

# Transformers-based architectures for stroke segmentation: A review

Yalda Zafari-Ghadim<sup>1</sup>, Essam A. Rashed<sup>2</sup>, Mohamed Mabrok<sup>1\*</sup>

<sup>1</sup>Department of Mathematics and Statistics, Qatar University, Doha,  
P.O.Box 2713, Qatar.

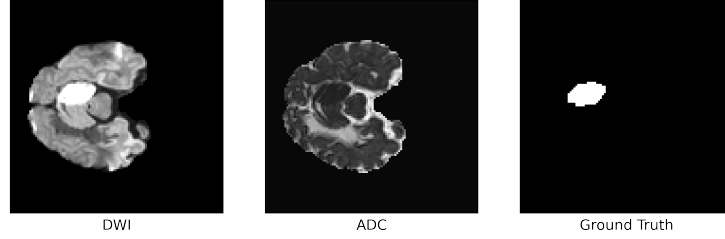
<sup>2</sup>Graduate School of Information Science, University of Hyogo, Kobe  
650-0047, Japan.

\*Corresponding author(s). E-mail(s): [m.a.mabrok@gmail.com](mailto:m.a.mabrok@gmail.com);

## Abstract

Stroke remains a significant global health concern, necessitating precise and efficient diagnostic tools for timely intervention and improved patient outcomes. The emergence of deep learning methodologies has transformed the landscape of medical image analysis. Recently, *Transformers*, initially designed for natural language processing, have exhibited remarkable capabilities in various computer vision applications, including medical image analysis. This comprehensive review aims to provide an in-depth exploration of the cutting-edge Transformer-based architectures applied in the context of stroke segmentation. It commences with an exploration of stroke pathology, imaging modalities, and the challenges associated with accurate diagnosis and segmentation. Subsequently, the review delves into the fundamental ideas of Transformers, offering detailed insights into their architectural intricacies and the underlying mechanisms that empower them to effectively capture complex spatial information within medical images. The existing literature is systematically categorized and analyzed, discussing various approaches that leverage Transformers for stroke segmentation. A critical assessment is provided, highlighting the strengths and limitations of these methods, including considerations of performance and computational efficiency. Additionally, this review explores potential avenues for future research and development.

**Keywords:** stroke segmentation, vision Transformer, deep learning, medical imaging



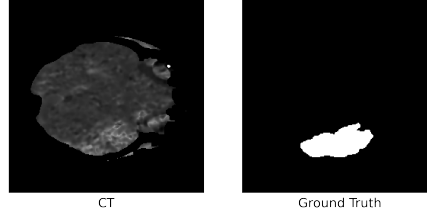
**Fig. 1** Sample of stroke infarct on Diffusion-weighted and Apparent diffusion coefficient MRI with the annotation, from [Hernandez Petzsche et al \(2022\)](#). The infarct appears hyperintense in DWI and hypointense in ADC.

## 1 Introduction

Stroke, a cerebrovascular disease, stands as the second leading cause of morbidity and mortality worldwide, impacting over 100 million people globally [Feigin et al \(2022\)](#). It transpires when there is an abrupt disruption in the blood supply to the brain, resulting in the damage or death of neuro cells. This occurrence can be attributed to two primary reasons: a blockage in the blood vessels, referred to as ischemic stroke, and the rupture of vessels leading to bleeding into surrounding tissues, known as hemorrhagic stroke [Grysiewicz et al \(2008\)](#). The consequences of stroke on patients can be profound, often resulting in physical disabilities and cognitive impairments [Meyer et al \(2015\)](#); [Dimyan and Cohen \(2011\)](#). This underscores the importance of accurate and timely diagnosis for effective treatment and improved patient outcomes.

Stroke patients typically undergo neuroimaging techniques to distinguish between ischemic and hemorrhagic strokes. This differentiation can be achieved through magnetic resonance imaging (MRI) and computed tomography (CT), each offering distinctive insight into the condition of the brain [Goldstein and Simel \(2005\)](#). MRI offers excellent soft tissue contrast for the brain, and when diagnosis is uncertain, it can be more informative than CT [Hwang et al \(2012\)](#); [Chalela et al \(2007\)](#); [Fiebach et al \(2002\)](#), providing information on stroke location [Flossmann et al \(2008\)](#), timing [Aoki et al \(2010\)](#), and mechanism [Wessels et al \(2006\)](#). Diffusion-weighted imaging (DWI) and perfusion-weighted imaging (PWI) within the MRI protocol offer valuable information on the extent and impact of stroke on brain tissue [Simonsen et al \(2015\)](#). Refer to Figure 1 for an illustration showing a stroke infarct sample in two distinct magnetic resonance modalities along with the corresponding annotation. Additionally, refer to Figure 2 for CT images accompanied by corresponding annotations.

Stroke segmentation plays an essential role in the diagnostic process as well as treatment planning by providing spatial information about affected areas of the brain and the extent of damage. Traditional methods of stroke diagnosis, often based on manual interpretation of medical images, prove to be time-consuming and susceptible to human error. The inherent variability in the size, shape, and location of strokes, compounded by artifacts and noise present in the imaging data, presents substantial challenges for automated analysis, rendering it a challenging task. Furthermore, the need for real-time or near-real-time diagnosis in stroke cases demands algorithms that are not only accurate but also computationally efficient. As such, the development of



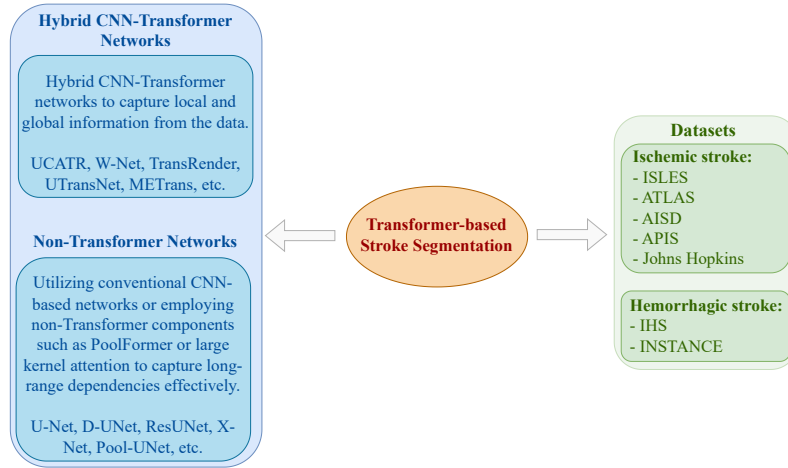
**Fig. 2** Sample of stroke infarct on CT images with the annotation, from [Cereda et al \(2016\)](#); [Hakim et al \(2021\)](#).

accurate and automatic methods for stroke segmentation remains a prominent focus in the research domain.

