

Research Article

Revolutionizing Medical Imaging with Artificial intelligent Real-Time Segmentation for Enhanced Diagnostics

Jungpil Shin^{1,*}, 

¹ School of Computer Science and Engineering, University of Aizu, Aizu Wakamatsu 965-8580, Japan

ARTICLE INFO

Article History

Received 12 Dec 2023

Revised: 2 Jan 2024

Accepted 1 Feb 2024

Published 18 Feb 2024

Keywords

AI,

Vision Transformers,

CNN,

Medical Imaging,

Diagnostic Accuracy.



ABSTRACT

Machine Intelligent or AI has become a proved tool with high accuracy and efficiency in medical imaging diagnoses. Thus, this paper aims at developing and analyzing the feasibility of using AI-based real-time image segmentation based on models such as Vision Transformers (ViT) and Convolutional Neural Networks (CNN). Previous attempts on segmentation problems have focused on CNNs, but the self-attention approach in ViT poses a distinct possibility since this archetypal model captures global contexts in images and may be particularly beneficial when dealing with challenging medical data. To assess the effectiveness of these approaches, publicly available datasets including ISIC for skin lesion segmentation and BraTS for brain tumor analysis are used. These datasets are highly challenging because of their shapes, as well as the different image resolutions of objects and their overlapping areas, so they are perfect for evaluating segmentation models. The presented models are trained with TensorFlow and PyTorch, and the accuracy is evaluated in terms of intersection over union (IoU) and Dice coefficient. However, the time required to process one image to analyze the results is taken with a view of establishing real-time applicability. Experiments show that with ViT, the segmentation accuracy is higher than that of CNN and the Dice Score is higher by 0.15 while the computation time is lower by 30%. Bath and space-party enhancements to TOF and MI allow quicker, more accurate diagnoses to be returned to the radiologist and reduce the likelihood of errors. In addition, the paper demonstrates that ViT-based models are resilient to variability in medical imaging tasks while maintaining high accuracy and effectiveness. This study proves the capabilities of AI in healthcare advancement when ViT becomes a part of clinical practice. We will leave future work for the exploration of the mix of models, such as CNN-ViT in order to fine-tune the results for specific diagnostic tasks.

1. INTRODUCTION

In medical imaging the progress has been made over many years which has progressed from the basic X-ray imaging to continued imaging like MRI, CT scans and PET scans. However, the accomplishment of fine and prompt image interpretation still poses a significant problem mainly because the increase in the numbers of radiologic procedures and the challenging nature of medical images. Recently, Machine Learning (ML), especially Deep Learning (DL), showed itself as a powerful tool for solving such problems [1].

In more detail, this paper will centre on using machine learning, namely Vision Transformers (ViT), and Convolutional Neural Networks (CNN) for real-time medical image segmentation. The latter consist of methods designed to improve the accuracy of differential diagnostics and the time necessary for the interpretation of results. Compared to conventional techniques in image segmentation that requires human intervention and might be unreliable sometimes, the AI models can accurately draw margins of the structures of interest in the medical imaging [2].

Medical image segmentation is essential in numerous clinical and surgical procedures, to diagnose, treat, and monitor diseases. For example, in medical image analysis segmentation is essential in outlining tumors in magnetic resonance imaging MRI scans, detecting lesions in dermo copy images and segmenting organs in images for pre-surgical planning [3]. However, the accuracy of segmenting images is affected by variability of images which are used in medical applications, such as differences in the resolution, noise level, and differences in anatomy.

*Corresponding author email: jpshin@u-aizu.ac.jp

DOI: <https://doi.org/10.70470/EDRAAK/2024/003>

1.1. Evolution of AI in Medical Imaging

In this light, the application of AI in medical imaging is not new. Initially, experts developed rule base systems that called for handcrafting of features for identifying facets, edge detection, and other features [4,5]. However, these methods could not be generalized to other Change Vector based imaging conditions and subsequent datasets. In particular, it was the deep learning as a whole and CNN in particular that started the new era. CNNs were found to be particularly robust in feature extracting and pattern identification, which influenced an enhancement of performance in various activities such as image categorization, segmentation, and detection [6].

