

D-Net: Dynamic Large Kernel with Dynamic Feature Fusion for Volumetric Medical Image Segmentation

Jin Yang¹, Peijie Qiu¹, Yichi Zhang², Daniel S. Marcus¹, and Aristeidis Sotiras^{1,3}

¹ Department of Radiology, Washington University School of Medicine, St. Louis, MO, USA

² School of Data Science, Fudan University, Shanghai, China

³ Institute for Informatics, Data Science and Biostatistics, Washington University School of Medicine, St. Louis, MO, USA
yang.jin@wustl.edu

Abstract. Hierarchical transformers have achieved significant success in medical image segmentation due to their large receptive field and capabilities of effectively leveraging global long-range contextual information. Convolutional neural networks (CNNs) can also deliver a large receptive field by using large kernels, enabling them to achieve competitive performance with fewer model parameters. However, CNNs incorporated with large convolutional kernels remain constrained in adaptively capturing multi-scale features from organs with large variations in shape and size due to the employment of fixed-sized kernels. Additionally, they are unable to utilize global contextual information efficiently. To address these limitations, we propose Dynamic Large Kernel (DLK) and Dynamic Feature Fusion (DFF) modules. The DLK module employs multiple large kernels with varying kernel sizes and dilation rates to capture multi-scale features. Subsequently, a dynamic selection mechanism is utilized to adaptively highlight the most important spatial features based on global information. Additionally, the DFF module is proposed to adaptively fuse multi-scale local feature maps based on their global information. We integrate DLK and DFF in a hierarchical transformer architecture to develop a novel architecture, termed D-Net. D-Net is able to effectively utilize a multi-scale large receptive field and adaptively harness global contextual information. Extensive experimental results demonstrate that D-Net outperforms other state-of-the-art models in the two volumetric segmentation tasks, including abdominal multi-organ segmentation and multi-modality brain tumor segmentation. Our code is available at <https://github.com/sotiraslab/DLK>.

1 Introduction

The development of vision transformers (ViTs) has led to significant improvements in computer vision tasks [8]. The key factor in the success of ViTs is

the attention mechanism, which empowers ViT-based models with large receptive fields to utilize global contextual information across the entire input image. However, ViTs face challenges in serving as a general-purpose backbone due to the high computational complexity of self-attention in high-resolution images. To reduce the complexity of ViTs, hierarchical ViTs have been proposed [16,20,22]. They are more efficient in modeling dense features at various scales, approximating self-attention with a linear complexity. Due to their superior performance, hierarchical ViTs have recently been utilized as backbones for medical image segmentation [3,9]. However, the attention mechanism often limits (hierarchical) ViT-based models to effectively extracting local contextual information.

Another widely used backbone, Convolutional neural networks (CNNs), is advantageous in local feature extraction. However, the receptive fields of CNNs are constrained by small convolutional kernels. To enlarge their receptive fields, large convolutional kernels (LCKs) were introduced and integrated into CNN architectures [7,15,17]. Currently, LCK-based CNNs are attracting attention in medical image segmentation [2,13]. However, these networks rely on single fixed-sized large kernels for feature extraction which limits their ability to capture multi-scale features from organs with large inter-organ and inter-subject variations in shape and size. Additionally, they lack mechanisms to enhance interactions between local features and global contextual information.

To address these limitations, we propose Dynamic Large Kernel (DLK) and Dynamic Feature Fusion (DFF) modules. In DLK, we propose to use multiple varying-sized large depthwise convolutional kernels. These kernels enable the networks to capture multi-scale contextual information, effectively handling large variations in shape and size. Instead of aggregating these kernels in parallel as in Atrous Spatial Pyramid Pooling (ASPP) [5] or other parallel designs [23], we sequentially aggregate multiple large kernels to enlarge receptive fields. Subsequently, following the idea of the dynamic mechanism [6,14,24], we introduced a spatial-wise dynamic selection mechanism to adaptively select the most informative local features based on global contextual information. Besides, the DFF module is adopted to adaptively fuse multi-scale features based on global information. During fusion, a channel-wise dynamic selection mechanism is used to preserve the important feature maps, and subsequently, a spatial-wise dynamic selection mechanism is utilized to highlight important spatial regions. We integrated the proposed DLK and DFF modules into a hierarchical transformer architecture, termed D-Net, for 3D volumetric medical image segmentation. We evaluated D-Net on two segmentation tasks: abdominal multi-organ segmentation and brain tumor segmentation. The proposed model outperformed the baseline models.