The field of medical image analysis has witnessed a transformative evolution with the advent of deep learning techniques [Zhou et al \(2019a\)](#). Deep learning models, with their ability to automatically learn intricate patterns from vast amounts of data, have shown promising results in various medical imaging tasks, including stroke segmentation [Zhang et al \(2022\)](#). Convolutional Neural Networks (CNNs) [O’Shea and Nash \(2015\)](#), a class of deep learning models, have demonstrated remarkable success in tasks such as image classification [Huang et al \(2017\)](#); [Hu et al \(2018\)](#), object detection [Wang et al \(2017\)](#), segmentation [Chen et al \(2017\)](#), and registration [Balakrishnan et al \(2019\)](#); [Jia et al \(2022\)](#). These models, with their hierarchical feature learning capabilities, have significantly improved the accuracy and efficiency of medical image analysis. However, the inherent limitations of CNNs in capturing long-range dependencies and contextual information in images have led to the exploration of alternative architectures, including Transformers [Li et al \(2023a\)](#).

Originally proposed for natural language processing tasks [Vaswani et al \(2017\)](#), Transformers have gained widespread attention in the computer vision community. Unlike traditional convolutional approaches, transformers process input data in a parallel and non-sequential manner, allowing them to capture complex spatial relationships and contextual dependencies effectively. The self-attention mechanism in Transformers enables them to weigh different parts of the input data differently, making them particularly suitable for tasks requiring a global understanding of the data, such as medical image analysis.

We conducted a comprehensive search across electronic databases, including PubMed, IEEE Xplore, and Google Scholar. Our search employed diverse queries to compile lists of published works, combining terms like "CNN," "deep learning," "Transformer," and "neural network" with stroke-specific terms such as "ischemic stroke," "stroke segmentation," and "stroke detection." Additionally, we manually explored the reference lists of relevant articles to identify additional studies. We included studies that satisfied the following criteria: (1) publication in English; (2) focus on stroke segmentation; (3) utilization of deep learning techniques, particularly emphasizing Transformer-based architectures; and (4) reporting of quantitative results on the model’s performance.



**Fig. 3** Overview of key aspects covered in this review paper.

We excluded studies that: (1) were unrelated to stroke segmentation; (2) did not achieve high segmentation performance; and (3) presented papers with identical methodologies where their contributions were negligible. However, in cases where the performance of all proposed pipelines was low for certain datasets, we kept the superior ones. Due to the relatively limited number of Transformer-based networks addressing stroke segmentation in the existing literature, we encompassed all available publications in our study. It is crucial to acknowledge that our review might have unintentionally omitted some noteworthy papers related to CNN-based architectures. Nevertheless, our primary objective was to offer an overview of the contributions in utilizing vision Transformers for stroke segmentation purposes.

In this review, we have discussed the applications of Transformers in stroke segmentation, exploring innovative methodologies developed to address the challenges posed by stroke diagnosis. We systematically reviewed the existing literature, analyzing different Transformer-based architectures, their integration with traditional deep learning techniques, and their performance in stroke-related tasks. Through this comprehensive review, we aimed to provide insight into the current state of the art, highlight the strengths and limitations of Transformers in stroke segmentation, and identify potential avenues for future research and development. Figure 3 provides a graphical abstract highlighting the key aspects discussed in this paper, encompassing available datasets for stroke segmentation and the diverse deep architectures employed in this context.

## 2 Fundamentals of Transformers

### 2.1 Architectural Components of Transformers

Transformers, introduced in the context of natural language processing, consist of fundamental architectural components that distinguish them from traditional CNN-based models. The core elements of Transformers include self-attention mechanisms and

position-wise feed-forward networks. Self-attention enables the model to weigh input elements differently, capturing contextual dependencies irrespective of their positions in the sequence. This mechanism allows Transformers to model long-range dependencies efficiently, making them well-suited for tasks requiring a global understanding of the data. Position-wise feed-forward networks introduce non-linear transformations, enhancing the model’s capacity to learn complex patterns within the data.

### 2.1.1 Self-Attention

The self-attention (SA) mechanism used in Transformers is a crucial component that empowers the model to capture the long-range dependencies between various parts of the input data. This is accomplished through a process in which each element (token) in the input sequence attends to every other element, calculating its representation based on the information from all other elements.

To compute the self-attention, the input sequence  $X \in \mathbb{R}^{N \times C}$  is projected into a query  $Q \in \mathbb{R}^{N \times D}$ , a key  $K \in \mathbb{R}^{N \times D}$ , and a value  $V \in \mathbb{R}^{N \times D_v}$  using three trainable projection layers  $W^Q \in \mathbb{R}^{C \times D}$ ,  $W^K \in \mathbb{R}^{C \times D}$ ,  $W^V \in \mathbb{R}^{C \times D_v}$ , respectively. Then, the corresponding attention matrix  $A \in \mathbb{R}^{N \times N}$ , which represents the affinity of the query and the key, can be calculated by:

$$A(Q, K) = \text{Softmax}\left(\frac{Q \times K^T}{\sqrt{D}}\right) \quad (1)$$

The attention matrix connects all elements, which allows the handling of long-range dependencies. Subsequently, the calculated attention matrix is applied to the value  $V$ , resulting in the output  $Z \in \mathbb{R}^{N \times D_v}$ :

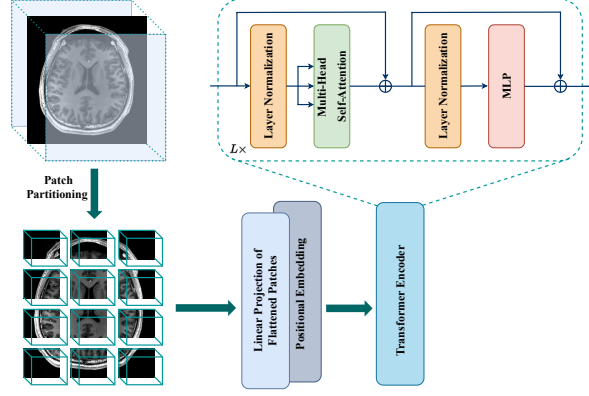
$$Z = SA(Q, K, V) = A(Q, K) \times V \quad (2)$$