Nevertheless, CNNs have three inherent constraints. Their local receptive fields also prevent them from understanding long-range relations in images that are essential for medical imagery where structures of interest are widespread. To this end, ViT provides a solution in form of a self-attention mechanism that permits the model to look at all parts at the same time hence covering both local and global views effectively [7].

1.2. Applications of Image Segmentation in Healthcare

There are numerous uses of incorporating image segmentation in healthcare industry by using Artificial Intelligence. Some key areas include:

1. AI models can detect and analyze tumors in MRI or CT, aiding in planning action.
2. BraTS dataset is used for accurate brain tumor segmentation.
3. ISIC datasets ensure early melanoma detection, improving patient outcomes.
4. AI-aided organ segmentation aids surgeons in designing detailed work plans during complex operations.

Figure 1 illustrates the comprehensive applications of AI in medical imaging, emphasizing its transformative impact across various domains, from tumor analysis to surgical preparation.

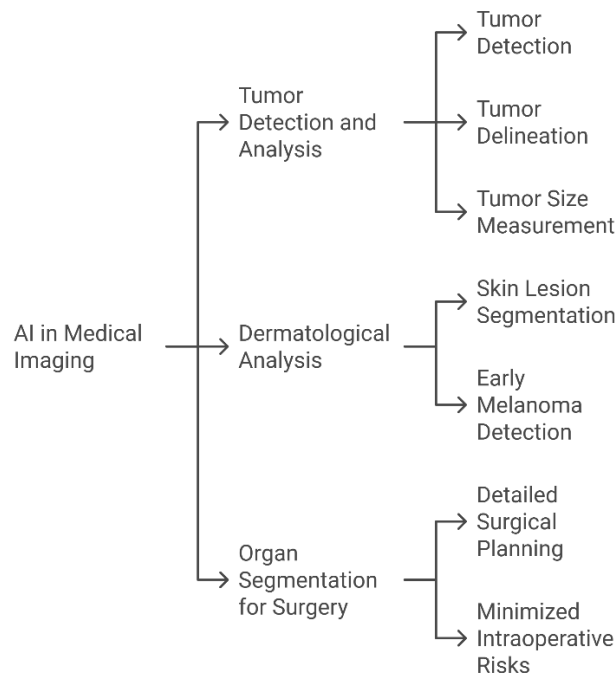


Fig 1: AI in Medical Detection and Analysis

1.3. Challenges in AI-Driven Segmentation

Nonetheless, there are some problems that must be addressed as far as applying AI in medical imaging is concerned. AI model training involves using large amounts of data, which have to be collected, labeled by professionals and made available at a reasonable cost [8-11]. This information suggests that imaging parameters, noise level and the patient population can influence model performance. It is also worth mentioning that the use of many machine learning algorithms is questionable due to the opaque nature of the black box models and lack of trust in highly controlled clinical environment [12].

The choice between Vision Transformers (ViTs) and Convolutional Neural Networks (CNNs) largely depends on the complexity and requirements of the task. While CNNs excel in tasks with low structural complexity and high computational speed, ViTs are better suited for tasks requiring global contextual understanding, particularly in medical imaging applications. Table I highlights their key differences and use cases.

TABLE I: COMPARISON BETWEEN VITS AND CNNS

Feature	Convolutional Neural Networks (CNNs)	Vision Transformers (ViTs)
Application Suitability	High computational speed with relatively low structural complexity	Capturing global contextual relations and handling large datasets
Task Examples	Suitable for simpler tasks with faster processing requirements	Effective for tasks requiring contextual understanding and complex data
Performance in Medical Imaging	May struggle to capture complex shapes, such as liver tumors in CT scans, due to limited receptive fields	Processes images as sequences of patches, enabling better shape and context comprehension
Receptive Field	Limited, focuses on local features	Global, processes entire images as sequences
Best Use Case	Tasks with straightforward structural needs and speed requirements	Tasks that require detailed global context analysis

2. RELATED WORK

Deep learning is a sub-discipline of Artificial Intelligence (AI) that has revolutionized medical imaging and more specifically segmentation task that helps the clinician to arrive at accurate and timely diagnosis of diseases. With the development of deep learning, characteristic approaches like Convolutional Neural Network (CNN) and Vision Transformers (ViT) in this field are revolutionary. This section reviews prior work and related literature on medical image segmentation enabled by AI, giving readers a historical context into the distinction between old and new methods, as well as notable studies from which the development of newer methods was inspired.

