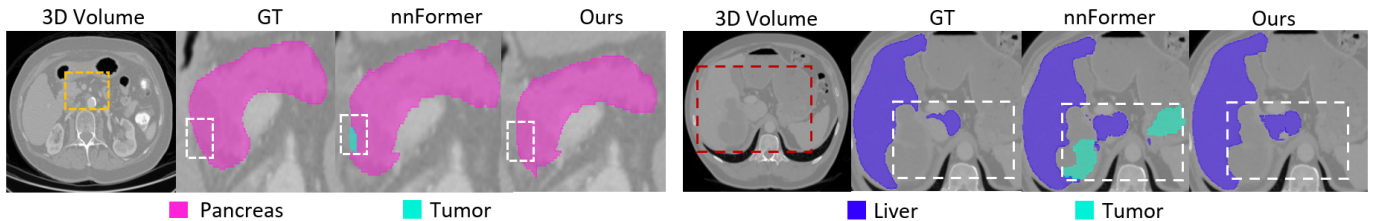


Fig. 4. Qualitative results on the MSD pancreas tumour (left) and MSD Liver tumour datasets. Our vMixer provides accurate segmentation of boundary regions.

TABLE V  
MEAN, MEDIAN, AND STANDARD DEVIATION OF PERFORMANCE SCORES OVER MSD LIVER AND MSD PANCREAS DATASETS. THE BEST RESULTS ARE IN BOLD.

	MSD Liver						MSD Pancreas					
	Mean DSC		Median DSC		Std Deviation DSC		Mean HD95		Median HD95		Std Deviation HD95	
	Liver	Tumor	Liver	Tumor	Liver	Tumor	Pancreas	Tumor	Pancreas	Tumor	Pancreas	Tumor
nnFormer [10]	94.13	55.16	94.25	54.90	0.89	1.06	77.91	43.14	76.92	44.11	0.93	1.21
vMixer (Ours)	<b>94.75</b>	<b>58.08</b>	<b>94.79</b>	<b>58.01</b>	<b>0.76</b>	<b>0.96</b>	<b>78.51</b>	<b>48.27</b>	<b>78.86</b>	<b>47.96</b>	<b>0.81</b>	<b>1.08</b>

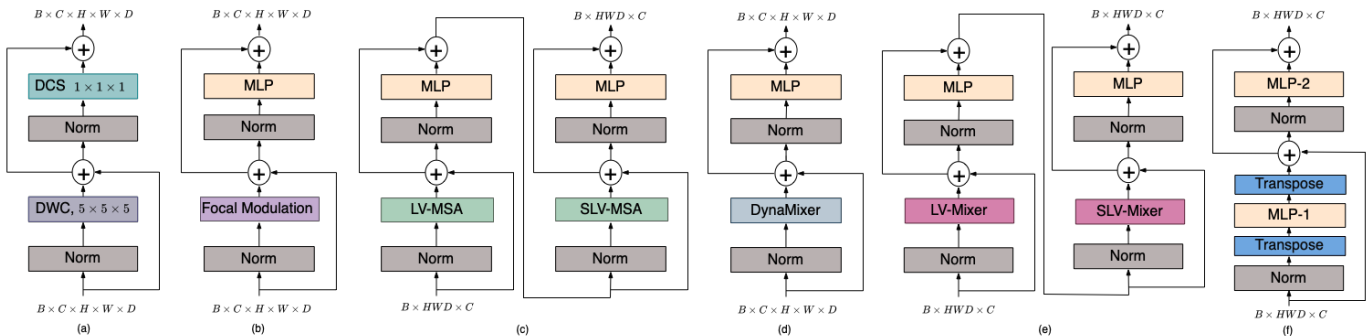


Fig. 5. Adaptation of different network architecture blocks for 3D medical image segmentation. The depth-wise convolution (DWC) and depth-wise scaling (DCS) based (a) ConvNeXt [46], (b) FocalNet [47] (see Fig. 6-b), and (d) DynaMixer [48] (see Fig. 6-a) operate on input 3D volume of size  $B \times C \times H \times W \times D$ . The (c) Swin Transformer [38], (e) Swin Mixer [38], and (f) GVM blocks reshape the 3D volume to  $B \times HWD \times C$  dimensional features. Here, B: batch, C: channel, H: height, W: width, and D: depth dimensions, respectively.