### 2.1.2 Multi-Head Self-Attention

In multi-head self-attention (MSA), multiple SA blocks (heads) are performed in parallel to produce multiple output maps. The final output is a concatenation and projection of all outputs of SA blocks. This enables better modeling of complex dependencies between different elements in the input. For  $H$  number of heads, each head has its learnable weight matrices,  $\{W^{(Q_i)}, W^{(K_i)}, W^{(V_i)}; i = 1, \dots, H\}$ .

$$\begin{aligned} Z_i &= SA(Q_i, K_i, V_i) = \text{Softmax}\left(\frac{Q_i \times K_i^T}{\sqrt{D/H}}\right), \\ MSA(Q, K, V) &= \text{Concat}(Z_1, Z_2, \dots, Z_H)W^0 \end{aligned} \quad (3)$$

where  $W^0$  is a linear projection that aggregates the outputs of all attention heads. It is noteworthy that a larger number of heads does not necessarily lead to better performance [Dosovitskiy et al \(2020\)](#).



**Fig. 4** Overview of Vision Transformer (ViT) and the Transformer encoder.

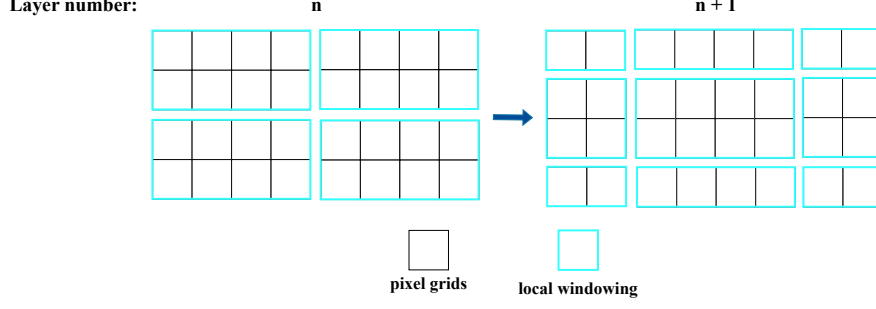
## 2.2 Vision Transformer Pipeline

The Vision Transformer [Dosovitskiy et al \(2020\)](#), or ViT, is a Transformer-like architecture introduced for image classification tasks. The main paradigm in the ViT is that tokens are created from the flattened patches of the image. Let  $X$  be a 3D image volume ( $X \in \mathbb{R}^{(H \times W \times L \times C)}$ ), where  $(H, W, L)$  represents the image dimensions, and  $C$  is the number of channels. The image is divided into  $N$  patches, which can overlap or not overlap, with each patch having a size of  $(P, P, P)$ . Then, a sequence is created from the flattened form of these patches  $x_P \in \mathbb{R}^{N \times P^3 C}$  and projected into a  $D$  dimensional space  $\hat{x}$ . To preserve positional information, a positional embedding was added, resulting in the input of the Transformer encoder, denoted as  $x$ :

$$x = \hat{x} + E_{pos}, E_{pos} \in \mathbb{R}^{N \times D} \quad (4)$$

The subsequent tokens were inputted into a Transformer encoder comprising  $L$  stacked base blocks. Each base block included multi-head self-attention and a multi-layer perceptron (MLP) with layer normalization (LN), and residual connections were employed following each block. A depiction of ViT and the Transformer encoder is shown in Fig. 4

It is noteworthy that the computational complexity of calculating the Softmax within the MSA blocks grows quadratically as the input sequence length increases [Dosovitskiy et al \(2020\)](#). This limitation could restrict its practical use, especially when dealing with high-resolution medical images. The introduction of the "Shifted Windows" idea in the Swin Transformer [Liu et al \(2021\)](#) improved the efficiency of MSA calculations. Unlike ViT [Dosovitskiy et al \(2020\)](#), which computes the relationship between one token and all others in every step of the self-attention calculations, Swin Transformer restricts self-attention calculations to non-overlapping local windows. It also enables cross-window connections and maintains linear computational complexity relative to the image size. Refer to Figure 5 for a visual representation of how the Shifted Windows concept partitions an input feature map with dimensions of



**Fig. 5** The cyclic shift of the local window for Shifted Windows-based self-attention computation. The self-attention is computed in each local window.

4×8 pixels, using a window size of 2×4. Additionally, it utilizes a hierarchical structure and generates feature maps at multiple resolutions through the incorporation of patch-merging layers.

### 2.3 Adaptations for Medical Image Analysis

Adapting Transformers for medical image analysis involves several considerations. One significant challenge is the high dimensionality of medical images, often requiring substantial computational resources for processing, especially for 3D images. Researchers have explored techniques such as patch-based processing [Liu et al \(2021\)](#) and efficient attention mechanisms [Xiong et al \(2021\)](#); [Rao et al \(2021\)](#) to mitigate this challenge. Another significant challenge lies in ensuring the generalizability of the proposed network. The pipeline must demonstrate robustness when tested on unseen data acquired from different imaging scanners or centers. The inherent variability among images obtained from different vendors, even when imaging the same subject, can result in a noticeable reduction in performance accuracy. Addressing and managing this variability is essential for handling the challenges associated with generalization. Various applications require specific conditions to be satisfied, and these conditions can vary significantly between different applications. The design of networks, especially Transformer-based architectures, should be approached with careful consideration based on the unique nature and requirements of each application.

Several architectures have been developed utilizing exclusively Transformer models [Cao et al \(2022\)](#); [Karimi et al \(2021\)](#); [Wang et al \(2021\)](#). DAE-Former [Azad et al \(2023\)](#) is a dedicated Transformer architecture proposed for medical image segmentation, featuring dual attention blocks in both the encoder and decoder, along with cross-attention blocks in the skip connections to optimize segmentation results. The key elements of dual attention blocks include efficient attention [Shen et al \(2021\)](#), employed to reduce computational complexity from quadratic to linear. Additionally, transpose attention [Ali et al \(2021\)](#) is incorporated into these blocks to capture channel attention. This architectural choice is based on empirical evidence suggesting that combining spatial and channel attention enhances the model’s capacity to capture more contextual features [Guo et al \(2022\)](#).