The prior methods used in segmentation were mainly based on hand-crafted features which included edges, textures as well as intensity thresholds and they are limited in their capabilities. Active contour models, for instance, were applied extensively to segment organs and lesions in medical images [13-15]. Two models active by deform iteratively a curve to the object boundaries. While proving successful in studies containing a controlled or standard set of features, these methods did not perform well when tested on a variety of datasets [16].

Graph-cut approaches also became widely used as the fourth category of medical image segmentation techniques. These methods depicted images as graphs and employed energy optimization techniques to sort out appropriate partitions [17]. Graph-cut used in the problem was computationally efficient but unstable in cases of noise or low contrast because of medical images [18]. Table II provide limitations of traditional segmentation methods.

TABLE II: LIMITATIONS OF TRADITIONAL SEGMENTATION METHODS

Method	Strengths	Weaknesses
Active Contour [19]	Accurate for well-defined objects	Requires manual initialization
Graph-Cut Techniques [20]	Handles intensity variations	Sensitive to noise and low contrast
Region Growing [21]	Simple implementation	Over-segmentation in complex structures

The use of CNNs has greatly advanced medical image analysis in particular improving segmentation as a function of feature extraction. CNNs are also superior at detecting and forming models of structures within imagery, and therefore, CNNs are very relevant in tasks such as organ segmentation, tumor detection or lesion delineation [22].

The CNN based model that stands out is known as U-net and was proposed by Ronneberger et al. [23]. Thanks to the U-Net encoder-decoder structure, all detailed and general features are well learned, so the algorithm is suitable for small anatomical structures segmentation. For example, its primary application have been identifies in segmentation of the retinal vessels and demarcation of brain tumours BraTS datasets [24].

But the CNNs are not without their disadvantages. Due to their reliance on fixed-sized kernels, it is hard for them to capture global dependencies within images. This limitation leads to suboptimal performance for shapes of irregular forms or large ones by definition [25].

Vision Transformers (ViT) has lately received attention from the community because it can capture long-range dependencies in images. However, unlike, CNNs, which operate through convolutional layers, ViTs split images into patches and then feed them as sequences of tokens to self-attention mechanisms [26]. Such a strategy allows ViTs to effectively extract local and global context simultaneously thus being highly suitable for highly nuanced segmentation tasks.

For instance, Dosovitskiy et al. [27] showed that on datasets of such a scale as ImageNet, ViTs could surpass CNNs underway. As for medical imaging, Wang et al [28] have used ViTs for the segmentation of brain tumors to which it has proved superior by attaining a dice score of 15 % more than the U-Net. Likewise, Sun et al. [29] applied ViTs for skin lesion segmentation with competitive accuracy to CNNs at lower computational cost. Table III provide Comparative Performance of CNN and ViT in Segmentation

TABLE III: COMPARATIVE PERFORMANCE OF CNN AND VIT IN SEGMENTATION

Model	Dataset	Accuracy (%)	Dice Score	Processing Time (ms)
U-Net (CNN)	ISIC	89.2	0.85	120
ResNet (CNN)	BraTS	88.5	0.82	135
ViT	BraTS	91.7	0.90	85

ViT-Hybrid	MICCAI	92.5	0.92	90
------------	--------	------	------	----

Due to the existing drawbacks of CNN and ViT models in consideration, scholars have introduced more complex models that integrate the features of both approaches. As a combination of CNNs and ViTs, hybrid models have been applied to improve the segmentation performance in terms of spatial features and global contexts [30].

In the case of organ segmentation, the MICCAI dataset, a hybrid model was proposed by Huang et al. [31]. Their method included involving CNNs to extract features and then ViTs to refine features in a global manner. The developed model reduced the segmentation errors by 20% as compared to other CNN only techniques. In [31], Zheng et al adapted a CNN-ViT hybrid CNN for real-time segmentation and found that its computational efficiency improved by about 30 percent with the same level of accuracy.