Our main contributions are threefold: (i) We propose a **Dynamic Large Kernel** module for generic feature extraction. DLK employs multiple large convolutional kernels to capture multi-scale features. It subsequently leverages a dynamic selection mechanism to adaptively highlight the most important spatial features based on global contextual information. (ii) We propose a **Dynamic Feature Fusion** module for adaptive feature fusion. DFF is designed to adap-

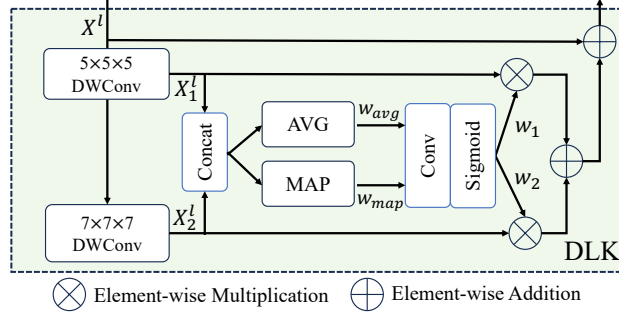


Fig. 1. The architecture of the DLK. Feature maps \mathbf{X}_1^l and \mathbf{X}_2^l are extracted by $5 \times 5 \times 5$ DWConv and $7 \times 7 \times 7$ DWConv from input features \mathbf{X}^l , respectively. The dynamic selection values w_1 and w_2 are generated to calibrate features \mathbf{X}_1^l and \mathbf{X}_2^l .

tively fuse multi-scale local features based on global information via dynamic selection mechanisms. (iii) We propose the **D-Net** for 3D volumetric medical image segmentation. D-Net is designed to adopt hierarchical transformer behaviors by incorporating DLK and DFF modules into a hierarchical ViT architecture, achieving superior segmentation accuracy with low model complexity.

2 Method

2.1 Dynamic Large Kernel (DLK)

DLK. We propose the Dynamic Large Kernel (DLK) to adaptively exploit spatial-wise contextual information via a large receptive field (Fig. 1). Specifically, multiple large depthwise kernels are used to extract multi-scale features. Additionally, instead of incorporating multiple kernels parallelly, we cascade these large kernels with growing kernel sizes and increasing dilation rates. This design has two advantages. First, contextual information is aggregated within receptive fields recursively, allowing the effective receptive fields to grow in size progressively [18]. Second, features extracted within deeper and larger receptive fields contribute more significantly to the output, enabling DLK to capture finer and more informative features. In our work, we use two depthwise convolution (DWConv) layers with large kernels: $\text{DWConv}_{(5,1)}$, featuring a $5 \times 5 \times 5$ kernel with dilation 1, and $\text{DWConv}_{(7,3)}$, featuring a $7 \times 7 \times 7$ kernel with dilation 3 for input features \mathbf{X}^l in the layer l :

$$\begin{aligned}\mathbf{X}_1^l &= \text{DWConv}_{(5,1)}(\mathbf{X}^l) \\ \mathbf{X}_2^l &= \text{DWConv}_{(7,3)}(\mathbf{X}_1^l).\end{aligned}$$

By cascading these kernels, DLK has the same effective receptive field with a $23 \times 23 \times 23$ kernel [7]. The global spatial relationship of these local features is

efficiently modeled by applying average pooling (AVP) and maximum pooling (MAP) along channels from concatenated features $[\mathbf{X}_1^l; \mathbf{X}_2^l]$

$$\begin{aligned} w_{avg} &= \text{AVP}([\mathbf{X}_1^l; \mathbf{X}_2^l]) \\ w_{map} &= \text{MAP}([\mathbf{X}_1^l; \mathbf{X}_2^l]). \end{aligned}$$