TABLE VI

EXPERIMENTAL RESULTS AFTER END-TO-END FINE-TUNING SYNAPSE WEIGHTS ON ZEBRAFISH 3D CELL MEMBRANE DATASET [44]. WE REPORT THE RESULTS IN TERMS OF OVERALL ACCURACY (JI AND DSC), AND CELL COUNT ACCURACY (JI/DSC GREATER THAN 50% OR 70%) METRICS. OUR vMIXER PERFORMS SIGNIFICANTLY BETTER AGAINST EXISTING METHODS AND ACHIEVES STATE-OF-THE-ART PERFORMANCE. THE BEST RESULTS ARE IN BOLD.

	Avg JI	Avg DSC	JI>70%	DSC>70%	JI>50%	DSC>50%
FocalNet [45]	51.98	62.65	39.01	51.7	55.97	71.80
ConvNeXt [46]	52.50	64.01	39.30	53.50	56.09	73.60
nnFormer [10]	52.17	63.80	38.53	53.08	<b>55.73</b>	73.77
Ours	<b>54.25</b>	<b>65.69</b>	<b>42.60</b>	<b>58.11</b>	<b>60.91</b>	<b>76.54</b>

fine-tuning of the entire network and achieve better JI and DSC scores (see Tab. IX column 2). This empirical study reveals that our model has a consistent performance gain, demonstrating a superior transfer learning ability compared to different architectural choices. Figure 7 shows the qualitative results for an example from ZebraFish 3D cell membrane dataset. We note that our method segments the foreground, background, and boundary of the cells as well as preserves the boundaries for cells. In addition to that, our method provides the 3D cell instance segmentation results.

### E. Discussion of Different Architectural Designs:

In the context of accurate segmentation, especially for precise boundary prediction, it is important to thoroughly capture both local and global information. To address this challenge, we meticulously investigate various context aggregation techniques utilized in 2D image classification literature by adapting them to the domain of 3D medical image segmentation. The CNNs and transformers have already proven to be efficient for volumetric 3D medical imaging as in nnFormer [10] and UNETR [36], the capabilities of MLP-mixer [13] were not explored for the 3D medical image segmentation tasks. Moreover, there was little exploration regarding the *hybrid* combinations of multiple context aggregator blocks such as convolution, attention, and MLP-mixer for 3D volumetric medical image segmentation. In this work, we strive to explore the inherent characteristics of MLP-mixer based architectures when combined with other context aggregator blocks traditionally used in 2D literature, by adapting them to 3D in a hybrid design. Specifically, we explored different context aggregators including CNN-based (3DConv, ConvNeXt [46], and FocalNet [47]), transformer-based (self-attention [39], Swin Transformer [38] as LSVA), and mixers-based (Swin Mixer [38], DynaMixer [48], MLP-mixer [13] as global volume mixer (GVM)) to