Pure Transformer architectures demonstrate certain limitations when compared with hybrid architectures that effectively capture both local and global information. Combining Transformers with CNNs in hybrid architectures leverages the strengths of both models, allowing Transformers to capture global context and CNNs to learn local features [Chen et al \(2021\)](#). These adaptations have paved the way for the successful application of Transformers in tasks like stroke segmentation, addressing the unique requirements of medical image analysis. Hybrid Transformer-CNN models offer flexibility in the placement of the Transformer component.

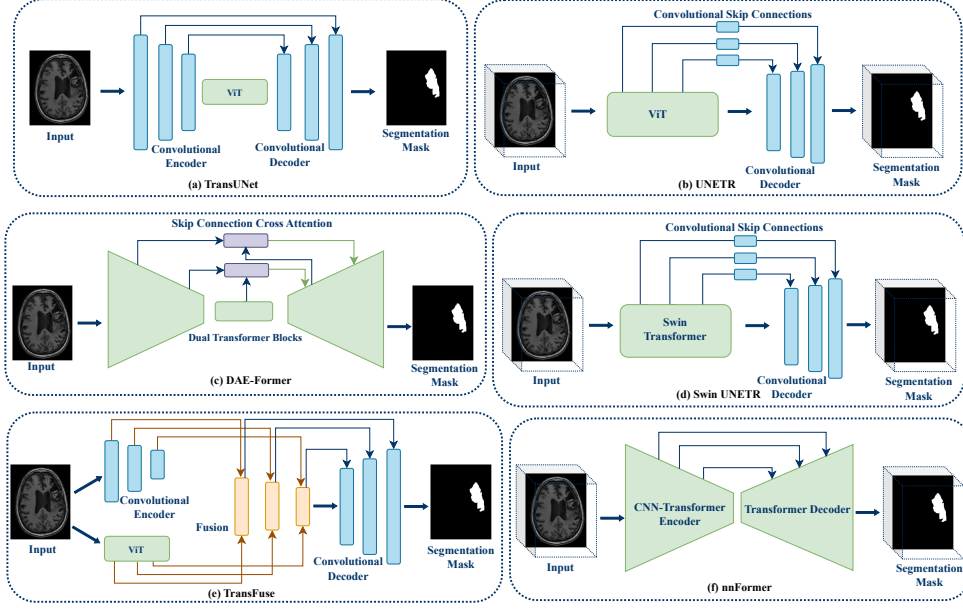
In Swin UNETR model [Hatamizadeh et al \(2021\)](#), the Swin Transformer took on the role of the encoder, and the encoded features from the Transformer were combined with the CNN-based decoder at various levels and resolutions. UNETR [Hatamizadeh et al \(2022\)](#) consisted of a Transformer-based encoder and a CNN-based decoder, featuring skip connections composed of convolutional-based blocks. In TransUNET model [Chen et al \(2021\)](#), the Transformer was integrated into the encoder, where it processed tokenized image patches derived from a CNN-generated feature map. TransFuse [Zhang et al \(2021\)](#) introduced a unique approach, employing a dual-branch encoder, one branch based on CNN and the other solely on the Transformer. A novel technique, called BiFusion, fused multi-level features extracted from both branches. In nnFormer [Zhou et al \(2021\)](#), a combination of interleaved convolution and self-attention operations was employed. Additionally, nnFormer utilized skip attention, akin to the traditional concatenation approach seen in skip connections within UNet-like architectures. Transformers have proven their effectiveness when utilized as the upsampling components within the decoder section [Li et al \(2022\)](#).

The Fully Convolutional Transformer (FCT) [Tragakis et al \(2023\)](#) was conceived to harness the strengths of CNNs for local feature representation and capitalize on Transformers’ proficiency in capturing long-range dependencies. The utilization of depth-wise convolutions in the projection layer obviates the necessity for positional embedding addition. Following the extraction of overlapping patches from an image, patch-based embeddings are incorporated, and MSA is subsequently calculated on these patches. In the FCT framework, a multi-branch convolutional paradigm is embraced to enhance spatial context. In this context, one layer applies a spatial convolution to the MSA output with a small kernel size, while other layers employ dilated convolutions with larger receptive fields. The integration of these outputs is facilitated by a fusion module known as Wide-Focus. Figure 6 illustrates several typical Transformer-based architectures for medical image segmentation, which have served as inspiration for many related models.

### 3 Datasets

In this section, we introduced available datasets for stroke infarct segmentation, encompassing both ischemic and hemorrhagic strokes across both CT and MRI modalities. Each dataset comprises various modalities with a different number of cases. Table 1 provides a summary of these datasets.





**Fig. 6** Transformer-based networks as segmentation model architectures. The schematic indication of whether the network is designed for a 2D or 3D image is presented in each figure. (a) TransUNet [Chen et al \(2021\)](#) incorporates Transformer blocks as additional encoders to model bottleneck features. (b) UNETR [Hatamizadeh et al \(2022\)](#) employs the Transformer as the primary encoder path, combining it with a CNN decoder and skip connections to create a hybrid model. (c) DAE-Former [Azad et al \(2023\)](#) integrates Transformers into dual attention blocks in both the encoder and decoder, enhancing the architecture with cross-attention in skip connections. (d) Swin UNETR [Hatamizadeh et al \(2021\)](#) features the Swin Transformer in the encoder section, complemented by a CNN-based decoder and skip connections to form a hybrid model. (e) TransFuse [Zhang et al \(2021\)](#) fuses the Transformer and CNN encoder to connect the decoder. (f) nnFormer [Zhou et al \(2021\)](#) incorporates Transformers in both the encoder and decoder for its architecture.

### 3.1 ISLES Dataset

The Ischemic Stroke Lesion Segmentation (ISLES) challenge offers publicly available datasets, released in 2015, 2017, 2018, and 2022. The objectives of the challenge varied each year, and distinct image modalities were provided in each one.

#### 3.1.1 ISLES 2015

ISLES 2015 [Maier et al \(2017\)](#) is a publicly available dataset comprising two distinct sub-challenges: Sub-Acute Ischemic Stroke Lesion Segmentation (SISS) and Stroke Perfusion Estimation (SPES).

ISLES 2015-SISS consisted of 64 sub-acute ischemic cases, with 28 cases allocated for training and 36 cases for testing with a voxel size of  $1mm^3$ . These cases were collected from two different medical centers with variations in image resolution. Each case within the dataset was accompanied by four MRI modalities, namely T1-weighted, T2-weighted, DWI, and Fluid attenuated inversion recovery (FLAIR). Preprocessing

steps, including skull-stripping and resampling to an isotropic space, were applied, and all modalities were registered to the FLAIR modality as a reference. Both the training and testing datasets included instances of single and multi-focal lesions, as well as large and small lesions.