Then a $7 \times 7 \times 7$ convolution layer (Conv₇) is used to allow such information to interact among different spatial descriptors and a Sigmoid activation function is used to obtain dynamic selection values w_1, w_2 :

$$[w_1; w_2] = \text{Sigmoid}(\text{Conv}_7([w_{avg}; w_{map}])).$$

Features from different large kernels are adaptively selected by utilizing these selection values to calibrate them. Finally, a residual connection is applied as

$$\mathbf{X}^l = ((w_1 \otimes \mathbf{X}_1^l) \oplus (w_2 \otimes \mathbf{X}_2^l)) + \mathbf{X}^l.$$

DLK module. The DLK module is built by integrating DLK into two linear layers ($1 \times 1 \times 1$ convolution layers; Conv₁) with a GELU activation in between. A residual connection is also applied. Accordingly, the output of the l -th layer in a DLK module can be computed as

$$\begin{aligned} \mathbf{X}^l &= \text{Conv}_1(\mathbf{X}^{l-1}) \\ \mathbf{X}^l &= \text{DLK}(\text{GELU}(\mathbf{X}^l)) \\ \hat{\mathbf{X}}^l &= \text{Conv}_1(\mathbf{X}^l) + \mathbf{X}^{l-1}. \end{aligned}$$

DLK block. To leverage the scaling capabilities of hierarchical ViTs, the DLK block is constructed by replacing the multi-head self-attention in a standard hierarchical ViT with the proposed DLK module. The yielded DLK block consists of a DLK module and an MLP module. Similar to hierarchical ViT blocks, a Layer Normalization (LN) layer is applied before each DLK module and MLP module, and a residual connection is applied after each module. Thus, two consecutive DLK blocks in the l -th and $(l+1)$ -th layer can be computed as

$$\begin{aligned} \hat{\mathbf{X}}^l &= \text{DLK}(\text{LN}(\mathbf{X}^{l-1})) + \mathbf{X}^{l-1} \\ \mathbf{X}^l &= \text{MLP}(\text{LN}(\hat{\mathbf{X}}^l)) + \hat{\mathbf{X}}^l \\ \hat{\mathbf{X}}^{l+1} &= \text{DLK}(\text{LN}(\mathbf{X}^l)) + \mathbf{X}^l \\ \mathbf{X}^{l+1} &= \text{MLP}(\text{LN}(\hat{\mathbf{X}}^{l+1})) + \hat{\mathbf{X}}^{l+1}. \end{aligned}$$

2.2 Dynamic Feature Fusion (DFF)

We propose a Dynamic Feature Fusion (DFF) module to adaptively fuse multi-scale local features based on global information (Fig. 2). It is achieved by dynamically selecting important features based on their global information during

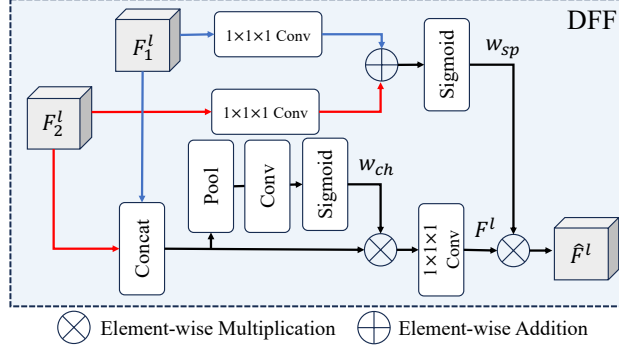


Fig. 2. The architecture of the DFF module. The global channel information w_{ch} is extracted from feature maps F_1^l and F_2^l . These feature maps are calibrated and only informative features are selected by a convolution layer to generate features F^l . In another path, the global spatial information w_{sp} is extracted from F_1^l and F_2^l , and is used to recalibrate features F^l to generate the adaptively fused features \hat{F}^l .

fusion. Specifically, feature maps F_1^l and F_2^l are concatenated along the channel. To ensure the following block can adopt fused features, a channel reduction mechanism is required to reduce the number of channels to the original one. Instead of simply using a $1 \times 1 \times 1$ convolution, channel reduction in DFF is guided by global channel information w_{ch} . This information is extracted to describe the importance of features by cascading an average pooling (AVGPool), a convolution layer (Conv₁), and a Sigmoid activation.