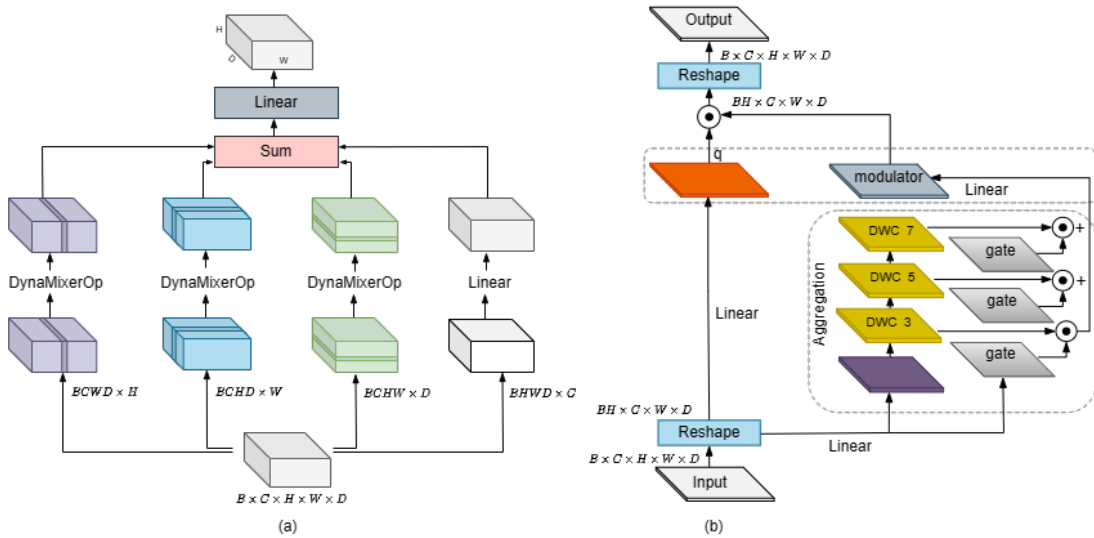


Fig. 6. (a) **DynaMixer [48] Architecture**. The architecture includes height, width, depth, and channel mixing. In height mixing, the tokens are mixed via DynaMixer operations (DynaMixerOp) along the height dimension using shared weights. Similar operations are performed for width and depth dimensions, respectively. The channel mixing performs linear transformation of the features. DynaMixer operation (DynaMixerOp) includes mixing elements along the dimension with a series of linear layers followed by a Softmax activation. (b) **Focal Modulation [47] Architecture**. Aggregation consists of hierarchical gated context aggregation (best viewed in Zoom).

TABLE VII

ANALYSIS OF VARIOUS ARCHITECTURES ADAPTED FOR 3D MEDICAL IMAGE SEGMENTATION IN A UNIFORM NETWORK WHERE THE SAME ARCHITECTURAL BLOCK IS USED IN ALL ENCODER AND DECODER STAGES OVER SYNAPSE MULTI-ORGAN DATASET. HERE, WE TAKE REPRESENTATIVE CONVOLUTIONAL, TRANSFORMER, AND MLP MIXER, ARCHITECTURES AND INTRODUCE THEM TO A UNIFORM nnFORMER-BASED ARCHITECTURE. THE INFERIOR **HD95** SCORES OF THESE *uniform* ARCHITECTURES COMPARED TO THE HD95 SCORE OBTAINED BY OUR *hybrid* VMIXER (GVM IN ALL STAGES: 9.90 VS OURS: **6.78**) DEMONSTRATE THE BENEFITS OF THE PROPOSED **HYBRID ARCHITECTURE** IN ENCODING LOCAL AND GLOBAL DETAILS FOR ACCURATE ORGAN SEGMENTATION. THE STRUCTURE OF THESE BLOCKS IS SHOWN IN FIG 5. THE BEST RESULTS ARE IN BOLD

Convolutional Networks			Transformers		Mixers		
3D Conv	ConvNeXt [46]	FocalNet [47]	Self Attention [39]	LVSA [10]	Swin Mixer [38]	DynaMixer [48]	GVM
10.78	32.80	20.66	16.00	10.63	10.10	11.80	<b>9.90</b>

TABLE VIII

ABLATION EXPERIMENTS ON VARIOUS **HYBRID ARCHITECTURE** DESIGN CHOICES OVER SYNAPSE MULTI-ORGAN DATASET. THE RESULTS SHOW THAT **LVSA FOR CAPTURING LOCAL INFORMATION** AT STAGE 1 AND USING **GVM TO CAPTURE GLOBAL INFORMATION** AT THE REMAINING STAGES PROVIDES THE BEST RESULT. THE BEST RESULT IS IN BOLD.

Stage	LVSA	LVSA	LVSA	LVSA	ConvNeXt	FocalNet	DynaMixer	Self-attention	Swin Mixer	LVSA
Stage 1	LVSA	LVSA	LVSA	LVSA	ConvNeXt	FocalNet	DynaMixer	Self-attention	Swin Mixer	LVSA
Stage 2	LVSA	LVSA	Swin Mixer	LVSA	GVM	GVM	GVM	GVM	GVM	GVM
Stage 3	LVSA	GVM	Swin Mixer	Swin Mixer	GVM	GVM	GVM	GVM	GVM	GVM
Stage 4	GVM	GVM	Swin Mixer	Swin Mixer	GVM	GVM	GVM	GVM	GVM	GVM
<b>HD95</b>	9.01	8.60	10.89	11.57	7.72	11.40	11.60	13.79	11.10	<b>6.78</b>

TABLE IX

ANALYSIS OF TRANSFER LEARNING EXPERIMENTS ON ZEBRAFISH 3D CELL MEMBRANE DATASET. THE BEST RESULTS ARE IN BOLD.