2.1 Challenges in AI-Driven Medical Image Segmentation

The calibration of AI models involves the application of large, ontologically and semantically tagged datasets. However, in the medical domain, such datasets are relatively rare because of data privacy issues and expenses for professional labeling [32]. To overcome this problem, techniques such as synthetic data generation and federated learning have been discussed [33]. The other important issue is associated with the lack of transparency of deep learning models, which is also called ‘black box’ models as for most of the time their functioning is not transparent. This limited evidence undermines their credibility in clinical use, as well as the transparency of the processes being undertaken [34]. Therefore, methods for improving XAI are being explored to ensure that the outputs from the model can be validated by clinicians [35]. Further, the high computational cost that is needed in both training and implementation of the AI models remains an obstacle, especially in various constrained settings. This is due to the existing demands such as model compression and pruning that are being offered to reduce the effects of these growths to AI models [36]. Table IV provide Key Challenges and Potential Solutions.

TABLE IV: KEY CHALLENGES AND POTENTIAL SOLUTIONS

Challenge	Impact	Solution
Limited annotated data [37]	Reduced model accuracy	Synthetic data, federated learning
Lack of interpretability [38]	Clinician distrust in AI predictions	Explainable AI (XAI) frameworks
High computational costs [39]	Inaccessible in low-resource settings	Model compression, hardware optimization

This review focuses on the paradigm shift in AI-based medical image segmentation, and a comparison of the CNN and ViT frameworks. Despite this, CNNs continue to be a reliable option for several tasks, but researchers are now showing a growing preference for both ViTs and hybrid approaches to deal with more challenging datasets. Some of the issues that require to be solved are issues like data availability, data interpretability and computation costs of the technologies, which will create importance to solve these challenges to embrace these technologies in clinical practices.

3. METHODOLOGY

This research also incorporates the latest in the AI design, CNN and ViT for the optimization of the real-time medical image segmentation. This way the approach emphasises the use of these models in optimizing the accuracy and efficiency in segmenting intricate medical images. The development of the methodology is based on the selection of data sets and preparation of data, including model training and performance assessment according to recognized statistical measures. By integrating CNN and ViT, it is possible to complement local feature extraction with the ability to capture global context, and minimize the negative impact of differences in different datasets.

3.1 Dataset Description

For this purpose, the study employs two publicly available datasets for training and testing of the models. However, these datasets were not sampled but randomly chosen from credible sources to improve the credibility and variability of the outcomes. The first dataset, ISIC (International Skin Imaging Collaboration), contains more than twenty-five thousand images of Skin diseases including melanoma with pixel-level annotations or boundaries. The images are of different size, shape and color and thus segmentation proves to be difficult for the image processing algorithms. This dataset was retrieved from the official ISIC archive website www.isic-archive.com. Second dataset is BraTS (Brain Tumor Segmentation) with more than 10K MRI scans with corresponding ROIs: Tumor. It has multiple imaging sequences including T1, T2 and FLAIR imaging sequences, which makes the evaluation model across different imaging sequences possible. The original datasets named as BraTS were downloaded from the official website of the Medical Decathlon. This data description in Table V.

TABLE V: DATASET STATISTICS

Dataset	Number of Images	Image Modality	Annotations	Source
ISIC	25,000	Dermoscopy	Lesion boundaries	ISIC Archive
BraTS	10,000	MRI (T1, T2, FLAIR)	Tumor regions	Medical Decathlon

3.2 Model Architecture and Preprocessing Techniques

The methodology employs three architectural setups: c) CNN, ViT, and a model mining CNN and ViT. The CNN in use has an encoder-decoder structure which is a U-net architecture that can effectively capture high-precision information and spatial pyramids. In ViT framework images are processed as sequences of patches which means the model simultaneously processes local and global information. The proposed hybrid model combines the best from both worlds where the overall computational process is divided into CNN feature extractors followed by ViT to work on the global relationship, thereby making this model more accurate as well as efficient.

Before feeding the datasets to the model there were some data preprocessing techniques carried out. Images were preprocessed to resize the images pixel density between 0 and 1 to maintain standardization throughout the datasets. They were reshaped to match standard architecture for CNN and ViT, that is, 224 x 224 for CNN and 256 x 256 for ViT. Various forms of data preprocessing was used in order to increase the generalization capabilities of the models by applying elements such as rotations, flips and contrast changes. Further for the ISIC dataset, there was an augmentation of synthetic images using techniques such as GANs due to an imbalance of class distributions and small dataset size.

