# SEGMENTATION OF 3D MEDICAL IMAGES USING TRANSFORMERS

**V. Narasimha[1], Neela Raju[2], Kandikatla Divya[3], Sathvik Sahi Poturaju[4]**

**Abstract**

In recent times, neural networks using convolutions have introduced the narcotic methods of segmenting a 3D image, which depend on representation of huge features and accomplish decent execution. This is because convolutional neural networks have limited regions in sensory periphery which cannot absolutely replicate the dependencies which are from long ranges present in the image. Newly, Transformer took advantage of worldwide reliant using some automatic-checking techniques and learning of different interpretations which are high. Some experiments were according to transformers, yet current transformers experience severe arithmetical along with memory problems and are unable to fully utilize powerful properties for segmentation of 3D medical pictures.As in the paper, we focus on the parallel combination of various resolution channels, and propose a unique system called High Definition Transformer Based Network using a efficient transformer block that has an adequate depicting attribute even at high attribute resolution. When a 3Dimensional picture is given, first encoder uses neural networks to takeout object manifestations that collect regional data, and intelligently constructs a variety of feature maps for tokens, which are fed in parallel to each Transformer channel, learns global information and iteratively information is shared. Unfortunately channels, the presented standard transformer-based framework requires a substantial computing, so we start a significant and efficient transformer that provide improved performance with fewer variables. The proposed Transformer based high resolution network is evaluated on a brain tumor task dataset.

[1]Assistant Professor , Computer Science and Engineering , CMR College of Engineering and Technology, Hyderabad, Telangana India.
[2,3,4]Student, Computer Science and Engineering, CMR College of Engineering and Technology, Hyderabad, Telangana India.

Email: [1]chinna.narasimha63@gmail.com , [2]rajuneela69@gmail.com

## 1.    Introduction

Analysis of medical images provides one approach for distribution of new experiment results in the field of organic and research of medical figures. It takes a crucial part during the process of image segmentation. Image segmentation is used to segment a medical image's subjects of interest and is regarded as an important process in diagnostic scanning. It is employed to calculate the shape and size of illness quantification, the region of interest (ROI), and computerized intervention and planning for treatment. However, Medical picture segmentation is a semi-automated or automatic procedure that promises to be faster, more reliable, and less expensive than the human method, which is exceedingly labor- and time-intensive.

U-Net is semantic segmentation architecture. It consists of a narrowing and widening pattern. The shrink path adheres to the standard convolutional network design. The architecture in a U shape consists of a special Encoder-decoder design: The encoder shrinks the dimensional information and adds channels in each layer. The decoder increases room dimming and reduces the number of channels. The tensor emitted in the decoder is usually called the bottleneck. Finally, the spatial blurs are fed back to make a prediction for every input image pixel.The current basic foundation for the majority of approaches is U-Net, which comprises an encoder-decoder network and pass-through links for detailed enhancement.

In U-Net planning, we use an encoder to capture the attribute representation, which is the most powerful module,which is demonstrated   in Figure 1. Convolution neural approaches have claimed performance, although they impose some basic restrictions because of the finite reception field. kernels for a convolution in CNNs, every thing focuses on only a few pixels of the entire picture and it makes the network work to acquire regional designs instead of general knowledge.
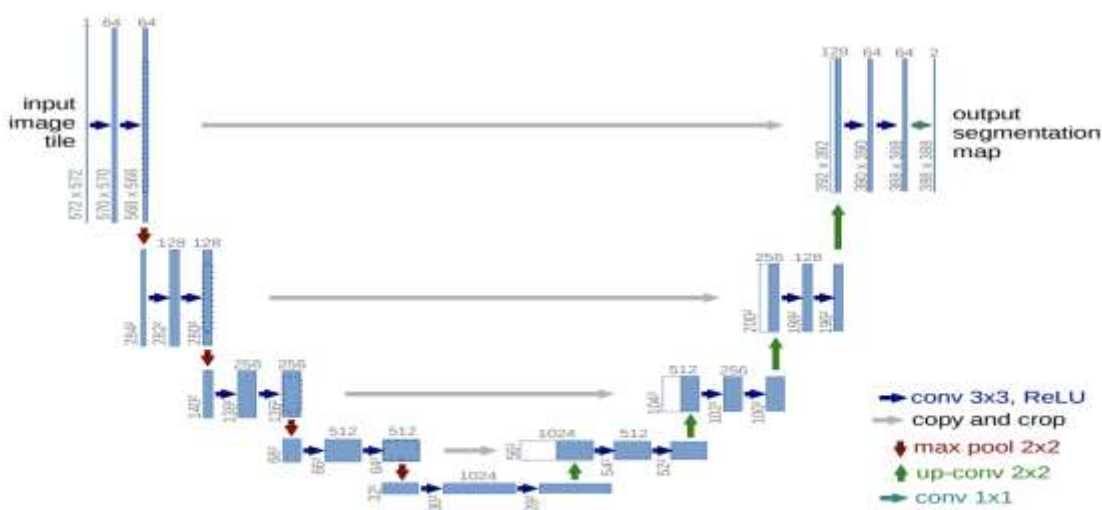