	Fine-tuning stem and expanding layers		End-to-end fine-tuning	
	Avg JI	Avg DSC	Avg JI	Avg DSC
FocalNet [47]	21.96	26.18	51.98	62.65
ConvNeXt [46]	24.12	29.42	52.5	64.01
nnFormer [10]	23.89	29.26	52.17	63.82
<b>vMixer (Ours)</b>	<b>25.17</b>	<b>30.13</b>	<b>54.25</b>	<b>65.69</b>

TABLE X

STUDY OF VMIXER STAGES SUITABLE FOR CAPTURING GLOBAL (GVM) AND LOCAL (LVSA) DEPENDENCIES OVER SYNAPSE MULTI-ORGAN DATASET. THE BEST RESULTS ARE IN BOLD.

	Stage 1	Stage 2	Stage 3	Stage 4	HD95
Uniform Architecture	LVSA	LVSA	LVSA	LVSA	10.60
	GVM	GVM	GVM	GVM	8.60
Hybrid Architecture	LVSA	LVSA	GVM	GVM	9.90
	LVSA	GVM	GVM	GVM	6.78

learn both local and global context information to learn better features for medical image segmentation. Figure 5 illustrates the adaptation of different network design blocks.

Our exploration begins with a detailed examination of CNN-based architectures. We observe that traditional 3DConv networks make use of 3D convolutions, which are indeed well-suited for local feature extraction but often prove inadequate in their ability to capture long-range dependencies. We also

evaluate the potential of ConvNeXt, an architecture that modernizes conventional convolutional networks by introducing larger kernel sizes and adopting fewer activation functions and normalization layers. We also analyze FocalNets from 2D image recognition literature that employ a focal modulation technique. FocalNets comprise a hierarchical contextualization and a gated aggregation, followed by a modulation strategy.

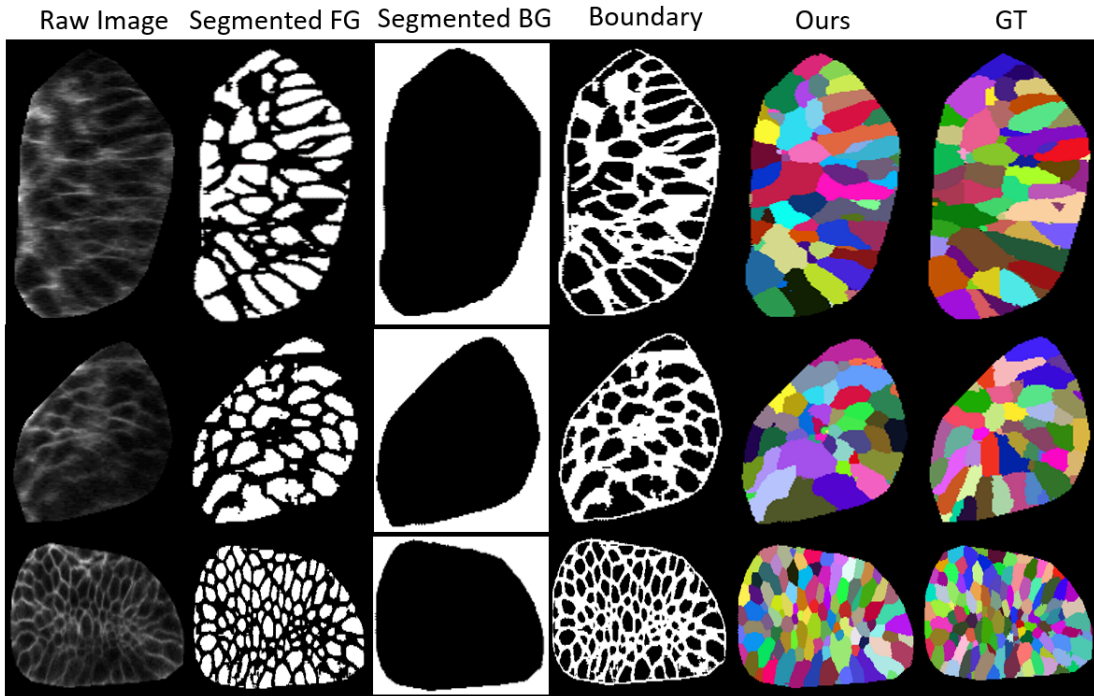


Fig. 7. Qualitative results of our vMixer fine-tuned on the Zebrafish 3D cell membrane dataset using pre-trained weights from multi-organ Synapse dataset. Rows 1,2 and 3 correspond to views from different planes. The proposed vMixer predicts foreground (FG), background (BG), and cell boundary regions (columns 2-4, respectively) which are post-processed using a watershed algorithm to accurately segment cell instances (column 5).