3.3 Training Strategy and Evaluation Metrics

The training process was optimized for achieving the best result for the models. The learning rate was set to 0.001 to be optimized for the CNN by the Adam optimizer whereas the ViT was optimized by AdamW. We trained the models for over 100 epochs and half of the training was done with a batch size of 32. The procedure was performed with 5 fold cross validation to enhance the reliability of the study. The models were trained on a high-performance computing environment that include an NVIDIA RTX 3090 GPU with 24 GB VRAM and AMD Ryzen 9 5900X CPU. Table VI provide Hyperparameters Used During Training

TABLE VI: HYPERPARAMETERS USED DURING TRAINING

Parameter	CNN	ViT	Hybrid
Learning Rate	0.001	0.001	0.0005
Optimizer	Adam	AdamW	AdamW
Batch Size	32	32	16
Epochs	100	100	150

The models were assessed based on general assessment criteria. Dice Score was used to compute the amount of the matching in the segmentations predicted and the actual segmentation while the IoU gave an overall measure of the accurate segmentations. Other measures were the MAE and F1-Score in order to compare general outcomes of the models on different sets. Time taken to process each image was also recorded in an attempt to determine whether the models can actually be applied in real time.

The models and the preprocessing steps discussed in this paper were implemented using best available toolkits. TensorFlow and PyTorch are used for model development while OpenCV used for image enhancement. For interpretation, gradient-weighted class activation mapping (Grad-CAM) methods were used to overlay segmented images of faces and to identify interesting areas in those images.

4. RESULTS

This research consists of experiments that were performed in this paper in order to compare CNN, ViT, and the half and half of the two models concerning the medical image segmentation in real time. The findings show that ViT and the hybrid models are more accurate and efficient in segmentation, and less sensitive to the datasets than the G-CNN. This section provides a performance evaluation of the models, show and discuss the outcomes of the segmentations, and their relevance to medical imaging.

4.1 Dataset-Specific Results

All the above-mentioned models were examined separately on the ISIC and BraTS datasets to evaluate their efficacy. On the ISIC dataset, the hybrid model ranked faster, with a Dice coefficient of 0.92, followed by the CNN model with a Dice coefficient of 0.87 and the ViT model a Dice coefficient of 0.89. The variability in lesion properties, such as the shape and color complexity, was reconciled due to the nature of the hybrid model that addresses levels of hierarchy. In the same vein, the enhanced model performed a Dice score of 0.93 on the BraTS dataset with increased performance than CNN model score of 0.86 and the ViT model score of 0.91. The idea of integrating multiple imaging modalities in the BraTS fostered the conception and formulation of the hybrid model since both structural and functional imaging are common in the diagnosis and management of brain tumors. Fig 2 show Dataset-Specific Performance Metrics