**Fig. 1.** U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

Transformers are used in order to get around this restriction.These are neither only effective for learning about global knowledge, nor show better replicability to subsequent tasks during large-scale pertaining. Today, an attractive part of the computer vision community is the development of Transformers. Top achieved excellent performance in Image Net, requiring pre-training on a large external Vision Transformer (ViT) dataset. Information distillation method is used to train the transformer. ViT a identification tasks require a massive set of training data performed using transformer-based methods.

Some recent approaches to reduce data dependency propose a simultaneous model consisting of a CNN and a transformer. Image features are captured using CNN further   Transformer encoder-decoder construction is utilised to produce the output of the last set of forecasting directly. Segmented transformer (SETR), created by Zheng et al., uses a transformer to encrypt a sequence of picture updates.TransUNet includes a transformer-based encoder that processes image update packets and a decoder using convolutions with bypass interfaces for final segmentation output, proposed by Chen et al. TransUNet uses CNN to extract all image features. The TransFuse architecture, a parallel-style two-pronged design that combines CNNs and transformers, was introduced by Zhang et al. They

can be considered preliminary research based on the transformer, of one resolution they apply functions. Transformer lacks detailed localization information, so it produces crude results because it mainly focuses on learning global information. Mehta et al., Wang et al. used the deep and efficient transform (EffTrans) with the Delight transform, which became efficient and more suitable for image analysis tasks. A Transformer block that uses group linear transformations (GLT) and therefore reduces the

computational problems of the Transformer.

## Background Work

Brain tumours are an unusual development of cells in human brain tissue. As the tumor gradually grows, it presses on the affected nerves, causing a number of negative and even life-threatening symptoms. Brain tumors are medically divided into benign and malignant tumors, the latter of which are further classified into primary and metastatic brain tumors. We process brain images using image processing techniques. Image pre-processing includes various steps. Taking a photo is as easy as taking a picture that is already in digital format. Typically, the image acquisition step includes an image scaling process. Applying image filters to the image is a part ofpre-processing of images. Several steps are used in image pre-processing.

First, picture filters can be used to highlight certain features or other features. A filter is used to change or improve a picture.Picture filters can be used to improve image edges and decrease noise in images. Noise from speckles and salt-and-pepper noise are two different forms of noise that may be seen in the picture. Dot noise is the noise that develops during imaging, whereas salt-and-pepper noise (which refers to infrequent white and black pixels) is brought on by abrupt changes in the picture signal. Image edge enhancement can help the model identify image features. Machine learning models' accuracy can be increased by performing an image pre-processing step.