This distinctive approach gathers context information, spanning from short to long-range, and the process of aggregation is influentially guided by the content of the query.

We also investigate different transformer-based architectures, commencing with the standard self-attention mechanism. While self-attention is capable of encoding long-range dependencies, it may encounter challenges when handling local information. Moreover, it is notably associated with quadratic complexity concerning token length. However, the Swin Transformer enhances locality by adopting a shifted window attention strategy, effectively refining the receptive field while concurrently reducing computational complexity. We evaluate the 3D self-attention and 3D Swin attention in our experiments.

Finally, we also evaluate MLP mixers-based architectures as context aggregators. MLP-Mixer leverages multi-layer perceptrons to globally mix features, demonstrating its inherent capability to capture complex relationships within the input data. Our proposed Global Volume Mixer (GVM) module is fundamentally derived from the MLP-Mixer architecture, incorporating token-mixing and channel-mixing operations to effectively extract global characteristics from 3D volumetric data. In addition, we also adapt DynaMixer, which is distinguished by its dynamic generation of mixing matrices, predicated on token content analysis, thereby improving computational efficiency and overall robustness. *To the best of our knowledge, this is the first time MLP-mixer based context aggregators are adapted for the problem of 3D medical image segmentation.*

To discern the most suitable architectural design, we thoroughly evaluate the performance of each context aggregator

when uniformly applied to all encoder-decoder stages. Significantly, GVM emerges as the top performer, yielding a reduced HD95 score of 9.9, as shown in Tab. VII. Further, to optimally capture both local and global information, we explore *hybrid* combinations of these context aggregators. Our experimental results endorse the implementation of LVSA at stage 1, complemented by GVM blocks in subsequent stages, as shown in Tab. VIII. Our comprehensive experiments demonstrate the advantages of the proposed hybrid architecture. This configuration, featuring LVSA at the highest resolution and GVM blocks at successive lower-resolution levels, firmly establishes itself as the preferred choice for accurate boundary prediction, as evidenced by the noteworthy improvements in HD95 scores.

#### F. Ablation Study

We perform an ablation study over the Synapse Multi-organ dataset to validate the effectiveness of the proposed vMixer. The proposed vMixer has four encoder stages and three decoder stages. In all our ablation experiments, we use the same architectural design choices for the encoder and decoder stages having the same resolution. For example, the first encoder ( $E1$ ) and the last decoder ( $D3$ ) use the same architecture and we refer to this as stage 1. Similarly, the remaining stages (stages 2-4) are also defined based on the respective encoder and decoder stages.

As discussed earlier, we take the different architectures designed for 2D detection tasks and adapt them for 3D medical image segmentation, including convolutional network designs (3D Convs, ConvNeXt [46], FocalNet [45]), transformers

(self-attention [39], Swin Transformer [38] as local volume-based self-attention (LVSA)), and mixers (Swin Mixer as local volume-based MLP-mixer, DynaMixer [48], and MLP-mixer [13] as global volume mixer (GVM)). In table VII, we present an analysis of different architectures adapted for 3D medical image segmentation, where the encoder-decoder has the same architectural block. We observe that GVM performs better with a reduced HD95 score of 9.90. We also perform an ablation study with various hybrid architecture design choices. In Table VIII, we set different design choices at different stages and observe that our vMixer presents a better HD95 score compared to all other hybrid design choices. Finally, we fix the GVM in all stages except the first stage of the encoder and the last stage of the decoder, as shown in table VIII. We observe that employing ConvNext or LVSA in the first stage of the encoder and the last stage of the decoder results in superior HD95 scores of 7.72 and 6.78, respectively. Based on the aforementioned ablation studies, we fix LVSA at the first stage and GVM at the last three stages of the proposed vMixer.

