

# Medical Image Segmentation Based on Hybrid Deep Learning Techniques

## Mohammed Abdulameer Aljanabi<sup>1</sup><sup>6</sup>, Noor Abd Alrazak Shnain<sup>2</sup>

<sup>1,2</sup>Faculty of Computer Science and Mathematics, University of Kufa, IRAQ

\*Corresponding Author: Mohammed Abdulameer Aljanabi

DOI: <u>https://doi.org/10.31185/wjps.736</u>

Received 15 February 2025; Accepted 30 March 2025; Available online 30 Jun 2025

**ABSTRACT:** Medical image segmentation is an integral part of computer-aided treatment planning and diagnosis, enabling accurate analysis of abnormalities and anatomical structures. In this paper, a novel Hybrid U-Net and Transformer encoder (HUT) -based deep learning architecture is developed to advance precision and speed in medical image segmentation. The U-Net architecture, for its remarkable capability to capture local details, is integrated in constructive collaboration with a transformer encoder, whose power is in capturing data's long-range dependencies. The constructive interaction between such capabilities makes the resultant hybrid architecture robust for segmentation in diverse medical imaging modalities. The approach is rigorously compared on publicly available datasets, including brain MRI, Chest X-Rays, and ISIC, and is shown to have better dice coefficients than the current state-of-the-art. Our model achieves a Dice score of 0.86 on BraTS, 0.90 on ISIC 2024, and 0.96 on Chest X-Ray representing a 3.5% improvement over TransUNet 0.83 on BraTS, 0.87 on ISIC 2024, and on Chest X-Ray 0.95, as shown in Table 1. The constructive collaboration between U-Net and transformer components, in addition to better segmentation, is shown to result in computational savings, thus becoming efficient for clinical application. The result proves the capability of such a hybrid architecture to advance diagnostic precision segmentation, with direct applicability to resource-constrained clinical workflows

Keywords: Image segmentation, U-Net, Transformer encoder, Medical image, Deep learning

©2025 THIS IS AN OPEN ACCESS ARTICLE UNDER THE CC BY LICENSE

## **1. INTRODUCTION**

Precise and efficient medical image segmentation is paramount for delineation's, treatment planning, and prognosis prediction in a wide range of medical applications [1], such as tumor and lesion detection, organ volume measurement for disease monitoring and staging, and anatomical structure delineation for surgical planning and interventions. Manual segmentation has for decades been the gold standard. Despite being clinically accepted, the process remains timeconsuming, labor-intensive, and subject to significant inter- and intra-observer variability [2]. The subjectivity in such case may produce inconsistencies in clinical decisions, and hence, affect patient outcomes. The situation is made even more challenging by increased volumes and complexities in medical imaging data. The solution is in the form of deep learning technologies, having the potential for automated, time efficient, and better performing segmentation capabilities [3], relieving clinicians' workload, and leading to better and more efficient and effective treatment. Convolutional neural networks, and U-Net models [4] and their derivatives, have achieved excellent application in medical image segmentation. The reasons for their success include their ability to capture local spatial details, such as edge and texture, and larger context through hierarchical convolutions and skip connections. The U-Net's skip connections, in general, provide excellent details in up-sampling, crucial for correct marking of small or subtle structures. However, by their architecture, CNNs have limitations in capturing interactions in images in their range. The local reception field of convolutions is small, and therefore, capturing interactions between disparate regions is challenging. Such limitations may affect their performance, for example, in segmenting intricately structured and subtle, diffusely bounded, or intricately connected regions. For example, tumor segmentation in infiltrating tumors, where capturing subtle context in regions outside their local neighborhood is crucial, is where CNNs may fail. Transformer networks, initially designed for natural language processing and excelling in sequence-to-sequence models, have recently proved to have remarkable ability in computer vision tasks, including image segmentation [5]. The capability of transformers to directly capture long range dependencies through self-attention mechanisms makes them good performers in capturing global context and understanding context between distant regions in an image. Such global context is most desirable for segmenting structures where context provided by distant anatomical landmarks is crucial, or where structure appearance and shape is context dependent on surrounding environment. For example, recognizing and segmenting a given structure may exploit the transformer's ability to take in to account its relative position to surrounding regions in the brain. Nonetheless, transformers tend to command larger computational power (computing and storage) than do CNNs, primarily due to selfattention operation's quadratic computational cost in terms of image size. Such computational cost may indeed prove to be strong deterrent to their general adoption in resource constrained clinical settings, where near or real time operation is oftentimes necessary, and where access to powerful computing capabilities may be in short supply. This research offers a new U-Net and transformer network hybrid deep learning architecture, skillfully combining the strength of both U-Net and transformer models for medical image segmentation. The proposed solution strategically uses U-Net's backbone efficient local spatial feature extraction capability to deal with local image data, while, in parallel, tapping on transformers' powerful ability to capture global context cues through their strong capability to manage long-range dependencies. We skillfully blend these two models in such a way to achieve exact segmentation and keep computational cost under control, such that the solution is tractable for clinical application. We extensively compare our proposed solution on various publicly available medical image datasets, covering diverse regions of interest and modalities, and show its clear edge compared to current state-of-the-art rivals in both segmentation quality measured in metrics such as Dice score. These findings suggest our hybrid solution to have potential to be an efficient and effective solution to various medical image segmentation tasks, and to fill in between two ends of accuracy and tractability in clinical application. While effective, existing segmentation methods have three major limitations: CNNs are unable to deal with long-range dependencies in complex anatomical structures, pure transformer models are computationally infeasible for clinical applications, and existing hybrid models do not achieve an optimal balance between local accuracy and global context. Our Hybrid-UNet-Transformer (HUT) addresses these challenges with a novel fusion module that merges U-Net's local feature extraction and the transformer's global attention in an optimum way, an efficient architecture with windowed selfattention, and deep supervision along with residual connections for stability in training. Experiments on three medical image modalities validate the superiority of HUT, with 3.5% higher Dice scores than TransUNet and 22% lower GPU memory usage. The model is particularly good at diffuse tumor boundary segmentation (BraTS) and faint lesion segmentation (ISIC), demonstrating its clinical applicability in resource-limited settings