$$w_{ch} = \text{Sigmoid}(\text{Conv}_1(\text{AVGPool}([F_1^l; F_2^l]))).$$

Fused features are calibrated by the global channel information. Subsequently, a $1 \times 1 \times 1$ convolutional layer (Conv₁) is utilized to select feature maps based on their importance. This channel information will guide the convolution layer to preserve the important features while dropping less informative ones.

$$F^l = \text{Conv}_1(w_{ch} \otimes [F_1^l; F_2^l]).$$

To model the spatial-wise inter-dependencies among local feature maps, the global spatial information w_{sp} is captured by $1 \times 1 \times 1$ convolution layers (Conv₁) and a Sigmoid activation from feature maps F_1^l and F_2^l . This information is used to calibrate feature maps and facilitate the emphasis on salient spatial regions.

$$w_{sp} = \text{Sigmoid}(\text{Conv}_1(F_1^l) \oplus \text{Conv}_1(F_2^l))$$

$$\hat{F}^l = w_{sp} \otimes F^l.$$

2.3 D-Net Architecture

The overall architecture of D-Net consists of an encoder, a bottleneck, a decoder, and a salience layer (Fig. 3). The salience layer is used to extract salient spatial

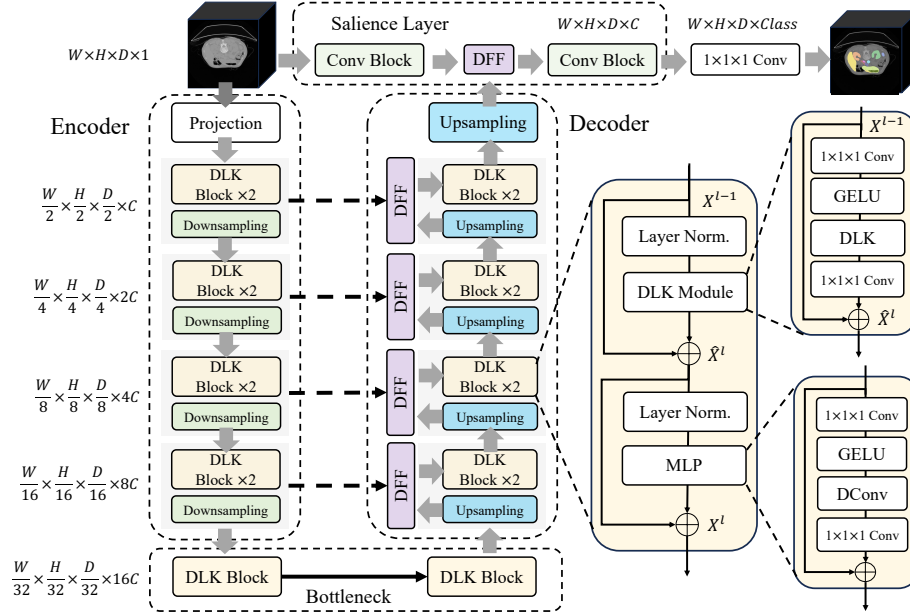


Fig. 3. The architecture of the D-Net. D-Net consists of an encoder, a bottleneck, a decoder, and a saliency layer. Two consecutive DLK blocks are used in each stage for feature extraction. Each DLK block consists of a DLK module and an MLP module.

features from original images, and the encoder-decoder architecture is responsible for learning hierarchical feature representations.

Encoder. Instead of flattening the patches and projecting them with linear layers, we utilize a large $7 \times 7 \times 7$ convolution with a stride of 2 to partition the image into feature embeddings with size $\frac{H}{2} \times \frac{W}{2} \times \frac{D}{2}$. These feature embeddings are then projected to C -dimensional vectors ($C = 48$). At each stage, two consecutive DLK blocks are combined to extract contextual information. To exchange the information across channels in the downsampling block, we use a convolution layer with a kernel size of $2 \times 2 \times 2$ and a stride of 2 to downscale the feature map and increase the number of channels by a factor of 2. The dimensions of the output feature maps at each stage are $\frac{H}{4} \times \frac{W}{4} \times \frac{D}{4} \times 2C$, $\frac{H}{8} \times \frac{W}{8} \times \frac{D}{8} \times 4C$, $\frac{H}{16} \times \frac{W}{16} \times \frac{D}{16} \times 8C$, and $\frac{H}{32} \times \frac{W}{32} \times \frac{D}{32} \times 16C$, respectively.