Next, we employ different combinations of LVSA and GVM blocks at different stages as shown in Table X. We can observe that employing GVM at low-resolution stages helps to better capture global information and optimum performance can be obtained by employing LVSA (similar to nnFormer) at the first stage and GVM at all remaining stages (row 3). The LVSA at stage 1 and GVM blocks at the last 3 stages provide a favorable HD95 score which indicates its better capability to preserve the shapes of different organs. We also present the qualitative results in Figure 3 and show that our architecture preserves better boundary information. When LVSA is uniformly employed in all stages, the method fails to capture global dependencies and hence loses characteristic details associated with the overall organ shape (as shown in row-2 of Figure 3-(ii)). On the other hand, our study demonstrates that our novel hybrid architecture employing LVSA in the highest resolution and GVM blocks in successive levels helps in a refined feature extraction from the 3D medical volume aiding improved boundary detection and facilitating better segmentation.

**Training Epochs vs Model Convergence:** We perform experiments with 1000 maximum epochs to be consistent with our baseline (nnFormer [10]). We have extended the experiments by increasing the number of epochs to 1200 as shown in Figure 8. It can be seen that our method has a clear dominance in performance scores and achieves faster convergence compared to the nnFormer, especially in terms of HD95 score (as shown in Fig. 8-right). We observe that there is enormous performance improvement for both nnFormer and vMixer until 900 training epochs, and later, optimal performance is achieved between the 900th and 1000th epoch. Thereafter, the model performance did not undergo a noticeable change, and hence the total number of epochs was set to 1000. We also observe from Fig. 8 that training for fewer epochs may not be conclusive about the exact performance trend as 3D datasets are complex and demand larger network architectures.

**Computational Cost Analysis:** We present the computational analysis of our method. We present the comparison of

Fig. 8. Comparison between nnFormer [10] and vMixer with respect to the number of epochs over Synapse Multi-organ dataset.

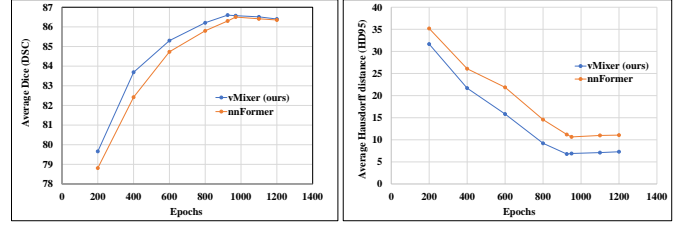


TABLE XI  
COMPUTATIONAL COST COMPARISON OVER SYNAPSE MULTI-ORGAN DATASET.

Model	Flops	Inference time (ms)	HD95
swin UNETR [37]	572	228.6	17.65
nnFormer [10]	212	148.0	10.6
vMixer	249	154.3	6.78

floating-point operations per second (FLOPs) and the inference time. The computational cost comparison for different methods over a multi-organ synapse dataset is presented in Table XI. It can be seen that our vMixer obtains a prominent dominance over other SOTA methods in terms of HD95 scores with comparable FLOPs and inference time. Therefore, there is a tradeoff between the accuracy and inference time.

## V. CONCLUSION

We propose a hierarchical encoder-decoder network to explicitly learn the local and global dependencies. We utilize the local volume-based self-attention to learn the local dependencies at high-resolution features and propose a volumetric global mixing mechanism to capture the global feature representations at low-resolution features. These explicit local-global feature representations benefit to capture the boundaries of the organs. Experimental study reveals that our approach provides favorable segmentation results at boundary regions compared to existing SOTA methods. Moreover, our experiments show that our vMixer provides promising 3D cell instance segmentation results on the Zebrafish cell 3D instance segmentation dataset.

## ACKNOWLEDGEMENT

This work is partially supported by the MBZUAI-WIS Joint Program for AI Research (Project grant number- WIS P008).

## REFERENCES

- [1] F. Isensee, J. Petersen, A. Klein, D. Zimmerer, P. F. Jaeger, S. Kohl, J. Wasserthal, G. Koehler, T. Norajitra, S. Wirkert *et al.*, “nnu-net: Self-adapting framework for u-net-based medical image segmentation,” *arXiv preprint arXiv:1809.10486*, 2018.
- [2] X. Xu, F. Zhou, B. Liu, D. Fu, and X. Bai, “Efficient multiple organ localization in ct image using 3d region proposal network,” *IEEE transactions on medical imaging*, vol. 38, no. 8, pp. 1885–1898, 2019.
- [3] A. Myronenko, “3d mri brain tumor segmentation using autoencoder regularization,” in *International MICCAI brainlesion workshop*. Springer, 2019, pp. 311–320.
- [4] W. Chen, B. Liu, S. Peng, J. Sun, and X. Qiao, “S3d-unet: separable 3d u-net for brain tumor segmentation,” in *International MICCAI Brainlesion Workshop*. Springer, 2019, pp. 358–368.