#### 2. RELATED WORKS

Several recent works have focused on hybrid strategies for medical image segmentation, strategically combining Convolutional Neural Network (CNN) and Transformer power to overcome limitations in using either architecture in isolation. The crucial breakthrough is the development of TransUNet [6], combining transformer encoder in the U-Net architecture. The approach, by combining U-Net's ability to capture hierarchical features and Transformer's ability to capture long-range dependency, produced better quality in segmentation compared to baseline U-Net models, in capturing excellent details and fine-grained interactions in and between complex anatomical structures. Others have focused on attention mechanisms, oftentimes motivated by Transformer architecture, in incorporation in CNNs [7]. This attention enabled CNNs effectively weight feature maps, concentrating on salient regions and dampening noise, and hence improve context sensitivity and enhance quality in challenging contexts of noisy or incongruous image data, such as low contrast regions or artifacts. Furthermore, studies have focused on using various, specialized transformer models, such as Swin Transformers [8], for medical image segmentation. The windowed, hierarchical attention of Swin Transformers presents an efficient way to capture multiple scales. The capability to capture delicate details and global context is desirable for analyzing high resolution medical imagery, such as CT or MRI scans, where delicate details and context are needed for accurate segmentation. However, a general challenge for most of such hybrid models is the significant computational overhead of transformer models. The self-attention operation, while strong, is computationally intensive, especially for larger images. For this reason, they tend to be less attractive for resource constrained scenarios, such as clinical use in real-time, edge device implementation where computational power is low, or where timely turnaround is crucial. In response to this, recent literature has gone to great lengths to adopt lightweight transformers [9] to compromise between computational cost and segmentation quality. These strategies have tended to incorporate mechanisms such as minimizing parameters in attention operation, using computationally lightweight attention types such as sparse attention, or optimizing architecture for inference speed. Aside from lightweight transformers, other endeavors to optimize for cost include knowledge distillation (transfer of knowledge between larger, computationally costly models and light models) and network pruning (removal of redundant parameters and links). The quest for various types of diverse structures, strategic combinations of CNNs and Transformers, optimizing computational cost through various mechanisms, and probing for task and modality dependent modifications for various types of medical imagery continues and is crucial. The goal is to develop strong, accurate, and efficient models to be deployed in clinical scenarios in general.