Bottleneck. Two consecutive DLK blocks are used for the bottleneck. The dimensions of both input and output features are $\frac{H}{32} \times \frac{W}{32} \times \frac{D}{32} \times 16C$.

Decoder. At each stage, a $2 \times 2 \times 2$ transposed convolution with a stride of 2 is used to upscale the feature map and decrease the number of channels by both

a factor of 2. These upsampled features are then fused with the features from the encoder via skip connections within DFF modules. Two consecutive DLK blocks are then used. The dimensions of the output feature maps at each stage are $\frac{H}{16} \times \frac{W}{16} \times \frac{D}{16} \times 8C$, $\frac{H}{8} \times \frac{W}{8} \times \frac{D}{8} \times 4C$, $\frac{H}{4} \times \frac{W}{4} \times \frac{D}{4} \times 2C$, and $\frac{H}{2} \times \frac{W}{2} \times \frac{D}{2} \times C$, respectively. Lastly, a transposed convolution layer is used to upsample the feature maps to a dimension of $H \times W \times D \times C$.

Salience layer. A convolution block, which consists of two consecutive $3 \times 3 \times 3$ convolutional layers, is used to generate features with the dimension $H \times W \times D \times C$ from the input image. These features are fused with features from the decoder within a DFF module. Another convolution block is then used to capture finer features. Lastly, a $1 \times 1 \times 1$ convolutional layer is used to produce the voxel-wise segmentation prediction.

3 Experiments and results

Datasets. We conducted experiments on two publicly available datasets. The first is the MICCAI 2022 AMOS Challenge dataset (AMOS 2022) [12]. It consists of 300 multi-contrast abdominal CT images with 15 anatomical organs manually annotated for abdominal multi-organ segmentation. 3D volumes were pre-processed and augmented to volumetric patches with a dimension of $96 \times 96 \times 96$ by the pipeline implemented by MONAI⁴. The second one is the Medical Segmentation Decathlon (MSD) Brain Tumours Challenge dataset [1]. It consists of 484 subjects, each with four 3D MRI modalities (FLAIR, T1w, T1gd, T2w) and three foreground annotations: Edema (ED), Enhancing Tumor (ET), and Non-Enhancing Tumor (NET). Data was preprocessed to volumetric patches with a dimension of $128 \times 128 \times 128$ by nnUNet pipeline [11].

Implementation details. The D-Net is implemented using PyTorch⁵. A combination of dice loss and cross-entropy loss was used as the loss function. In abdominal multi-organ segmentation, An AdamW was used as the optimizer. The initial learning rate was set to 0.0001, and a learning rate decay strategy (ReduceLROnPlateau) was applied. For brain tumor segmentation, we followed protocols in nnUNet [11]. The SGD was used as the optimizer. The initial learning rate was set to 0.001 and was decayed with a poly learning rate scheduler. For a fair comparison, all experiments are implemented with the same setups and implemented by us.

Main results. We compared the performance of D-Net with recent state-of-the-art segmentation models, including 3D U-Net (nnUNet) [11,19], TransUNet [4],

⁴ <https://monai.io/>

⁵ <http://pytorch.org/>

Table 1. Comparison of segmentation performance and model complexity among D-Net, DLK-Net, and other models on the 2022 AMOS multi-organ segmentation task (**Bold** represents the best results, and an underline represents the second best results).