ISLES 2015-SPES comprised 50 acute ischemic cases, with 30 cases designated for training and 20 cases for testing, with a voxel size of  $2mm^3$ . Each case was accompanied by seven distinct image modalities, namely T1 contrast-enhanced (T1c), T2, DWI, cerebral blood flow (CBF), cerebral blood volume (CBV), time-to-peak (TTP), and time-to-max (Tmax). To ensure consistency, all modalities were registered on the T1c modality.

### 3.1.2 ISLES 2017

This dataset [Winzeck et al \(2018\)](#) represents an extension of the ISLES 2016 stroke lesion segmentation challenge, notable for its expansion in the number of acute ischemic cases. Specifically, the dataset increased from 35 training and 19 testing cases to 43 training and 32 testing cases. Furthermore, this dataset introduced a new set of MRI modalities, including ADC, rBF, rBV, MTT, Tmax, TTP, and raw PWI, distinguishing it from the previously utilized modalities in the 2015 version. The pre-processing steps in ISLES 2017 were more concise, primarily encompassing registration and skull-stripping, while the voxel size and image resolutions exhibited variations.

### 3.1.3 ISLES 2018

The primary objective of the ISLES 2018 dataset [Cereda et al \(2016\)](#); [Hakim et al \(2021\)](#) was to perform the segmentation of stroke lesions using computed tomography perfusion (CTP) images, guided by annotations derived from DWI images, which are considered the standard image modalities. The dataset encompasses information from 103 acute ischemic cases with MRI images acquired within a 3-hour window of CTP. For training, 63 cases were designated, while the remaining 40 were reserved for testing. The input data for the algorithms consisted of various perfusion maps, including cerebral blood volume (CBV), cerebral blood flow (CBF), MTT, and Tmax.

### 3.1.4 ISLES 2022

The dataset [Hernandez Petzsche et al \(2022\)](#), sourced from three distinct medical centers, comprised information from 400 acute and sub-acute ischemic cases. Notably, this dataset featured a diverse range of infarct patterns, and a high degree of variability in terms of lesion size and location, with a mean number of 9.289 and a maximum of 126 unconnected ischemic regions per scan. The dataset exhibited heterogeneity due to the utilization of three different imaging devices, which serves as a valuable criterion for assessing the generalization of proposed methods. Modalities included in this dataset encompass DWI, ADC, and FLAIR images.

## 3.2 ATLAS Dataset

The Anatomical Tracings of Lesions After Stroke (ATLAS) v2.0 dataset [Liew et al \(2022\)](#) contained high-resolution T1-weighted images for the segmentation of acute,

sub-acute, and chronic stroke lesions. Significantly, this dataset was more extensive, containing more than four times the data volume of its predecessor, ATLAS v1.2. Aggregated from multiple centers worldwide, the dataset included data from 1,271 cases. Of these, 655 cases were allocated for training, 300 cases offered images only with hidden segmentation masks, and an additional 316 cases were entirely withheld to assess the generalizability of the proposed methods. The preprocessing steps applied to this dataset involved intensity normalization and registration on the MNI-152 template.

### 3.3 AISD Dataset

The Acute ischemic stroke dataset (AISD) [Liang et al \(2021\)](#) comprised paired CT-MRI data for 397 acute ischemic stroke cases. The dataset included Non-Contrast-enhanced CT (NCCT) scans and DWI scans, which were acquired within 24 hours of the CT images. The segmentation labels were derived from the MRI scans, which served as the standard for this purpose.

### 3.4 APIS Dataset

The APIS dataset [Gómez et al \(2023\)](#) was designed as a paired CT-MRI dataset with the objective of ischemic stroke lesion segmentation, utilizing NCCT images and annotations from ADC scans. The training set comprised 60 pairs of CT-MRI data, while the testing phase involved 36 NCCT scans exclusively. All cases underwent preprocessing steps, including skull-stripping and registration onto the ADC scans, to ensure data consistency and alignment.

### 3.5 Johns Hopkins University’s Dataset

This dataset comprised 2888 MRI datasets from cases involving acute and early sub-acute stroke patients, along with corresponding annotations [Liu et al \(2023a\)](#). The data were collected over 10 years using 11 MRI scanners. For all patients, DWI images, B0, and ADC were provided, and nearly 98.8% of patients had additional MRI modalities, including T1, high-resolution T1 MPRAGE, T2, FLAIR, SWI, and PWI. DWI images were registered onto the standard MNI space, subjected to skull-stripping, and resampled to a voxel size of  $1mm^3$ . The considerable diversity within this dataset renders it an excellent benchmark for evaluating proposed methods in the context of stroke segmentation. Nevertheless, it is important to note that access to this dataset is subject to certain restrictions.

### 3.6 Intracranial Hemorrhage Segmentation (IHS) Dataset

The dataset comprised non-contrast CT scans from 36 patients diagnosed with various types of intracranial hemorrhage, including intraventricular, intraparenchymal, sub-arachnoid, epidural, and subdural hemorrhage [Hssayeni et al \(2020\)](#). Each scan was characterized by an average of 30 slices with a thickness of  $5mm$ . Annotations for the dataset were provided by two radiologists. Collected in 2018, this dataset is accessible from PhysioNet [Goldberger et al \(2000\)](#) with certain restrictions.

**Table 1** Summary of available datasets information for stroke segmentation

Dataset	Disease	Number of Cases	Target Modality
ISLES 2015-SPES	Acute ischemic stroke	50	MRI
ISLES 2015-SISS	Subacute ischemic stroke	64	MRI
ISLES 2017	Acute ischemic stroke	75	MRI
ISLES 2022	Acute and subacute ischemic stroke	400	MRI
ATLAS v2.0	Acute, subacute, and chronic ischemic stroke	1271	MRI
Johns Hopkins University	Acute and early subacute ischemic stroke	2888	MRI
ISLES 2018	Acute ischemic stroke	103	CT
AISD	Acute ischemic stroke	397	CT
APIS	Ischemic stroke	96	CT
IHS	Hemorrhagic stroke	36	CT
INSTANCE	Hemorrhagic stroke	200	CT

### 3.7 INSTANCE Dataset

For the Intracranial Hemorrhage Segmentation on Non-Contrast Head CT (INSTANCE) challenge [Li et al \(2023b\)](#), a dataset comprising non-contrast CT scans from 200 patients diagnosed with various types of intracranial hemorrhage was assembled. The dataset allocation for different phases was as follows: 100 scans were designated for the training phase, 30 cases without ground truth labels were set aside for validation, and the remaining 70 cases were reserved for the final evaluation. The image size for each slice was  $512 \times 512$ , with the number of slices varying from 20 to 70 for each case. While the pixel size in each slice was  $0.42mm^2$ , providing a good inter-slice resolution, the slice thickness was  $5mm$ , resulting in a lower inter-slice resolution.