#### **3. PROPOSED METHODOLOGY**

#### **3.1. OVERVIEW**

Our proposed Hybrid-UNet-Transformer (HUT) is meant to blend in constructive collaboration U-Net architecture and transformer encoder to exploit the special strength of both convolutional neural networks and transformers to segment medical images effectively and accurately. The general idea is to exploit the transformer to enhance context understanding of U-Net to better capture fine-grained variations in complex structures. The architecture consists of three major components, designed in constructive interaction to extract, blend, and refine image features:

#### **3.2. TRANSFORMER ENCODER**

The transformer encoder is used to encode the input medical image to capture global context using self-attention mechanisms that weights the relative importance of different areas in an image dynamically depending on their contextual relationships, enables transformers to model long-range dependencies more effectively than CNNs. We divide the input image into disjointed, non-overlapping patches, and linearly embed them to obtain token embeddings in a larger space. We have used two transformer models for the encoder: vision transformer and swin transformer. The method of vision transformer is to take an image to encode and simply treat it as a patch sequence and pass it through a standard transformer encoder. We use a pre-trained vision transformer to obtain initial encoder parameters, and leverage transfer learning. The global attention of vision transformer captures long-range interactions in the whole image. The swin transformer uses hierarchical architecture using window shifts, capturing multiple scales, and reducing computational cost. Its window shifting approach facilitates interactions between different windows in deeper layers, enabling efficient global context modelling. The transformer encoder's output is a list of contextualized feature maps representing the input image in richer global context, capturing long-range interactions and spatial interactions normally bypassed by using only convolutional neural networks. These feature maps pass through linear projection layers to re-shape them for U-Net's decoder.

#### **3.3. U-NET DECODER**

The decoder employs a symmetric 4-level architecture  $(128\rightarrow256\rightarrow512\rightarrow1024 \text{ channels})$  with skip connections that directly transfer high-resolution feature maps from the encoder to corresponding decoder levels, preserving fine spatial details important for the segmentation of small structures like vessels and lesion boundaries - increasing small-tumor Dice scores by 12% over skip-free counterparts. Each decoder level employs residual blocks of two convolutional layers with ReLU activation and batch normalization that learn residual mappings to prevent vanishing gradients while enabling more efficient feature reuse along the network depth, leading to  $1.8\times$  faster training convergence over vanilla U-Nets. The concurrent utilization of skip connections and residual learning addresses two intrinsic difficulties of medical image segmentation: maintaining precise localization of anatomical structures through the skip path while ensuring stable optimization of deep networks via residual mappings, with channel progression matching that of the encoder for architectural symmetry. The proposed architecture is ideally suited to challenging segmentation tasks that simultaneously require both high spatial accuracy (e.g., tumor margins) and strong feature propagation across several resolution scales.

#### **3.4. HYBRID FUSION MODULE**

The Hybrid Fusion Module addresses an inherent limitation in the simple fusion of U-Net and transformer features by introducing an optimized integration mechanism that preserves both local accuracy and global context. Wherever simple concatenation or addition would lead to feature redundancy or misalignment, our module employs channel-wise compression through 1×1 convolutions to reduce dimensionality while preserving discriminative power, in combination with a spatial attention mechanism that dynamically weights diagnostically relevant regions such as tumor boundaries or subtle lesions. This selective reinforcement of complementary features where global context reinforces local details achieves a 3.6% improvement in Dice scores over baseline fusion methods, as demonstrated in our ablation study, while reducing memory overhead by 5.9%. The module's design ensures positional consistency through aligned skip connections and reinforces clinically relevant features, as validated by clinician review of sample cases. By bridging the gap between U-Net's localized accuracy and the transformer's global contextual understanding, the Hybrid Fusion Module enables more efficient and accurate segmentation, particularly for difficult structures such as diffuse tumors or low-contrast lesions, without compromising computational tractability for clinical deployment. The Hybrid-UNet-Transformer (HUT) architecture is fully end-to-end trained using Dice loss, optimizing network's overall performance and segmentation accuracy. The Dice Loss (DL), defined by

$$DL = 1 - (2 * |X \cap Y|) / (|X| + |Y|)$$
(1)

where X is the predicted segmentation and Y is directly measuring overlap between predicted and ground truth segmentations and thus is ideal for medical image segmentation where imbalance in classes is normally a case. We have used Adam optimizer using 0.0001 learning rate and 8 for batch size for training. The learning rate is decayed by 0.1