Here we use a median filter to blur the image, which removes the noise in the image. In order to do this, a $n \times n$ core's pixel values must be averaged. The mean is then used to replace the middle element's pixel intensity. It smoothes the image's edges and eliminates image noise. One of the most straightforward non-linear filters, it is a function for a local average. The average of all values in the immediate vicinity is used to replace each pixel's value. The smoothed picture, g(x,y), may be produced using the formula, with S as the(x, y) region, n is the quantity of pixels in S.Assume that f(i,j) is a noisy image

$$g(x, y) = 1/n \, (i, j) \, \epsilon \, S \, \Sigma \, f(i, j) \quad (1)$$

A digital image is segmented by breaking it up into several groups, or segments, which aids reduce the image's complexity and simplify the image's

For our framework, we have to combine different functions; it's computational and the expense of memory still makes medical image analysis unaffordable.In order to lower the computational expense of learning various resolution attributes, EffTrans employed a Spatial Attention Reduction (SRA) layer. For medical 3D segmentation, the suggested EffTrans TransHRNet can handle multidimensional features well and produce more varied, detailed representations. Our method is better for extensive experiments.

additional processing or analysis. Thresholding is a technique that creates a binary image from a given grayscale image by dividing it into two regions based on a threshold value. Therefore, pixels with intensity values greater than the mentioned threshold are considered white or 1 in the output image and others are black or 0. In other words, if we have a threshold T, the segmented image g(x,y) is calculated as follows: Thus, the printed segmented image has only two pixel classes - one with a value of 1 and the other with a value of 0.

$$g(x, y) = 1 \quad if \ f(x, y) > T \quad (2)$$
$$g(x, y) = 0 \quad if \ f(x, y) \leq T \quad (3)$$

As a pixel-based image segmentation technique that uses initial seed points, region growth is likewise categorised under this category. In order to determine whether to expand the region, this segmentation technique looks at the pixels surrounding the original origin points. As with typical data clustering techniques, the procedure is repeated. Image expansion, which is the act or process of making or increasing a size or quantity, is the following procedure. The process of adding an image involves producing fresh images for our deep learning model to learn from. We don't need to manually gather these new images because they are generated utilising the training images that already exist.

The use of CNN for image segmentation has additional drawbacks. Some methods classify images based on locations, but they are unable to encode the position and orientation of the objects in the image. If the images have any tilt or rotation, CNNs frequently have trouble classifying them. While the human visual system sees an image with some noise as the same image with noise, a CNN recognises it as something entirely different, demonstrating that CNNs employ information that is quite distinct from that used by the regular visual system. CNN is essentially a mental model that tracks the direction of the current object and numerous picture properties but lacks coordinate frames, a crucial aspect of human vision.

## 2. Results & Discussions

In recent times, a new technology called as Vision Transformers has challenged Convolutional neural

networks (CNNs), the most advanced technology available today for a number of applications for computer vision that identify images.By way of ViT models outperform the relation of computing efficiency and veracity by approximately four times. Natural language processing has adopted transformer models as the de facto standard (NLP). When used in computer vision applications like semantic picture segmentation, object identification, and image classification, vision transformers perform exceptionally well.

They perform better than convolutional neural networks (CNN) while using a lot less computing power during pre-training. During training on smaller datasets, the inductive bias of Vision Transformer is frequently weaker than that of convolutional neural networks (CNN), which increases the need for model regularisation or data augmentation. The ViT is an image of a transformer whose primary function was to perform chores involving text.The ViT model directly predicts the class labels for the input picture by encoding it as a collection of image patches, much to the collection of word embeddings used when employing transformers to translate text. ViT outperforms a comparable state-of-the-art CNN while using four times less CPU power when there is enough training data.
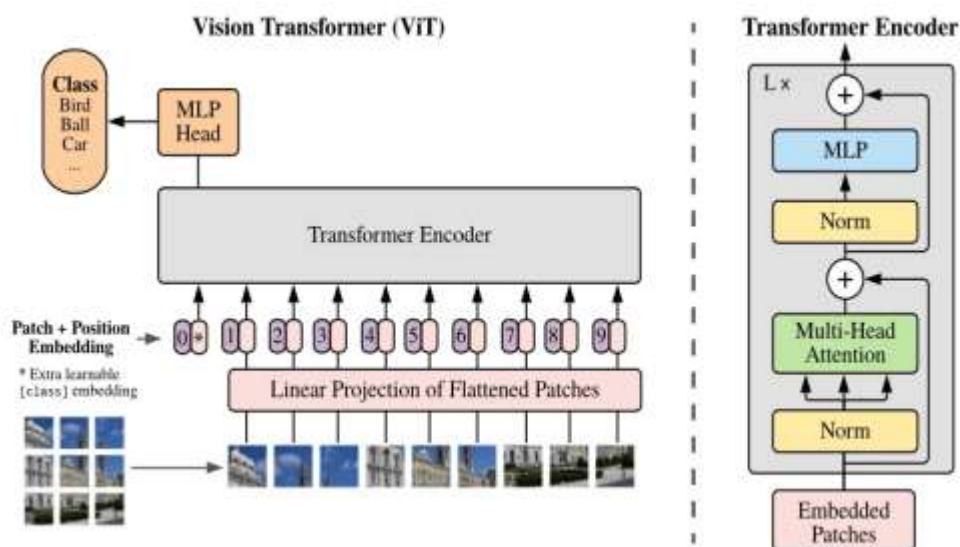