- [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] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 fourth international conference on 3D vision (3DV)*. IEEE, 2016, pp. 565–571.
- [7] N. Nuechterlein and S. Mehta, "3d-espnet with pyramidal refinement for volumetric brain tumor image segmentation," in *International MICCAI Brainlesion Workshop*. Springer, 2019, pp. 245–253.
- [8] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," 2021. [Online]. Available: <https://arxiv.org/abs/2105.05537>
- [9] X. Huang, Z. Deng, D. Li, and X. Yuan, "Missformer: An effective medical image segmentation transformer," 2021. [Online]. Available: <https://arxiv.org/abs/2109.07162>
- [10] H.-Y. Zhou, J. Guo, Y. Zhang, L. Yu, L. Wang, and Y. Yu, "nnformer: Interleaved transformer for volumetric segmentation," 2021. [Online]. Available: <https://arxiv.org/abs/2109.03201>
- [11] M. Fiaz, M. Heidari, R. M. Anwar, and H. Cholakkal, "Sa2-net: Scale-aware attention network for microscopic image segmentation," *arXiv preprint arXiv:2309.16661*, 2023.
- [12] D. N. A. Kareem, M. Fiaz, N. Novershtern, and H. Cholakkal, "Medical image segmentation using directional window attention," 2024. [Online]. Available: <https://arxiv.org/abs/2406.17471>
- [13] I. O. Tolstikhin, N. Houlsby *et al.*, "Mlp-mixer: An all-mlp architecture for vision," *Advances in Neural Information Processing Systems*, vol. 34, pp. 24 261–24 272, 2021.
- [14] D. Karimi and S. E. Salcudean, "Reducing the hausdorff distance in medical image segmentation with convolutional neural networks," *IEEE Transactions on medical imaging*, vol. 39, no. 2, pp. 499–513, 2019.
- [15] A. A. Taha and A. Hanbury, "Metrics for evaluating 3d medical image segmentation: analysis, selection, and tool," *BMC medical imaging*, vol. 15, no. 1, pp. 1–28, 2015.
- [16] Y. Guo, H. Shi, A. Kumar, K. Grauman, T. Rosing, and R. Feris, "Spot-tune: transfer learning through adaptive fine-tuning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4805–4814.
- [17] Y.-L. Sung, J. Cho, and M. Bansal, "VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5227–5237.
- [18] S. Chen, K. Ma, and Y. Zheng, "Med3d: Transfer learning for 3d medical image analysis," *arXiv preprint arXiv:1904.00625*, 2019.
- [19] B. Cui, X. Chen, and Y. Lu, "Semantic segmentation of remote sensing images using transfer learning and deep convolutional neural network with dense connection," *Ieee Access*, vol. 8, pp. 116 744–116 755, 2020.
- [20] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," 2014. [Online]. Available: <https://arxiv.org/abs/1411.4038>
- [21] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3d u-net: Learning dense volumetric segmentation from sparse annotation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:2164893>
- [22] A. Rehman and A. Khan, "Maxvit-unet: Multi-axis attention for medical image segmentation," *arXiv preprint arXiv:2305.08396*, 2023.
- [23] Z. Tu, H. Talebi, H. Zhang, F. Yang, P. Milanfar, A. Bovik, and Y. Li, "Maxvit: Multi-axis vision transformer," in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIV*. Springer, 2022, pp. 459–479.
- [24] R. C. Aralikatti, S. Pawan, and J. Rajan, "A dual-stage semi-supervised pre-training approach for medical image segmentation," *IEEE Transactions on Artificial Intelligence*, 2023.
- [25] A. Saif, T. Imtiaz, S. Rifat, C. Shahnaz, W.-P. Zhu, and M. O. Ahmad, "Capscovnet: A modified capsule network to diagnose covid-19 from multimodal medical imaging," *IEEE Transactions on Artificial Intelligence*, vol. 2, no. 6, pp. 608–617, 2021.
- [26] T. Mahmud, M. A. Rahman, S. A. Fattah, and S.-Y. Kung, "Covsegnet: A multi encoder-decoder architecture for improved lesion segmentation of covid-19 chest ct scans," *IEEE Transactions on Artificial Intelligence*, vol. 2, no. 3, pp. 283–297, 2021.