## 4 Performance Evaluation for Stroke Segmentation

Quantitative analysis of a segmentation process, evaluating its effectiveness in categorizing pixels or voxels into desired classes, is a crucial component of model evaluation. Many commonly used metrics rely on pixel-wise or voxel-wise calculations. The simplest method for assessing performance is through overall accuracy, defined as:

$$OverallAccuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

where, TP, TN, FP, and FN represent true positives, true negatives, false positives, and false negatives, respectively. However, overall accuracy may not provide sufficient insights, especially in imbalanced tasks such as stroke segmentation. To address this limitation, alternative metrics are widely used for a more nuanced evaluation of imbalanced semantic segmentation performance. One prominent metric is the Dice Similarity Coefficient (DSC), which measures the overlap between the predicted segmentation and the ground truth, ranging from 0 (no overlap) to 1 (perfect overlap). It is calculated as:

$$DSC = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (6)$$

Another valuable metric is the Intersection over Union (IoU), also known as the Jaccard Index, which assesses the ratio of the overlap to the total combined area,

ranging from 0 (no similarity) to 1 (perfect similarity). The IoU is computed as:

$$IoU = \frac{TP}{TP + FP + FN} \quad (7)$$

Both DSC and IoU provide comprehensive insights into the agreement between the predicted segmentation and the ground truth, making them particularly useful in scenarios with an imbalanced class distribution. There are also other commonly used metrics to measure stroke segmentation performance, including Precision, which assesses the accuracy of the positive predictions, Recall/Sensitivity, which gauges the ability to capture positive instances, and F1-score, a metric that strikes a balance between precision and recall by calculating their harmonic mean.

$$Precision = \frac{TP}{TP + FP} \quad (8)$$

$$Recall = \frac{TP}{TP + FN} \quad (9)$$

$$F1 - score = \frac{2 \times Recall \times Precision}{Recall + Precision} \quad (10)$$

F1-score is a specific form of the general  $F_\beta - score$ , in which the parameter  $\beta$  controls the trade-off between recall and precision, particularly useful when there is uneven importance assigned to precision and recall. The formula for the  $F_\beta - score$  is as follows:

$$F_\beta - score = \frac{(1 + \beta^2) \times Recall \times Precision}{Recall + \beta^2 \times Precision} \quad (11)$$

Another metric worth considering is the Hausdorff Distance (HD), which measures the maximum distance between two segmentation sets. Utilizing the HD metric reflects the level of dissimilarity between predicted and ground truth boundaries. The HD is computed as follows:

$$HD(A, B) = \max \left( \sup_{a \in A} \inf_{b \in B} d(a, b), \sup_{b \in B} \inf_{a \in A} d(a, b) \right) \quad (12)$$

Here,  $A$  and  $B$  represent two sets, and  $d(a, b)$  is the distance function between points  $a$  in set  $A$  and  $b$  in set  $B$ . The formula calculates the HD by determining the maximum of the infimum of distances from points in set  $A$  to the nearest point in set  $B$  and vice versa. This measurement can be represented in mm or voxel/pixel-based units. Some papers also utilized additional measurements, including Simple Lesion Count (SLC), Volume Difference (VD), Average Volume Difference (AVD), Volumetric Overlap Error (VOE), Relative Volume Difference (RVD), and average symmetric surface distance (ASSD).

## 5 Stroke Segmentation using Transformers

### 5.1 Earlier Approaches for Stroke Segmentation

The majority of prior studies of stroke lesion segmentation have primarily focused on CNN models, with an emphasis on U-Net-based architectures [Clerigues et al \(2019\)](#); [Basak and Rana \(2020\)](#); [Kheezrpour et al \(2022\)](#); [Liu et al \(2020\)](#); [Kadry et al \(2021\)](#). Notably, the U-Net [Ronneberger et al \(2015\)](#) architecture, specifically designed for biomedical image segmentation, exhibits a distinctive U-shaped configuration, comprising an encoder segment dedicated to contextual feature extraction and a decoder segment tailored for accurate localization. The incorporation of skip connections within this architecture facilitated the seamless integration of high-level feature maps derived from the encoder path with fine-grained details from the decoder path, enhancing its segmentation performance. Some research studies have drawn inspiration from the DenseNet architecture for their utilized neural network design [Zhang et al \(2018a\)](#). Table 2 summarizes the performance of selected CNN-based methods for stroke segmentation across various datasets. The inclusion of papers in the table is based on their demonstrated superior performance.

The bilateral quasi-symmetry property of the brain has been utilized in some research studies [Liang et al \(2021\)](#); [Wang et al \(2016\)](#); [Clèrigues et al \(2020\)](#); [Vuputuri et al \(2018\)](#). In [Clèrigues et al \(2020\)](#) a patch-based deep learning pipeline was proposed, wherein the extraction of patches was carried out with a significant degree of overlap. These patches were subsequently fed into the neural network using a well-balanced sampling strategy, to mitigate issues associated with class imbalance. In [Praveen et al \(2018\)](#) a stacked sparse autoencoder network for unsupervised feature learning was employed, which was subsequently coupled with a Support Vector Machine (SVM) classifier to classify patches into normal and lesion categories.

The D-UNet [Zhou et al \(2019b\)](#) model was introduced for chronic stroke segmentation, and it incorporated a fusion of 2D and 3D convolutions in the encoder stage via a dimension transform block. This combination of 2D and 3D information facilitated more effective lesion identification. Additionally, a novel loss function called Enhance Mixing Loss (EML) was employed, which is a composite of the Focal loss [Lin et al \(2017\)](#) and Dice coefficient loss. By evaluating their proposed method on the ATLAS dataset, they achieved a mean Dice coefficient of 0.5349 for pixel-wise calculation and 0.7231 for voxel-wise calculation, respectively. In [Kumar et al \(2020\)](#) a hybrid Classifier-Segmenter network (CSNet) was introduced. Initially, the images were input into a classifier that distinguished slices with lesions. The chosen slices were then fed into a segmenter network, which utilized a fractal U-Net model for segmentation, and a final voting mechanism was employed to improve segmentation performance.