Figure 1: ViT Architecture

The location construction is included as input to the transform encoder by the visual transformer, which also separates the image into chunks of a defined size and correctly embeds them. Moreover, ViT models outperform CNNs by approximately four times in relation to correctness , computation streak. Moreover, ViT models are employed in multi-model activities including visual argumentation, visual discussion, and visual grounding as well as generative modelling tasks. The model can autonomously learn picture structure since images are provided as sequences and image class labels are anticipated.
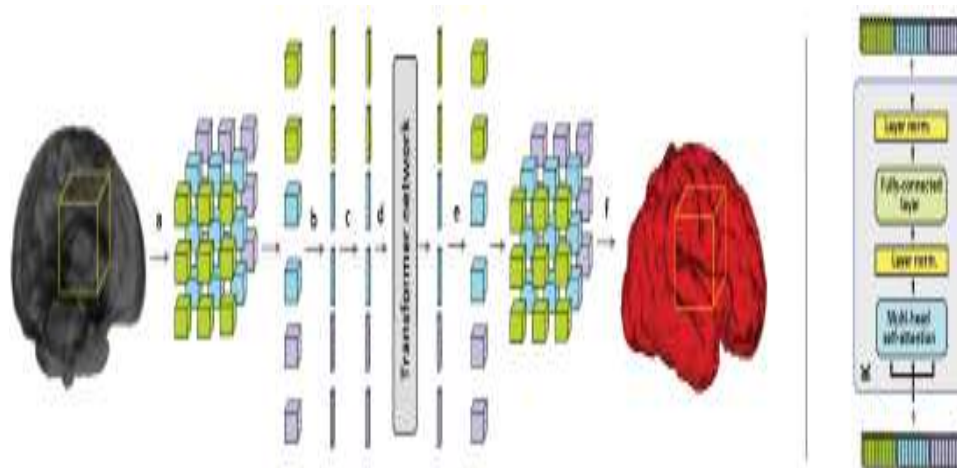


Figure 2: Image segmentation by using transformer network

The deep learning model Transformer employs a self-awareness technique that places a different emphasis on the relevance of each incoming data element. Computer Vision as well as Natural Language Processing are the first domain to use it. Recurrent neural networks (RNNs), a type of converter, are designed primarily to handle sequential input data, such language useful in jobs like translation and text summarization. Transformers process each input concurrently, in contrast to RNNs. The attention mechanism allows for context to be present at any point in the input stream.