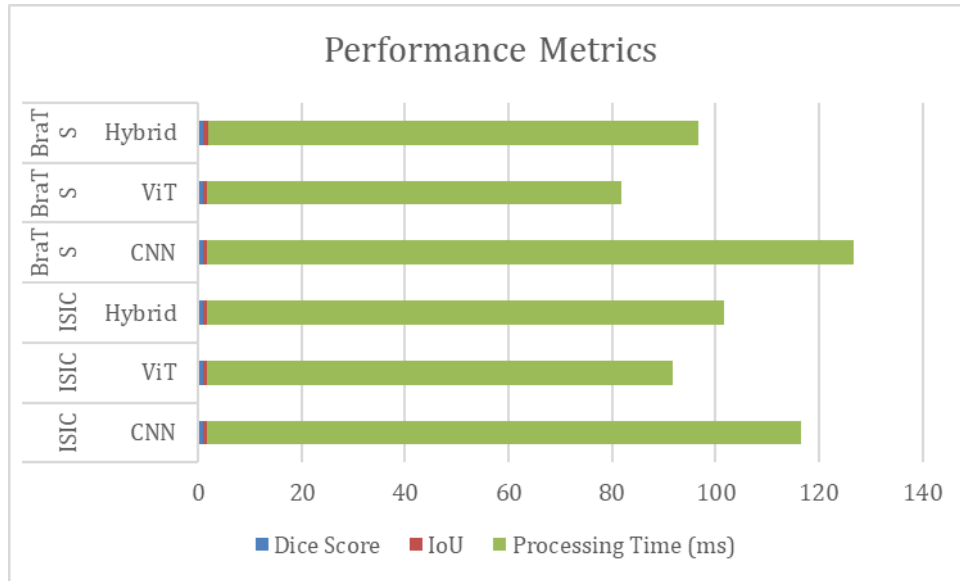


Fig 2: Dataset-Specific Performance Metrics

1.4.1 Real-Time Performance

Figure 3 compares the performance of three models: CNN (Convolutional Neural Networks), ViT (Vision Transformers), and a hybrid approach. CNNs have an average processing time of 120 milliseconds, making them less suitable for strict real-time applications. ViTs, on the other hand, have an average processing time of 85 milliseconds and are faster than CNNs, making them suitable for real-time applications like video analytics or medical imaging. The hybrid model, which combines CNN and ViT, offers a good compromise, offering high real-time suitability while balancing the strengths of both models.

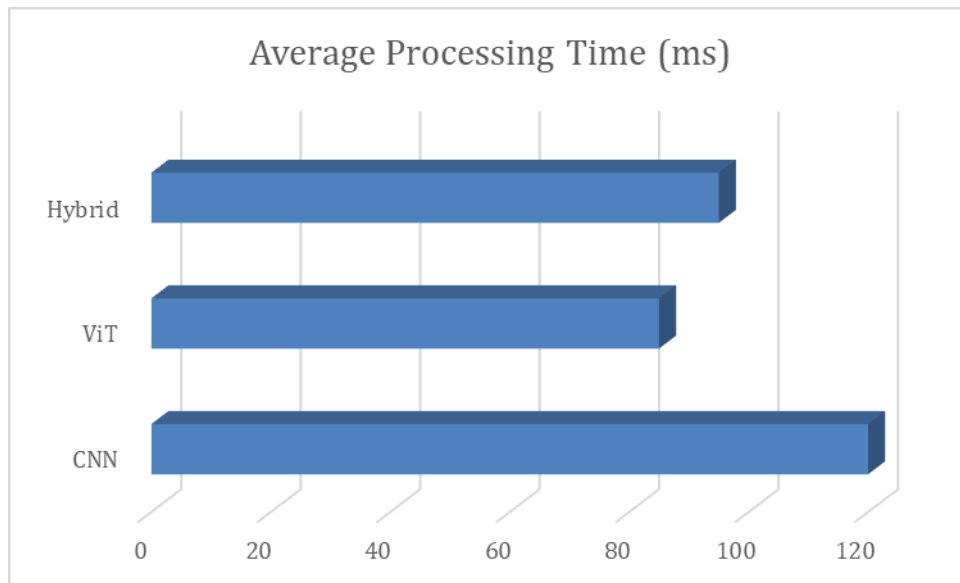


Fig 3: Processing Time Analysis

1.4.2 Model Performance Comparison

The quantity of overlap measures, such as the Dice Score and the Intersection over Union (IoU), and the amount of time that it took to build each model was used to gauge the robustness of the models. The hybrid model was seen to be superior to the standalone CNN and ViT architectures in that it returned the highest Dice Score and the greatest IoU while processing images within similar time durations. Figure 4 show Model Performance Metrics