- [27] S. Peng, W. Chen, J. Sun, and B. Liu, "Multi-scale 3d u-nets: an approach to automatic segmentation of brain tumor," *International Journal of Imaging Systems and Technology*, vol. 30, no. 1, pp. 5–17, 2020.
- [28] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker, "Efficient multi-scale 3d CNN with fully connected CRF for accurate brain lesion segmentation," *Medical Image Analysis*, vol. 36, pp. 61–78, feb 2017. [Online]. Available: <https://doi.org/10.1016%2Fj.media.2016.10.004>
- [29] Z. Fang, J. Bai, X. Guo, X. Wang, F. Gao, H.-Y. Yang, B. Kong, Y. Hou, K. Cao, Q. Song *et al.*, "Annotation-efficient covid-19 pneumonia lesion segmentation using error-aware unified semisupervised and active learning," *IEEE Transactions on Artificial Intelligence*, vol. 4, no. 2, pp. 255–267, 2022.
- [30] D. Karimi and A. Gholipour, "Improving calibration and out-of-distribution detection in deep models for medical image segmentation," *IEEE Transactions on Artificial Intelligence*, vol. 4, no. 2, pp. 383–397, 2022.
- [31] D. Karimi, S. Vasylechko, and A. Gholipour, "Convolution-free medical image segmentation using transformers," 2021. [Online]. Available: <https://arxiv.org/abs/2102.13645>
- [32] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," 2021. [Online]. Available: <https://arxiv.org/abs/2103.14030>
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017, pp. 6000–6010. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
- [34] A. Khan, Z. Rauf, A. Sohail, A. Rehman, H. Asif, A. Asif, and U. Farooq, "A survey of the vision transformers and its cnn-transformer based variants," *arXiv preprint arXiv:2305.09880*, 2023.
- [35] 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," 2021. [Online]. Available: <https://arxiv.org/abs/2102.04306>
- [36] A. Hatamizadeh, Y. Tang, V. Nath, D. Yang, A. Myronenko, B. Landman, H. Roth, and D. Xu, "Unetr: Transformers for 3d medical image segmentation," 2021. [Online]. Available: <https://arxiv.org/abs/2103.10504>
- [37] A. Hatamizadeh, V. Nath, Y. Tang, D. Yang, H. Roth, and D. Xu, "Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images," 2022. [Online]. Available: <https://arxiv.org/abs/2201.01266>
- [38] Z. Liu, Y. Lin, and *et al.*, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 012–10 022.
- [39] A. Dosovitskiy, L. Beyer *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [40] L. Kaiser, A. N. Gomez, and F. Chollet, "Depthwise separable convolutions for neural machine translation," *arXiv preprint arXiv:1706.03059*, 2017.
- [41] P. Gao, J. Lu, H. Li, R. Mottaghi, and A. Kembhavi, "Container: Context aggregation network," *arXiv preprint arXiv:2106.01401*, 2021.
- [42] B. Landman and *et al.*, "Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge," in *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*, vol. 5, 2015, p. 12.
- [43] M. Antonelli, A. Reinke, and *et al.*, "The medical segmentation decathlon," *Nature Communications*, vol. 13, no. 1, pp. 1512–1518, jul 2022. [Online]. Available: <https://doi.org/10.1038%2F41467-022-30695-9>
- [44] A. Wang, Q. Zhang, Y. Han, S. Megason, S. Hormoz, K. R. Mosaliganti, J. C. Lam, and V. O. Li, "A novel deep learning-based 3d cell segmentation framework for future image-based disease detection," *Scientific Reports*, vol. 12, no. 1, pp. 1–15, 2022.
- [45] J. Yang, C. Li, and J. Gao, "Focal modulation networks," *arXiv preprint arXiv:2203.11926*, 2022.
- [46] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," 2022. [Online]. Available: <https://arxiv.org/abs/2201.03545>
- [47] J. Yang, C. Li, X. Dai, L. Yuan, and J. Gao, "Focal modulation networks," 2022. [Online]. Available: <https://arxiv.org/abs/2203.11926>
- [48] Z. Wang, W. Jiang, Y. Zhu, L. Yuan, Y. Song, and W. Liu, "Dynamixer: A vision mlp architecture with dynamic mixing," 2022. [Online]. Available: <https://arxiv.org/abs/2201.12083>