| Practice | Differential Scanning calorimeter (%) | Huntington's Disease95 (mm) | Thoracic Aorta | GB(GallBladder) | LK(Left Kidney) | (RK)Right Kidney | Hepato(Liver) | Pancreas | Spleen | Stomach |
|---|---|---|---|---|---|---|---|---|---|---|
| V-Net | 69.78 | – | 74.67 | 52.77 | 76.18 | 81.57 | 88.48 | 41.05 | 82.96 | 57.89 |
| DARR | 70.67 | – | 75.57 | 54.27 | 73.21 | 72.42 | 95.80 | 55.81 | 88.05 | 43.69 |
| R50 ViT | 72.72 | 34.78 | 75.93 | 57.23 | 75.81 | 71.09 | 90.49 | 47.25 | 83.43 | 75.38 |
| R50 AttnUNet | 77.56 | 38.54 | 65.27 | 65.19 | 82.15 | 76.17 | 92.45 | 50.34 | 77.25 | 74.65 |
| Att-UNet | 76.75 | 37.01 | 88.59 | 69.01 | 78.00 | 70.98 | 92.96 | 57.90 | 87.41 | 75.67 |
| UNet ++ | 79.17 | 26.78 | 88.75 | 82.67 | 84.05 | 79.46 | 93.65 | 57.07 | 84.45 | 75.37 |
| TransUNet | 76.45 | 30.68 | 88.02 | 62.13 | 82.01 | 76.98 | 93.98 | 56.02 | 84.92 | 74.96 |
| UNet | 78.88 | 40.89 | 79.33 | 70.28 | 78.46 | 76.17 | 92.45 | 50.34 | 77.25 | 74.65 |
| R50 UNet | 73.65 | 37.76 | 85.23 | 64.48 | 80.19 | 72.92 | 95.85 | 68.73 | 74.01 | 75.50 |
| ViT | 63.54 | 40.51 | 88.34 | 65.13 | 82.73 | 78.34 | 95.80 | 56.35 | 86.78 | 76.26 |
| UNet3+ | 75.86 | 31.45 | 87.56 | 60.98 | 76.45 | 72.64 | 94.68 | 49.27 | 85.66 | 72.87 |
| SwinUNet | 79.98 | 21.65 | 84.97 | 66.45 | 82.91 | 78.11 | 93.69 | 57.03 | 91.66 | 75.66 |
| MT-UNet | 78.62 | 26.59 | 87.92 | 64.99 | 81.47 | 77.29 | 93.06 | 59.46 | 87.75 | 76.81 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| TransH RNet (Ours) | 91.34 | 27.23 | 89.84 | 68.90 | 85.76 | 81.27 | 95.34 | 67.99 | 89.28 | 79.67 |

Table 1. Results of experiments with the Synapse Dataset. Additionally included is the DSC (%) for each particular class. Significant values are shown in bold.
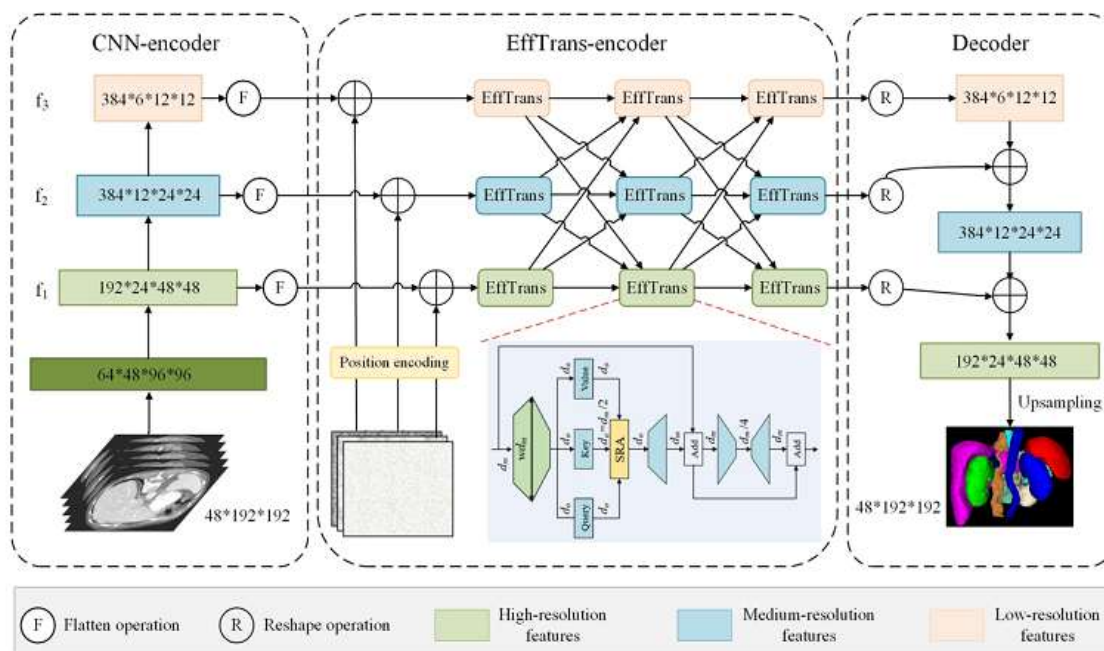