Tasks	3D UNet [†]	TransBTS	UNETR	nnFormer	UX-Net	DLK-Net	D-Net
Spleen	95.86	95.51	95.14	89.44	<u>96.65</u>	96.37	96.74
R. kidney	95.89	95.28	95.46	88.12	96.22	<u>96.25</u>	96.31
L. kidney	96.08	95.22	94.72	86.89	96.03	<u>96.09</u>	96.17
Gall bladder	83.21	82.33	76.90	69.20	82.60	<u>83.89</u>	84.08
Esophagus	81.80	80.00	77.53	58.92	80.98	<u>82.19</u>	82.30
Liver	97.27	96.79	96.64	94.14	97.03	<u>97.31</u>	97.59
Stomach	87.96	88.18	85.27	75.36	87.78	<u>88.67</u>	93.51
Arota	93.78	93.38	92.76	85.77	93.99	<u>94.41</u>	95.19
Postcava	88.75	88.66	85.76	75.03	88.92	<u>89.85</u>	91.28
Pancreas	84.28	82.11	80.65	66.74	84.12	<u>85.99</u>	86.16
R. adrenal gland	76.02	72.50	72.64	55.41	74.50	<u>78.16</u>	79.60
L. adrenal gland	74.43	70.74	68.30	47.91	73.17	<u>77.33</u>	78.49
Duodenum	77.90	76.10	69.76	53.85	77.68	<u>80.46</u>	81.12
Bladder	90.23	88.96	85.35	74.07	<u>91.14</u>	90.73	91.40
Prostate	83.87	80.69	81.22	55.51	79.75	<u>84.94</u>	85.23
Average \uparrow	87.16	85.76	83.87	71.76	86.70	<u>88.18</u>	89.01 *
Params [‡] \downarrow	107.71M	31.58M	92.78M	149.33M	53.01M	29.17M	<u>29.96M</u>
FLOPs [‡] \downarrow	1046.39G	110.71G	<u>82.73G</u>	284.28G	632.33G	47.14G	236.90G

[†]: we implemented a standard 6-layer 3D UNet here, not nnUNet. However, our 3D UNet achieved a similar accuracy as the one reported for nnUNet [13] (87.16 vs. 87.8).

*: $p < 0.01$ with Wilcoxon signed-rank test between D-Net and other baselines.

[‡]: the number of parameters and FLOPs were calculated by the python package ptflops.

Table 2. Comparison of segmentation performance among D-Net, DLK-Net, and other models on the MSD multi-modality brain tumor segmentation task (**Bold** represents the best results, and an underline represents the second-best results).

Tasks	3D UNet [†]	TransUNet	TransBTS	UNETR	nnFormer	DLK-Net	D-Net
ET	78.54	74.98	78.44	78.45	<u>79.53</u>	78.68	80.52
ED	82.26	78.70	79.58	81.17	83.05	<u>83.29</u>	83.86
NET	<u>61.75</u>	60.41	60.40	59.75	60.53	61.73	62.72
Average \uparrow	74.18	71.36	72.81	73.12	74.37	<u>74.57</u>	75.70 *

[†]: we implemented 3D nnUNet here.

*: $p < 0.01$ with Wilcoxon signed-rank test between D-Net and other baselines.

TransBTS [21], UNETR [10], nnFormer [25], and 3D UX-Net [13] on two segmentation tasks. Table 1 shows the performance comparison on the AMOS abdominal multi-organ segmentation task. D-Net achieved the best overall performance with comparably fewer FLOPs and the lowest number of parameters. Additionally, D-Net showed significant improvement in Dice score across all organ-specific segmentation tasks. Table 2 shows the results for the MSD brain tumor segmentation task. D-Net demonstrated superior performance across all segmentation tasks compared to other segmentation methods.

Ablation study. For the ablation study, we deconstructed a D-Net to a DLK-Net by removing the Saliency layer from the D-Net and replacing each DFF module with a concatenation layer followed by a $1 \times 1 \times 1$ convolutional layer. Compared with other baselines, DLK-Net demonstrated higher segmentation accuracy in both segmentation tasks, while at the same time having the lowest model complexity (Table 1 and Table 2).