In [Abulnaga and Rubin \(2019\)](#) a CNN-based model was introduced that incorporated the pyramid pooling module, as introduced in PSPNet [Zhao et al \(2017\)](#). This module was utilized to extract global and local contextual information, enhancing the accuracy of stroke lesion segmentation. It achieves this by capturing global information through the use of varying kernel sizes and aggregating multi-scale region-based context. The best results were achieved by incorporating pretraining and utilizing the Focal Loss on the ISLES 2018 dataset.

In X-Net [Qi et al \(2019\)](#), a Feature Similarity Module (FSM) was implemented to capture long-range dependencies, thereby enhancing the segmentation process. This module was employed at the bottleneck between the encoder and decoder to investigate dense long-range contextual information. To reduce network size and control the number of parameters, mitigating the risk of overfitting, depthwise convolutions were also integrated. Their proposed network demonstrated superior performance compared to other architectures, including U-Net, SegNet [Badrinarayanan et al \(2017\)](#), PSPNet [Zhao et al \(2017\)](#), ResUNet [Zhang et al \(2018b\)](#), 2D Dense-UNet [Li et al \(2018\)](#), and DeepLabv3+ [Chen et al \(2018\)](#), as evaluated on the ATLAS dataset.

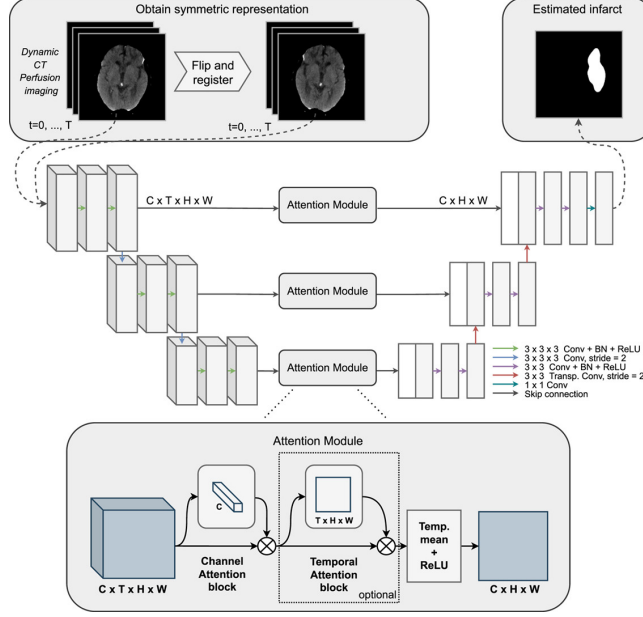
In [Liu et al \(2019a\)](#) a CNN-based pipeline was proposed that divides the network into two subnetworks. This approach involved using multi-kernels of various sizes to extract feature maps across different receptive fields. Post-processing techniques were also applied to retain edge details in the images and reduce noise. The optimal performance for their proposed pipeline was achieved by employing a dropout rate of 0.1. [Bal et al \(2023\)](#) followed a similar approach and incorporated a local pathway and a global pathway within their proposed model, with larger kernel sizes employed in the global pathway to expand the receptive field for extracting long-range dependencies and global information. The best results in their proposed pipeline were obtained through the inclusion of preprocessing and data augmentation for the ISLES 2015 dataset.

[Zhang et al \(2020\)](#) devised a pipeline centered around a Detection and Segmentation Network (DSN). They utilized a triple-branch architecture to extract predictions for slices in the axial, sagittal, and coronal planes separately. Subsequently, the predicted labels from different slices within each plane were fine-tuned and passed through a fusion module to obtain the final segmentation label. Their approach outperformed architectures such as U-Net, V-Net [Milletari et al \(2016\)](#), and DeepMedic [Kamnitsas et al \(2016\)](#) in validation using the ISLES-SSIS dataset.

In [Huo et al \(2022\)](#) a model within the nnU-Net framework [Isensee et al \(2021\)](#) was introduced. Their approach incorporated four schemes: a generic U-Net, utilizing a TopK10 loss to improve performance in small lesion segmentation, a residual U-Net, and a self-training U-Net to enhance model diversity. An ensemble method was employed to combine the predicted results from these four networks, followed by post-processing techniques to enhance the results.

In [Abramova et al \(2021\)](#) a 3D U-Net-based network for hemorrhagic stroke segmentation in CT scans was utilized. To enhance the representation of informative features, they incorporated Squeeze and Excitation (SE) blocks [Hu et al \(2018\)](#) in the bottleneck and last layers of their network. Through symmetric data augmentation and the implementation of a restrictive patch sampling approach, their proposed architecture achieved a mean Dice coefficient of 0.86 on a clinical dataset consisting of 76 cases.

Pool-UNet [Liu et al \(2022c\)](#) incorporated SE blocks in a novel module called DSE-ResNet placed in the bottleneck. This module captures interdependencies between channels to provide the most informative features for the decoder. Additionally, they combined the Poolformer structure [Yu et al \(2022\)](#), a transformer-like structure utilizing pooling operations, with CNNs to capture both local and global information.



**Fig. 7** PerfU-Net architecture [de Vries et al \(2023\)](#).

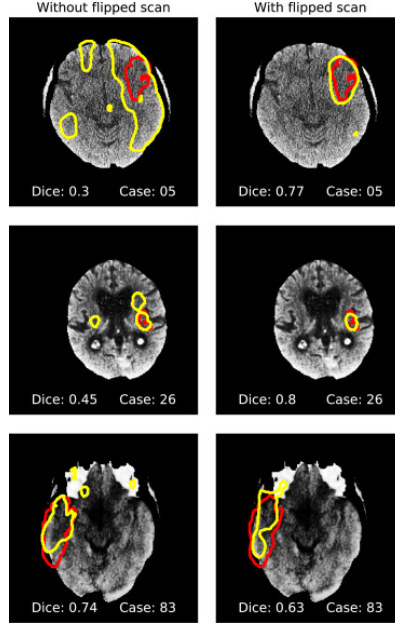
Evaluations on the ISLES 2018 dataset demonstrated the superior performance of their proposed architecture compared to architectures such as U-Net, R2UNet [Alom et al \(2018\)](#), and TransUNet. [Chalcroft et al \(2023\)](#) employed Large Kernel Attention [Guo et al \(2023\)](#) to capture long-range dependencies, capitalizing on the inherent biases of convolutions. The Large Kernel Attention mechanism comprises a sequence of depth-wise convolutions, dilated depth-wise convolutions, and pointwise convolutions.