Figure 3: TransHR Network

A prominent form of neural network architecture is one using transformers. Recently, they were employed by DeepMind in its AlphaStar programme to defeat the top StarCraft player, as well as lately by OpenAI in their language models. Transformers were created to address the issues with nervous system sequence transfer, often known as machine translation. As a result, every task is transformed from an input sequence to an output sequence. This covers text-to-speech, speech recognition, etc.
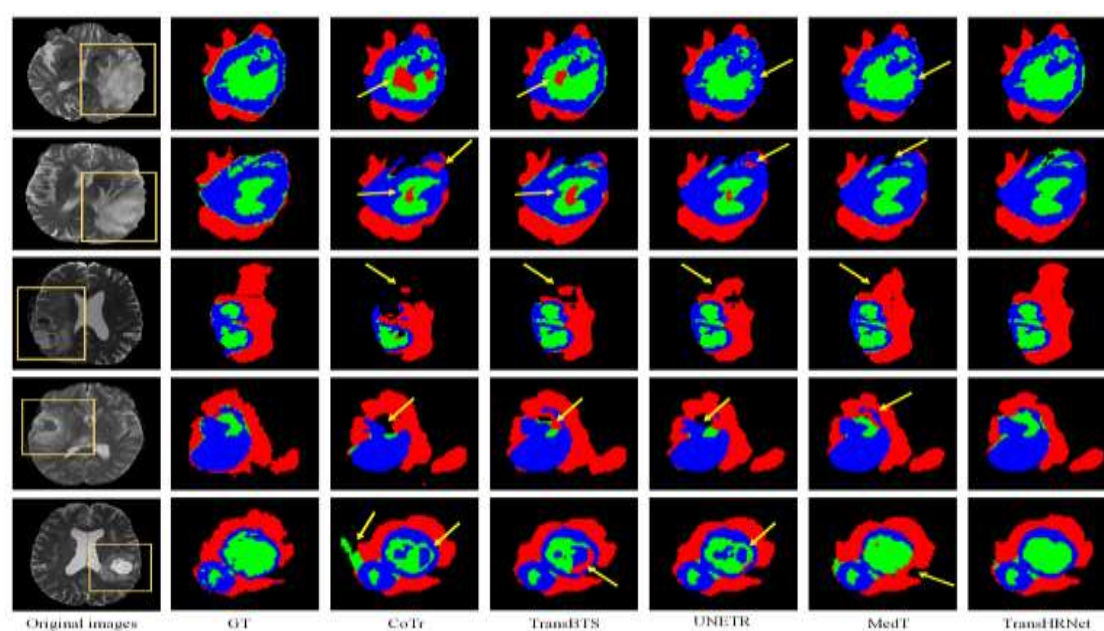


Figure 4: Performance evaluated for brain tumor

Table 2. Experimental findings utilising the ACDC Dataset.