4 Conclusion

We introduced D-Net for volumetric medical image segmentation by incorporating a Dynamic Large Kernel module and a Dynamic Feature Fusion module into a hierarchical transformer architecture. The Dynamic Large Kernel block was adopted as the basic block for generic multi-scale local feature extraction and adaptive global spatial information utilization. Furthermore, the Dynamic Feature Fusion module was proposed for adaptive feature fusion. D-Net performed better than current popular baselines on two segmentation tasks: abdominal multi-organ segmentation and brain tumor segmentation. We believe that D-Net has the potential to achieve promising segmentation performance on various medical image segmentation tasks.

Acknowledgements Computations were performed using the facilities of the Washington University Center for High Performance Computing (CHPC), which was partially funded by National Institutes of Health (NIH) grants S10OD025200, 1S10RR022984-01A1, and 1S10OD018091-01.

References

1. Antonelli, M., Reinke, A., Bakas, S., Farahani, K., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B., Ronneberger, O., Summers, R.M., et al.: The medical segmentation decathlon. *Nature communications* **13**(1), 4128 (2022)
2. Azad, R., Niggemeier, L., Hüttemann, M., Kazerouni, A., Aghdam, E.K., Velichko, Y., Bagci, U., Merhof, D.: Beyond self-attention: Deformable large kernel attention for medical image segmentation. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 1287–1297 (2024)
3. Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-unet: Unet-like pure transformer for medical image segmentation. In: *European conference on computer vision*. pp. 205–218. Springer (2022)
4. Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A.L., Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306* (2021)
5. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* **40**(4), 834–848 (2017)
6. Chen, Y., Dai, X., Liu, M., Chen, D., Yuan, L., Liu, Z.: Dynamic convolution: Attention over convolution kernels. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 11030–11039 (2020)