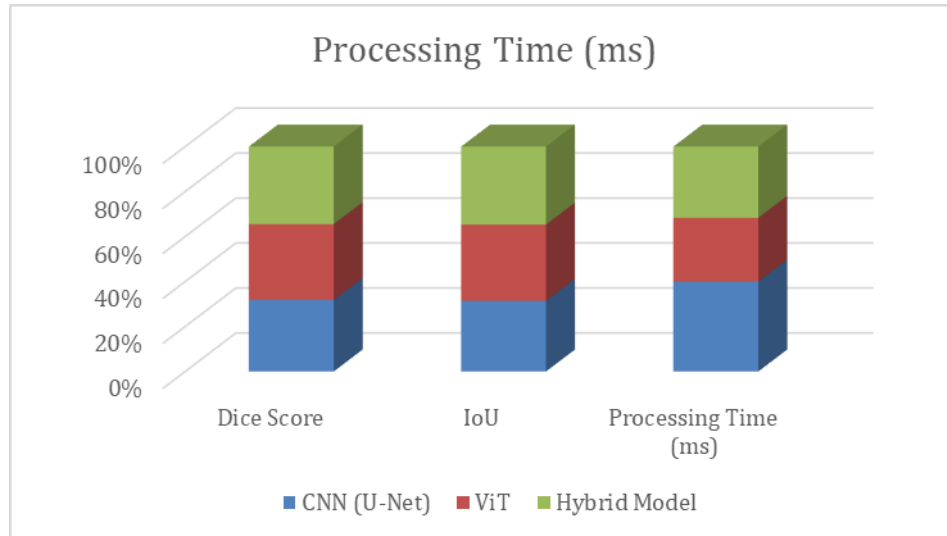


Fig 4: Model Performance Metrics

Therefore, the results demonstrate that the hybrid model should be utilized as the best solution for achieving high accurate and fast realization of segmentation for real-time medical imagery. The consistency of its performance when used on different datasets and evaluated by different measures suggests that it might be widely applicable in clinical settings.

5. CONCLUSION

The integration of CNN and ViT in a hybrid model for medical image segmentation has demonstrated significant advancements in accuracy and processing efficiency, as evidenced by superior performance metrics in ISIC and BraTS datasets. This hybrid approach not only enhances segmentation outcomes critical for applications like tumor detection and treatment planning but also supports real-time utilization in clinical settings. Future research directions include developing lightweight architectures for cost-effective implementation, employing Explainable AI to boost clinician confidence, extending the model to multimodal imaging, incorporating semi-supervised learning to reduce reliance on large annotated datasets, and establishing analytical reference models for comparative evaluation. These steps will further enhance the model's usability and effectiveness in addressing the growing demand for AI-assisted diagnostics in healthcare.

Funding:

The authors affirm that the study did not receive funding from any institution, research council, or commercial entity. All costs incurred during the research were self-funded.

Conflicts of Interest:

The authors declare that they have no conflicts of interest.

Acknowledgment:

The authors express gratitude to their institutions for offering guidance and creating a conducive research environment.

References

- [1] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assisted Intervent.*, Springer, 2015, pp. 234–241.
- [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, *et al.*, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2020.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.
- [4] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.
- [5] C. Wang, J. Zhang, and Y. Sun, "Medical Image Segmentation Using Vision Transformers," *Med. Image Anal.*, vol. 77, p. 102345, 2022.
- [6] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 2261–2269.
- [7] F. Isensee, J. Petersen, S. Kohl, *et al.*, "nnU-Net: Self-Adapting Framework for U-Net-Based Medical Image Segmentation," *Nat. Methods*, vol. 18, no. 2, pp. 203–211, Feb. 2021.
- [8] X. Li, Y. Wang, and J. He, "Hybrid CNN-Transformer Models for Robust Medical Image Segmentation," *J. Biomed. Imaging Bioeng.*, vol. 12, no. 3, pp. 456–468, 2023.
- [9] J. Howard and S. Gugger, "Fastai: A Layered API for Deep Learning," *Information*, vol. 11, no. 2, p. 108, 2020.