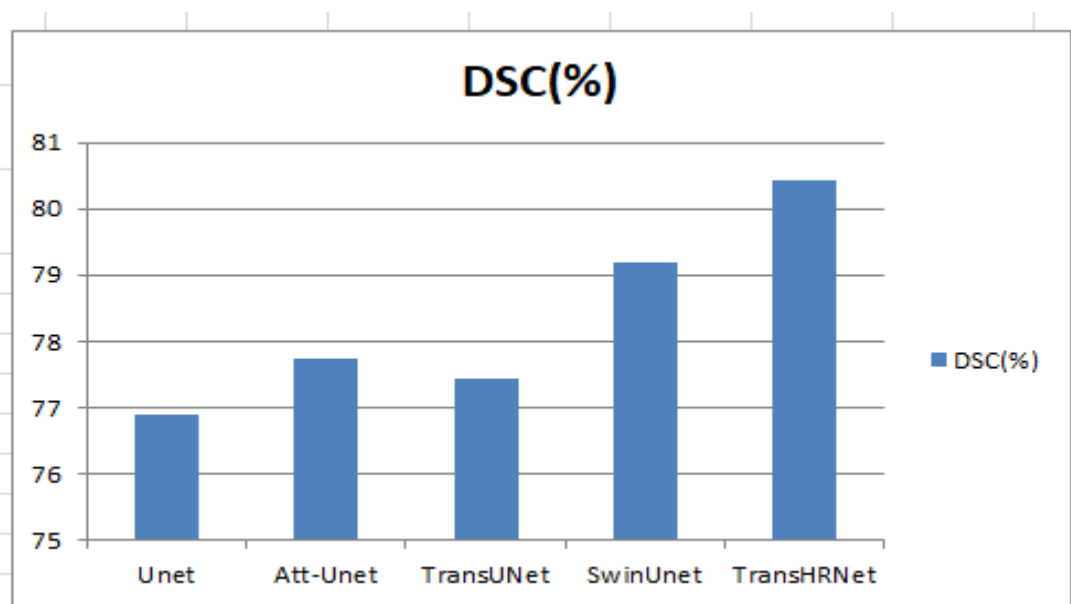| Practice | Performance comparisons DSC (%) | Right Ventricle | MYO | Left Ventricle |
|---|---|---|---|---|
| **UNet** | **88.28** | **86.08** | **86.04** | **92.72** |
| **R50 Att-UNet** | **85.90** | **86.56** | **80.34** | **92.78** |
| **CE-Net** | **85.21** | **85.68** | **83.97** | **91.98** |
| **R50 UNet** | **86.45** | **86.89** | **81.26** | **95.29** |
| **UNet + +** | **90.23** | **88.79** | **87.78** | **92.57** |
| **R50 ViT** | **86.60** | **87.77** | **80.98** | **93.98** |
| UNet3+ | 89.34 | 87.45 | 88.67 | 93.67 |
| **UNETR** | **88.61** | **85.29** | **86.52** | **94.02** |
| **SwinUNet** | **90.10** | **89.65** | **84.92** | **95.84** |
| **TransUNet** | **88.77** | **89.88** | **85.58** | **95.46** |
| **TransHRNet (Ours)** | **91.23** | **90.36** | **88.45** | **93.56** |



Figure 5: For the Synapse dataset, a comparison of the mean DSC (%) using several approaches. Comparison of the segmentation performance of classical designs using a variety of techniques.

Table 3. The Synapse dataset was used for an ablation investigation. Bold values indicate significant values.

| Method | DSC(%) | HD95(mm) |
|---|---|---|

| | | |
|---|---|---|
| TransUNet | 76.34 | 32.54 |
| TransUNet + Group Parallel Axial Attention | 79.05 | 25.34 |
| TransUNet + sMLP | 79.02 | 28.78 |
| TransHRNet (Ours) | **81.73** | **21.43** |

Table 4. Analysis of the ACDC dataset. Bold values indicate significant values.

| Method | DSC(%) | RV | Myo | LV |
|---|---|---|---|---|
| TransUNet | 89.57 | 88.86 | 84.53 | **95.93** |
| TransUNet + Group Parallel Axial Attention | 89.66 | 88.14 | 87.37 | 92.49 |
| TransUNet + sMLP | 89.66 | 89.30 | 86.08 | 93.56 |
| TransHRNet (Ours) | **91.37** | **90.56** | **88.98** | 93.68 |