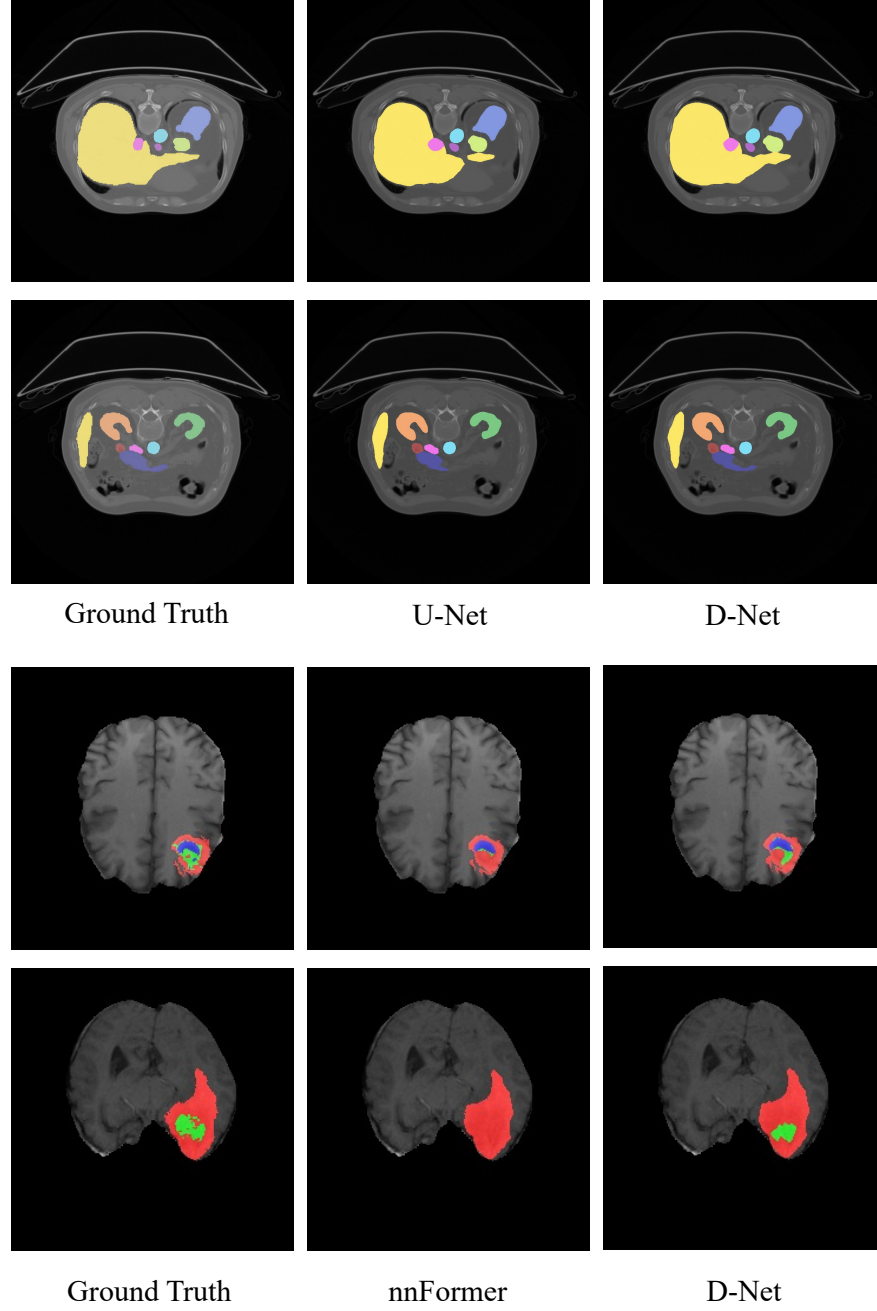


Fig. 4. Qualitative representations of multi-organ segmentation on 2022 AMOS and brain tumor segmentation on MSD BraTS dataset. D-Net shows better segmentation quality than U-Net and nnFormer.

7. Ding, X., Zhang, X., Han, J., Ding, G.: Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11963–11975 (2022)
8. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
9. Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth, H.R., Xu, D.: Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In: International MICCAI Brainlesion Workshop. pp. 272–284. Springer (2021)
10. Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., Roth, H.R., Xu, D.: Unetr: Transformers for 3d medical image segmentation. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 574–584 (2022)
11. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* **18**(2), 203–211 (2021)
12. Ji, Y., Bai, H., Ge, C., Yang, J., Zhu, Y., Zhang, R., Li, Z., Zhanng, L., Ma, W., Wan, X., et al.: Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *Advances in Neural Information Processing Systems* **35**, 36722–36732 (2022)
13. Lee, H.H., Bao, S., Huo, Y., Landman, B.A.: 3d ux-net: A large kernel volumetric convnet modernizing hierarchical transformer for medical image segmentation. arXiv preprint arXiv:2209.15076 (2022)
14. Li, X., Wang, W., Hu, X., Yang, J.: Selective kernel networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 510–519 (2019)
15. Li, Y., Hou, Q., Zheng, Z., Cheng, M.M., Yang, J., Li, X.: Large selective kernel network for remote sensing object detection. arXiv preprint arXiv:2303.09030 (2023)
16. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10012–10022 (2021)
17. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11976–11986 (2022)
18. Luo, W., Li, Y., Urtasun, R., Zemel, R.: Understanding the effective receptive field in deep convolutional neural networks. *Advances in neural information processing systems* **29** (2016)
19. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015)
20. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 568–578 (2021)
21. Wang, W., Chen, C., Ding, M., Yu, H., Zha, S., Li, J.: Transbts: Multimodal brain tumor segmentation using transformer. In: Medical Image Computing and Com-

- puter Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24. pp. 109–119. Springer (2021)
22. Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., Zhang, L.: Cvt: Introducing convolutions to vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 22–31 (2021)
 23. Yang, J., Marcus, D.S., Sotiras, A.: Abdominal ct pancreas segmentation using multi-scale convolution with aggregated transformations. In: Medical Imaging 2023: Computer-Aided Diagnosis. vol. 12465, pp. 418–426. SPIE (2023)
 24. Yang, J., Marcus, D.S., Sotiras, A.: Dynamic u-net: Adaptively calibrate features for abdominal multi-organ segmentation. arXiv preprint arXiv:2403.07303 (2024)
 25. Zhou, H.Y., Guo, J., Zhang, Y., Yu, L., Wang, L., Yu, Y.: nnformer: Interleaved transformer for volumetric segmentation. arXiv preprint arXiv:2109.03201 (2021)