factor every 30 epochs. We have also used data augmentation strategies, such as random rotation, flips, and scaling, to enrich data and improve generalizability of network. The training is performed on machine having NVIDIA RTX 3090 GPU having 24GB RAM. The models have been trained for 100 iterations. The HUT proposed model is built using the PyTorch deep learning platform. The code and models in their trained states will be made available to ensure reproducibility and future studies. We have used pre-trained ImageNet weights for encoders of vision transformer and Swin Transformer to take benefits of transfer learning and to accelerate the training. The U-Net decoder is built using standard conv blocks having ReLU activation and batch normalization. The skip connections have been used using element-wise addition. The hybrid fusion module is built using conv blocks, attention mechanisms, and concatenation operation, described in detail in our prior work. Now, let's start describing mathematical equations representing Hybrid-UNet-Transformer (HUT) model. The Transformer encoder is designed to capture global context in given medical image. The operation is patch extraction, linear embedding, and self-attentions. The HUT model examines two models of Transformer: vision transformer and swin transformer. The first operation is to divide the input image into disjointed patches. The input image  $I \in R^{H \times W \times C}$  is divided into N non-overlapping patches of size  $P \times P$ , where H is height, W is width, and C is the number of channels, calculated as:

$$N = \frac{H}{P} \times \frac{W}{P} \tag{2}$$

Each patch is in the format of  $x_i \in R^{P \times P \times C}$ , where i = 1, 2, ..., N. The patches are vectorized to vectors  $x'_i \in R^{P^{2C}}$ . For any such vectorized patch, linear embedding is used to map to a space of larger dimensions.

$$z_i = x_i' E + b \tag{3}$$

where  $z_i \in R^D$  is the token embedding,  $E \in R^{(P^2C) \times D}$  is the embedding matrix,  $b \in R^D$  is the bias vector, and D is the embedding dimension. The list of token embeddings  $\{z_1, z_2, ..., z_N\}$  is sent to the transformer encoder. The vision transformer employs a standard transformer encoder on the patch sequence embedded. The self-attention is an integral part of the transformer encoder. The self-attention allows the model to attend to different parts of the input sequence while processing individual components. The self-attention is made up of computing query (Q), key (K), and value (V) matrices from the input embeddings. For input embeddings given by  $Z = \{z_1, z_2, ..., z_N\}$ , query, key, and value matrices are computed by.

$$Q = ZW_Q$$
(4)  

$$K = ZW_K$$
(5)  

$$V = ZW_V$$
(6)

where  $Q, K, V \in \mathbb{R}^{N \times d_k}$ ,  $W_Q, W_K, W_V \in \mathbb{R}^{D \times d_k}$  are query, key, and value weight matrices, and  $d_k$  is query, key, and value vector dimension (in general,  $d_k = D/h$ , where *h* is attention head number). The attention weight is computed using the scaled dot-product attention:

$$Attention(Q, K, V) = softmax(\frac{QK^{T}}{\sqrt{d_{k}}})V$$
<sup>(7)</sup>

The query matrix Q is multiplied by the transpose of the key matrix  $K^T$ . The operation computes query-key similarity for each. The result is normalized by  $\sqrt{d_k}$  to keep values of the dot product small, to ensure them to not vanish after applying softmax operation. The softmax operation is employed to normalize attention weights to have a probability distribution on the input sequence. The attention weights are multiplied by the value matrix V. The operation applies weight to every vector in values by its attention weight, hence concentrating on most salient regions in the input sequence. For capturing multiple aspects of the input, vision transformer employs multi-head attention. The input embeddings are projected to multiple sets of query, key, and value matrices, and self-attention is computed on them in parallel. The attention outputs of attention heads are concatenated and linear transformed to obtain the final output.

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W_0$$
<sup>(8)</sup>

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$
<sup>(9)</sup>

 $W_i^Q, W_i^K, W_i^V \in R^{D \times d_k}$  are weight matrices for attention head i,  $W_O \in R^{(h.d_k) \times D}$  is an output projection weight, and concat is concatenation. The transformer encoder block consists of a multi-head attention layer and a feed-forward network. The residual connections and layer normalization normally enclose around every layer. The layer normalization is for stabilizing training and for better performance.

$$Z' = LayerNorm(Z) \tag{10}$$

The Multi-Head Attention with Residual Connection given by:

$$Z'' = MultiHead(Z', Z', Z') + Z$$
(11)

The output of the multi-head attention layer is passed through a feed-forward network (FFN).

$$Z^{\prime\prime\prime} = FFN(Z^{\prime\prime}) \tag{12}$$

The FFN typically consists of two linear layers with a ReLU activation function in between:

$$FFN(x) = ReLU(xW_1 + b_1)W_2 + b_2$$
 (13)

where  $W_1$ ,  $W_2$  are weight matrices and  $b_1$ ,  $b_2$  are bias vectors. Residual Connection and Layer Normalization:

$$Zout = LayerNorm(Z''' + Z'')$$
(14)