## 3.    Conclusion

The circuit operation has a solid foundation in the mammalian primary visual cortex's structure and is ideally suited to the creation of effective methods for modelling and comprehending images. CNNs have recently shown to be quite efficient for a variety of computer vision issues. However, there is no reason to rule out the possibility that another model may outperform CNNs in a certain visual task.Due to the 3D nature of imagery and the scarcity of labelled images, the evaluation of medical images especially encounters special challenges. Other models might perform better in certain applications than CNNs. We developed a novel method for 3D segmentation medical images in this work. Unlike all current models that rely mostly on circulants.Our approach is based on self-adhesion between adjacent 3D regions as its fundamental building pieces. Our findings demonstrate that, on three medical picture segmentation datasets, the proposed network can beat state-of-the-art FCNs. While only some labelled training images were available, our network beat FCN both noise-free and unlabeled image colouring tasks. With regard to additional medical image analysis tasks, such as anomaly detection and classification, we anticipate that the network suggested in this research would perform well.

Even at high resolution, the block of EffTrans appropriate depiction of a feature and offers reduced parameters, better performance. Furthermore, we utilised feature fusion to concurrently merge high-resolution Transformer streams. Several tests using MSD datasets have shown that TransHRNet is more reliable than a well-built CNN system. In conclusion, brain tumour segmentation is essential to the procedure for diagnosis. Correct segmentation helps with clinical diagnosis and prolongs patient life. TransHRNet, a hybrid CNN Transformer model constructed using the Effective Transformer (EffTrans) block, was used in this study to segment brain tumours. The algorithm combines local and global features in order to accurately execute segmentation.

## 4.    References

Transformer and group parallel axial attention co-encoder for medical image segmentation Chaoqun Li, Liejun Wang , Yongming Li.

Qingsen Yan, Shengqiang Liu, Songhua Xu, Caixia Dong, Zongfang Li, Javen Qinfeng Shi, Yanning Zhang, Duwei Dai. "3D Medical image segmentation using parallel transformers" , Pattern Recognition, 2023.

"Near-Optimal Bounds for Generalised Orthogonal Procrustes Problem via Generalised Power Method",arXiv:2112.13725v1 [cs.IT] 27 Dec 2021

Davood Karimi, Haoran Dou, Ali Gholipour. "Medical Image Segmentation using Transformer Networks" , IEEE Access, 2022.

Yiyang Zhao, Jinjiang Li, Zhen Hua. "MPSH: Multiple Progressive Sampling Hybrid Model Multi-Organ Segmentation" , IEEE Journal of Translational Engineering in Health and Medicine, 2022.

Shaila S G, Shivamma D, Monica U S, Tejashree K. "Facial Expression Recognition for Compound Emotions using Mobile Net

Architecture" , 2022 International Conference on Artificial Intelligence and Data Engineering (AIDE), 2022.

"Image Processing and Capsule Networks" , Springer Science and Business Media LLC, 2021.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. 2016.

Spyridon Bakas, Mauricio Reyes, et Int, and Bjoern Menze. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge. In arXiv:1811.02629, 2018.

Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254, 2021.

Irwan Bello. Lambdanetworks: Modeling long-range interactions without attention. In International Conference on Learning Representations, 2020.

Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In European Conference on Computer Vision, pages 213–229. Springer, 2020.

Mathilde Caron, Hugo Touvron, Ishan Misra, Herve J ́egou, ́ Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. arXiv preprint arXiv:2104.14294, 2021.

Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306, 2021.

Jianxu Chen, Lin Yang, Yizhe Zhang, Mark Alber, and Danny Z Chen. Combining fully convolutional and recurrent neural networks for 3d biomedical image segmentation. In Proceedings of the 30th International Conference on Neural Information Processing Systems, pages 3044–3052, 2016.

Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 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.

Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In Proceedings of the European conference on computer vision (ECCV), pages 801–818, 2018.

Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. arXiv:1802.02611, 2018.

Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. arXiv preprint arXiv:2104.13840, 1(2):3, 2021.