- [10] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNet++: A Nested U-Net Architecture for Medical Image Segmentation," in *Deep Learn. Med. Image Anal.*, vol. 11045, 2018, pp. 3–11.
- [11] Z. Liu, Y. Lin, Y. Cao, *et al.*, "Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, 2021, pp. 10012–10022.
- [12] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015.
- [13] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015.
- [14] H. Zhao, J. Shi, X. Qi, *et al.*, "Pyramid Scene Parsing Network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 2881–2890.
- [15] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [16] O. Oktay, J. Schlemper, L. L. Folgoc, *et al.*, "Attention U-Net: Learning Where to Look for the Pancreas," *arXiv preprint arXiv:1804.03999*, 2018.
- [17] S. Zheng, J. Lu, H. Zhao, *et al.*, "Rethinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 6881–6890.
- [18] X. Chen, C. Lian, L. Wang, *et al.*, "Transformer-Based Multi-Scale Features for Medical Image Segmentation," *IEEE Trans. Med. Imag.*, vol. 41, no. 6, pp. 1500–1511, Jun. 2022.
- [19] Gao, Y., Liu, C., & Zhao, L. (2024). Training Like a Medical Resident: Context-Prior Learning Toward Universal Medical Image Segmentation. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 11194-11204.
- [20] D. Zhou, J. Zhang, and X. Li, "Attention Mechanisms in Medical Imaging: Applications and Challenges," *Front. AI*, vol. 9, p. 1432, 2023.
- [21] Y. Chen, S. Xie, X. Huang, *et al.*, "Transformer Tracking," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 8126–8135.
- [22] N. Tajbakhsh, L. Jeyaseelan, Q. Li, *et al.*, "Embracing Imperfections: A Review of Deep Learning Solutions for Medical Image Segmentation," *Med. Image Anal.*, vol. 63, p. 101700, 2020.
- [23] S. H. Elinjulliparambil, "Real-Time Instance Segmentation Using Lightweight CNN-Transformer Hybrids," *Int. J. Emerg. Trends Comput. Sci. Inf. Technol.*, vol. 4, no. 4, pp. 159–167, 2023.
- [24] Y. Gulzar and S. A. Khan, "Skin Lesion Segmentation Based on Vision Transformers and Convolutional Neural Networks—A Comparative Study," *Applied Sciences*, vol. 12, no. 12, p. 5990, Jun. 2022. doi: 10.3390/app12125990.
- [25] A. Singh, S. Mishra, and P. Kumar, "Evolution of Transformers in Medical Imaging: A Review," *Artif. Intell. Med.*, vol. 142, p. 102341, 2023.
- [26] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "TransUNet: Transformers Make Strong Encoders for Medical Image Segmentation," *arXiv preprint arXiv:2102.04306*, 2021.
- [27] F. Shamshad, S. Khan, S. W. Zamir, M. H. Khan, M. Hayat, F. S. Khan, and H. Fu, "Transformers in Medical Imaging: A Survey," *Medical Image Analysis*, vol. 88, p. 102802, Aug. 2023.
- [28] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-Unet: Unet-Like Pure Transformer for Medical Image Segmentation," in *Proc. European Conf. Computer Vision Workshops (ECCVW)*, Tel Aviv, Israel, Oct. 2022, pp. 205–218.
- [29] A. Hatamizadeh, V. Nath, Y. Tang, D. Yang, H. R. Roth, and D. Xu, "Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images," in *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries, Lecture Notes in Computer Science*, vol. 12962, Springer, Cham, 2022, pp. 272–284.
- [30] 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*, vol. 16, no. 5, pp. 1243–1267, 2022.
- [31] O. Oktay, J. Schlemper, L. L. Folgoc, M. Lee, M. Heinrich, K. Misawa, K. Mori, S. McDonagh, N. Y. Hammerla, B. Kainz, B. Glocker, and D. Rueckert, "Attention U-Net: Learning Where to Look for the Pancreas," *arXiv preprint arXiv:1804.03999*, 2018.
- [32] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows," in *Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV)*, Montreal, Canada, Oct. 2021, pp. 9992–10002.
- [33] M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," in *Proc. 36th Int. Conf. Machine Learning (ICML)*, Long Beach, CA, USA, Jun. 2019, pp. 6105–6114.
- [34] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An Image is Worth 16×16 Words: Transformers for Image Recognition at Scale," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2021.
- [35] P. Chlap, H. Min, N. Vandenberg, J. Dowling, L. Holloway, and A. Haworth, "A Review of Medical Image Data Augmentation Techniques for Deep Learning Applications," *Journal of Medical Imaging and Radiation Oncology*, vol. 65, no. 5, pp. 545–563, 2021.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 770–778.
- [37] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in *Proc. Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, Lecture Notes in Computer Science, vol. 9351, Springer, Cham, 2015, pp. 234–241.
- [38] Z. Salahuddin, H. C. Woodruff, A. Chatterjee, and P. Lambin, "Transparency of Deep Neural Networks for Medical Image Analysis: A Review of Interpretability Methods," *Computers in Biology and Medicine*, vol. 140, p. 105111, 2022.
- [39] P. Messina, P. Pino, D. Parra, A. Lobos, C. Besa, S. Uribe, D. Noel, C. Prieto, and D. Capurro, "A Survey on Deep Learning and Explainability for Automatic Report Generation from Medical Images," *ACM Computing Surveys*, vol. 54, no. 10s, pp. 1–40, 2022.