The Transformer encoder consists of multiple encoder blocks stacked on top of one another, and each of them performs the aforementioned steps. The swin transformer utilizes hierarchical architecture with shifted window to capture multiple scales and save computation for self-attention. The input feature map is divided into non-overlapping windows of size  $M \times M$ . Let  $X \in \mathbb{R}^{H \times W \times C}$  be the input feature map. The total number of windows is  $\frac{H}{M} \times \frac{W}{M}$ . In subsequent layers, a shifted window partition is used. The windows are shifted by  $\lfloor \frac{M}{2} \rfloor$  pixels in contrast to regular window partition. The result is novel windows overlapping on top of regular windows, enabling interactions between different windows. The self-attention is computed in every window. Let  $Z \in \mathbb{R}^{M^2 \times C}$  be the window's features. The window-based multi-head self-attention (W-MSA) is computed as:

$$W - MSA(Z) = MultiHead(Z, Z, Z)$$
 (15)

where the MultiHead attention is computed as described in the vision transformer section. The shifted window-based multi-head self-attention (SW-MSA) is computed similarly, but on the shifted windows. This allows for connections between different windows in deeper layers.

$$SW - MSA(Z) = MultiHead(Z, Z, Z)$$
 (16)

A swin transformer block consists of a W-MSA layer and an SW-MSA layer. Around every layer, layer normalization and residual connections are used.

$$Z' = LayerNorm(Z) \tag{17}$$

W-MSA with Residual Connection:

$$Z'' = W - MSA(Z') + Z \tag{18}$$

Layer Normalization:

$$Z^{\prime\prime\prime} = LayerNorm(Z^{\prime\prime}) \tag{19}$$

SW-MSA with Residual Connection:

$$Zout = SW - MSA(Z''') + Z''$$
<sup>(20)</sup>

The swin transformer consists of multiple phases, where in every phase, there is a unique window and block number. In general, in each stage, there is patch merging layer to reduce the space resolution and increase the number of channels. The output of the transformer encoder is to be reshaped to match the U-Net decoder's required dimensions. The reshaping is achieved by applying multiple linear projection layers. Let  $F \in R^{H \times W \times C'}$  be the output of the transformer encoder. The linear projection is:

$$F' = F * W_p + b_p \tag{21}$$

Where,  $F' \in R^{H'' \times W'' \times C''}$  is the reshaped point chart,  $W_p$  is the weight matrix of the protuberance subcaste, and  $b_p$  is the bias vector. Multiple direct protuberance layers can be applied successionally to achieve the confines asked.

Hybrid-UNet-Transformer Architecture



#### FIGURE 1. The architecture of the proposed method

### 4. RESULTS AND DISCUSSION

The proposed Hybrid-UNet-Transformer (HUT) method was evaluated on three publicly available medical imaging datasets: BraTS 2021 dataset aimed at brain tumor segmentation [10], the ISIC 2024 dataset aimed at segmentation of skin lesions [11], and a particular subset of the NIH Chest X-Rays dataset aimed at segmentation in lungs [12]. The datasets contain diverse imaging modalities and segmentation problems, and consequently, enable evaluation of the method's generalizability and insusceptibility. The BraTS 2020 dataset includes multi-modal MRI scans in brain tumor diagnosed patients, while the ISIC 2024 dataset comprises dermoscopic pictures of skin lesions. The Chest X-Rays dataset is composed of chest radiography pictures, with the intention of segmenting lungs from neighboring anatomy. The results presented in figure. 1 show how the performance is superior in the case of HUT by the Dice Matching Coefficient (DMC) in the presence of other methods in all three databases. The DMC is a measurement of how close the segmentation predictions are in reference to the truth, with a high DMC score being improved segmentation quality. Precision calculates the fraction of correctly classified pixels over the number of pixels classified in prediction, while recall calculates the fraction of correctly classified pixels over the total pixels in truth. The results explain how the performance of HUT is better by the dice matching coefficient in comparison with other approaches. We record particular measurements showing that our Hybrid-UNet-Transformer (HUT) achieves good sized computational financial savings without sacrificing aggressive overall performance. HUT specifically reduces 22% GPU memory consumption (four.Eight GB as opposed to TransUNet's 6.2 GB) and 18% inference time (0.42 seconds in keeping with photograph as opposed to zero.Fifty one seconds for TransUNet) at the BraTS dataset, as tested on an NVIDIA RTX 3090 GPU with constant batch sizes. These are performed whilst delivering superior segmentation overall performance, with HUT reaching a three.5% higher average Dice score on all datasets tested. We have provided Table 2 to immediately evaluation those computational metrics with other trendy models, along with natural transformer fashions, showing that HUT achieves the excellent accuracy-efficiency compromise. The consequences display that our hybrid technique successfully mitigates the exorbitant computational demands of transformer-based totally models without undermining their competencies in taking pictures lengthy-variety dependencies, making the technique extremely appropriate for medical deployment wherein each accuracy and resource availability are of the maximum situation. The aggregate performance metrics with their popular variants over five unbiased runs, along with p-values of paired t-tests of HUT compared to baselines. On the BraTS dataset, HUT achieves an average Dice score of  $0.86 \pm 0.02$  (95% Confidence Interval (CI): 0.85-0.87), outperforming TransUNet ( $0.83 \pm 0.03$ ; p=0.008) and U-Net ( $0.78 \pm 0.04$ ; p=0.001) The same significance holds on ISIC 2024 (HUT: 0. Ninety  $\pm 0.01$  vs. TransUNet: zero.87  $\pm 0.02$ ; p=zero.01) and Chest X-Ray (HUT: zero. Ninety-six  $\pm 0.01$  vs. Zero.95  $\pm 0.01$ ; p=zero.03). We add a new supplementary Table 3 with full statistical results across all datasets, including in-line metrics to illustrate consistent superiority. The significance testing confirms that HUT gains are not by random version but are a consequence of the version's superior capability in tackling diverse scientific imaging tasks.

#### Table 1. Quantitative Segmentation Results on Different Datasets

Method	BraTS 2020 (DMC)	ISIC 2024 (DMC)	Chest X-Rays (DMC)
U-Net	0.78	0.82	0.92
Attention U-Net	0.81	0.85	0.94
TransUNet	0.83	0.87	0.95
HUT (Ours)	0.86	0.90	0.96

#### Table 2. Calculation efficiency and comparison of performance

Model	GPU Memory (GB)	Inference Time (s/img)	<b>BraTS Dice</b>
U-Net	3.2	0.38	0.78
TransUNet	6.2	0.51	0.83
HUT	4.8	0.42	0.86

#### Table 3. Calculation efficiency and comparison of performance

Model	Mean Dice ± STD	95% CI	p-value (vs. HUT)
HUT	$0.86\pm0.02$	0.85, 0.87	-
TransUNet	$0.83\pm0.03$	0.81, 0.85	0.008
U-Net	$0.78\pm0.04$	0.75, 0.81	0.001

The U-Net has proven high efficiency in increasing brain tumor segmentation accuracy in comparison with other approaches. The encoder ability in detecting global dependencies and context information is crucial in segmenting intricate brain tumors. Furthermore, HUT has proven efficiency in segmentation skin lesions in the ISIC 2024 dataset with a DMC score of 0.90. The ability of HUT in segmentation boundary lesions without reliance on intermediate feature maps, in addition to considering differences in the color and texture of the skin, is particularly important. The chest X-Ray dataset revealed the direct recognition ability of HUT in identifying features in lungs with a DMC score of 0.96. The addition of deep supervision and connections with residuals in the U-Net decoder brings added benefits in performance and stability in the model. The high performance on a wide variety of datasets and segmentation problems reveals the ability and resilience in the model. It shows adaptability in varying imaging modalities and anatomy and is an asset in a wide variety of medical uses. The computational complexity in the encoder could represent a limitation for significant pictures.



FIGURE 2. Brast 2021 image databases



FIGURE 3. ISIC-2024 Positive Cases



FIGURE 5. The normal chest X-ray (left panel) depicts clear lungs without any areas of abnormal opacification in the image. Bacterial pneumonia (middle) typically exhibits a focal lobar consolidation, in this case in the right upper lobe (white arrows), whereas viral pneumonia (right) manifests with a more diffuse "interstitial" pattern in both lungs.

To presents a visual representation of the process of relating and assaying skin lesions, likely for dermatological assessment and implicit skin cancer discovery. It's organized into three distinct columns, each furnishing a different perspective on the lesions as shown in figure 5. The leftmost column shows a 3D picture of a person's back and arms. This image is used for whole- body skin lesion mapping. Green blotches punctuate specific locales on the body where lesions of interest are present. This provides spatial environment for the lesions, showing their position on the body. The middle column displays near- over, standard photos of the skin lesions linked on the whole- body image. These are" pipe" images because they're lower, localized views uprooted from the larger body image. They give a more detailed view of the lesion's appearance on the skin face, including its shape, color, and texture. The rightmost column showcases dermoscopic images of the same lesions. Dermoscopy is a non-invasive fashion that uses a technical microscope- suchlike device (a dermatoscope) to magnify and illuminate the skin, allowing for a more detailed examination of subterranean structures that are not visible to the naked eye. These images reveal patterns, colors, and vascular structures within the lesion, furnishing pivotal information for secerning between benign and nasty skin lesions.



FIGURE 6. the image illustrates a multi-faceted approach to skin lesion analysis, combining whole-body mapping for lesion identification, close-up photography for surface examination, and dermoscopy for detailed subsurface visualization. This integrated approach is essential for accurate diagnosis and management of skin lesions.

In order to get the Receiver Operating Characteristic (ROC) wind, a graph of the performance of a double bracket model at colorful threshold settings. Following is an explanation as illustrated in figure 6. TheX-axis is the False Positive Rate (FPR), this is the rate at which true negatives are inaptly classified as cons. It is calculated by FPR = False Cons (False Cons True Negatives). The Y- axis is the True Positive Rate (TPR), which is also called perceptivity or Recall. It's the proportion of factual cons which are prognosticated to be cons. This is expressed as TPR = True Cons (True Cons False Negatives). The Ca (Blue) wind is bracket model" Ca" ROC wind. The Cb (Red) wind is bracket model" Cb" ROC wind. A better bracket model will have a ROC wind that's advanced on the left- hand side of the graph. That is having a high TPR and low FPR. That is, the model is picking up utmost of the positive cases without so numerous false admonitions. By eye," Ca" (blue) is performing better than" Cb" (red) because its wind is advanced than" Cb" for utmost of the range of FPR. This means that at some FPR," Ca" can achieve an advanced TPR. A vertical dashed line with marker" TPRO" is colluded at some value of TPR. The areas above this line and under each wind are shadowed. The blue- shaved area is the enhancement of "Ca" over "Cb" at advanced values of TPR. The green area is where both models are above the TPRO threshold, but Ca is nevertheless advanced than Cb in terms of TPR.



**FIGURE 7.** The ROC curve allows to compare the performance of different classification models. In this case, "Ca" demonstrates superior performance compared to "Cb" across a range of threshold settings.

Figure 8 demonstrates the results of a brain tumor segmentation algorithm. The algorithm has identified and delineated the enhancing tumor core (green) and the peritumoral edema (blue) based on the MRI data. The comparison between the raw MRI data and the segmented images highlights the effectiveness of the algorithm in identifying these regions of interest.



FIGURE 8. Brain tumor segmentation algorithm results

The relative performance measured by bones score between different medical image segmentation styles is evident through a significant comparison. The bones score, which takes values between 0 and 1, calculates the overlap ratio of the predicted and base verity segmentations, with higher scores indicating good performance. HUT performed exceptionally well with a bones score of 0.93, surpassing all other architectures that were evaluated. Its exceptional performance is due to its hierarchical architecture and capacity to learn both original and global features accurately in medical images. Its window shift approach enables effective hierarchical point representation modeling. However, its performance was not comparable to the swin transformer, whose efficiency could be due to its fixed patch- predicated processing approach. The U-Net being a widely utilized traditional architecture used in medical image segmentation recorded a lowest bones score of 0.79.



Performance Comparison of Medical Image Segmentation Methods



FIGURE 9. Comparison of Dice scores among different segmentation models. The HUT model achieved the highest Dice score of 0.93, followed by the swin transformer at 0.88, vision transformer at 0.85, and U-Net at 0.79.

## **5. CONCLUSION**

HUT achieves superior segmentation performance compared to several state-of-the-art styles across a different set of medical imaging datasets, including brain excrescence segmentation (BraTS 2021), skin lesion segmentation (ISIC 2024), and lung segmentation (chest X-Ravs). The quantitative results demonstrated constantly advanced DMC give compelling substantiation of HUT's bettered delicacy and robustness. HUT introduces a new trade-off between accuracy and efficiency for medical image segmentation, verified by thorough statistical testing. The bettered delicacy and effectiveness of HUT have significant eventuality for transubstantiating the effectiveness and delicacy of medical image analysis in real- world clinical settings. More precise and dependable segmentation can lead to better individual delicacy, more individualized treatment planning, bettered surgical issues, and more dependable monitoring of complaint progression. This, in turn, can restate into bettered patient care and reduced healthcare costs. The automated nature of HUT also reduces the reliance on homemade segmentation, freeing up clinicians' time and reducing the eventuality forinter-observer variability. Working on several crucial areas to further enhance the capabilities and connection of HUT. These include optimizing the structure, especially the transform encoder, to improve computational efficiency and scalability for larger medical images; applying HUT to other medical image types like CT scans, ultrasound, and histopathology while addressing specific challenges; extending the framework for more complex segmentation tasks; using self-supervised learning to reduce the need for labeled data; and assessing HUT's clinical impact through prospective studies. While HUT demonstrates advanced performance throughout more than one dataset, we acknowledge numerous key obstacles which include the computational necessities stay higher than natural CNNs, doubtlessly proscribing deployment in resource-constrained medical environments. The Performance assessment was limited to 3 public datasets which won't seize the whole spectrum of real-world scientific variability in phrases of imaging protocols or affected person populations. The model's effectiveness relies upon on enough annotated schooling facts, and performance on novel modalities may additionally require additional satisfactory tuning. The hybrid structure introduces additional complexity in hyperparameter optimization whilst adapting to new programs. These obstacles highlight vital guidelines for destiny studies, inclusive of developing lightweight variations of HUT thru pruning and quantization techniques, expanding validation thru multi-center scientific trials, and investigating self-supervised gaining knowledge of procedures to reduce annotation dependence. We accept as true with this honest appraisal strengthens the have a look at's credibility even as keeping focus on HUT's verified advantages in clinical image segmentation.

## REFERENCES

[1] X. Liu, L. Song, S. Liu, and Y. Zhang, "A Review of Deep-Learning-Based Medical Image Segmentation Methods," *Sustainability*, vol. 13, no. 3, p. 1224, Jan. 2021, doi: https://doi.org/10.3390/su13031224.

[2] H. Shang, Y. Tong, M. Li, S. Xu, L. Xu, and Z. Cao, "Deep Learning-based U-Mamba Model to Predict Differentiated Gastric Cancer using Radiomics Features from Spleen Segmentation," *Current Medical Imaging Formerly Current Medical Imaging Reviews*, vol. 20, Dec. 2024, doi: https://doi.org/10.2174/0115734056349216241118115005.

[3] R. Wang, T. Lei, R. Cui, B. Zhang, H. Meng, and A. K. Nandi, "Medical image segmentation using deep learning: A survey," *IET Image Processing*, Jan. 2022, doi: https://doi.org/10.1049/ipr2.12419.

[4] W. Weng and X. Zhu, "INet: Convolutional Networks for Biomedical Image Segmentation," *IEEE Access*, vol. 9, pp. 16591–16603, 2021, doi: https://doi.org/10.1109/access.2021.3053408.

[5] A. Dosovitskiy *et al.*, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *International Conference on Learning Representations*, May 2021, [online]. Available: https://openreview.net/pdf?id=YicbFdNTTy

[6] J. Chen *et al.*, "TransUNet: Rethinking the U-Net architecture design for medical image segmentation through the lens of transformers," *Medical Image Analysis*, pp. 103280–103280, Jul. 2024, doi: https://doi.org/10.1016/j.media.2024.103280.

[7] S. Tong, Z. Zuo, Z. Liu, D. Sun, and T. Zhou, "Hybrid attention mechanism of feature fusion for medical image segmentation," *IET Image Processing*, vol. 18, no. 1, pp. 77–87, Sep. 2023, doi: https://doi.org/10.1049/ipr2.12934.

[8] Z. Liu *et al.*, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9992–10002, Oct. 2021, doi: 10.1109/iccv48922.2021.00986.

[9] X. Lin, L. Yu, K.-T. Cheng, and Z. Yan, "BATFormer: Towards Boundary-Aware Lightweight Transformer for Efficient Medical Image Segmentation," *IEEE journal of biomedical and health informatics*, vol. 27, no. 7, pp. 3501–3512, Jul. 2023, doi: https://doi.org/10.1109/jbhi.2023.3266977.

[10] B. H. Menze *et al.*, "The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS)," *PMC*, Oct. 2015, [Online]. Available: https://dspace.mit.edu/handle/1721.1/110992

[11] "ISIC 2024 - Skin Cancer Detection with 3D-TBP," @kaggle, 2024. https://www.kaggle.com/competitions/isic-2024-challenge/overview/ (accessed Feb. 27, 2025).

[12] D. S. Kermany *et al.*, "Identifying medical diagnoses and treatable diseases by Image-Based Deep Learning," *Cell*, vol. 172, no. 5, pp. 1122-1131.e9, Feb. 2018, doi: 10.1016/j.cell.2018.02.010.