PerfU-Net [de Vries et al \(2023\)](#) was introduced for stroke segmentation from CT images. This architecture incorporated an attention module placed in the skip connections, with two variations tested: one featuring channel attention and the other incorporating both channel attention and temporal attention. The training process utilized the generalized Dice loss [Sudre et al \(2017\)](#) as the loss function. PerfU-Net achieved a mean Dice coefficient of 0.564 when evaluated on the ISLES 2018 dataset, utilizing 32 frames as the input to the model and employing channel attention as the attention module. Refer to Figure 7 for an illustration of the proposed pipeline in PerfU-Net, and refer to Figure 8 for a qualitative analysis of performance using the ISLES 2018 dataset under two conditions: with and without flipped scans. In their observations, they noted that the use of flipped scans contributed to a reduction in the number of false positives.

## 5.2 Transformer-Based Architectures for Stroke Segmentation

Vision Transformers have been employed in recent years for stroke segmentation, leveraging their capabilities, especially when combined with CNNs to capture local and global information from the input data. Table 3 summarizes the performance of





**Fig. 8** Qualitative analysis of PerfU-Net performance for three cases from ISLES 2018 dataset with and without flipped representation [de Vries et al \(2023\)](#). Yellow lines indicate the predicted segmentation mask and red lines indicate the ground truth.

**Table 2** Performance Analysis of CNN-Based Approaches for Stroke Segmentation. \* indicates the median.

Reference	Year	Dataset	Performance				
			Dice	Sensitivity/Recall	Precision	HD	Others
<a href="#">Praveen et al (2018)</a>	2018	ISLES 2015-SISS training dataset	$0.943 \pm 0.057$	$0.924 \pm 0.072$	$0.968 \pm 0.074$	-	-
<a href="#">Liu et al (2019a)</a>	2019	ISLES 2015-SISS	$0.57 \pm 0.29$	-	-	$43.02 \pm 30.48$	ASSD: $8.22 \pm 16.54$
		ISLES 2015-SPES	$0.76 \pm 0.11$	-	-	$36.93 \pm 25.42$	ASSD: $1.79 \pm 0.54$
<a href="#">Clérigues et al (2020)</a>	2020	ISLES 2015-SISS	$0.59 \pm 0.31$	$0.60 \pm 0.30$	$0.65 \pm 0.35$	$34.7 \pm 28.9$	-
		ISLES 2015-SPES	$0.84 \pm 0.10$	$0.89 \pm 0.06$	$0.82 \pm 0.15$	$20.7 \pm 13.9$	-
<a href="#">Zhang et al (2020)</a>	2020	ISLES 2015-SISS	0.622	0.541	-	-	IoU: 0.4514
<a href="#">Bal et al (2023)</a>	2023	ISLES 2015-SISS	$0.87 \pm 0.10$	$0.86 \pm 0.09$	$0.88 \pm 0.10$	-	-
		ISLES 2015-SPES	$0.90 \pm 0.08$	$0.89 \pm 0.08$	$0.90 \pm 0.07$	-	-
<a href="#">Pereira et al (2019)</a>	2019	ISLES 2015-SPES	$0.82 \pm 0.09$	-	-	-	ASSD: $1.27 \pm 0.72$
		ISLES 2017	$0.34 \pm 0.20$	$0.55 \pm 0.30$	$0.36 \pm 0.25$	-	-
<a href="#">Islam and Ren (2018)</a>	2018	ISLES 2017	$0.29 \pm 0.2$	$0.6 \pm 0.24$	$0.35 \pm 0.23$	$46.9 \pm 12.8$	-
<a href="#">Lucas et al (2018)</a>	2018	ISLES 2017	0.35	0.35	0.52	21.48	ASSD: 3.45
<a href="#">Hu et al (2020)</a>	2020	ISLES 2017	$0.30 \pm 0.22$	$0.43 \pm 0.27$	$0.35 \pm 0.27$	-	-
<a href="#">Abulnaga and Rubin (2019)</a>	2019	ISLES 2018	$0.54 \pm 0.09$	-	-	-	-
<a href="#">Clérigues et al (2019)</a>	2019	ISLES 2018 testing dataset	$0.49 \pm 0.31$	$0.57 \pm 0.35$	$0.51 \pm 0.36$	$11.3 \pm 31.6$	-
<a href="#">Rubin and Abulnaga (2019)</a>	2019	ISLES 2018	$0.54 \pm 0.23$	$0.63 \pm 0.25$	$0.56 \pm 0.25$	$27.88 \pm 21.00$	-
<a href="#">Liu et al (2022c)</a>	2022	ISLES 2018	$0.565 \pm 0.205$	$0.565 \pm 0.218$	$0.678 \pm 0.223$	$21.14 \pm 13.61$	-
<a href="#">de Vries et al (2023)</a>	2023	ISLES 2018	$0.564 \pm 0.009$	$0.644 \pm 0.008$	$0.565 \pm 0.019$	$22.3 \pm 1.8$	AVD: $9.9 \pm 0.4$
<a href="#">Qi et al (2019)</a>	2019	ATLAS v1.2	0.4867	0.4752	0.6000	-	IoU: 0.3723
<a href="#">Zhou et al (2019b)</a>	2019	ATLAS v1.2	$0.534 \pm 0.276$	$0.524 \pm 0.291$	$0.6331 \pm 0.2958$	-	-
<a href="#">Liu et al (2019b)</a>	2019	ATLAS v1.2	0.5578	0.8291*	-	-	VOE: 0.1403*
<a href="#">Yang et al (2019)</a>	2019	ATLAS v1.2	0.581	0.581	0.649	-	VOE: 54.6 RVD: 25.4
<a href="#">Basak and Rana (2020)</a>	2020	ATLAS v1.2	$0.535 \pm 0.276$	$0.523 \pm 0.293$	$0.634 \pm 0.287$	-	-
<a href="#">Hui et al (2020)</a>	2020	ATLAS v1.2	0.593	0.62	0.691	-	-
<a href="#">Tomita et al (2020)</a>	2020	ATLAS v1.2	0.64	-	0.62	20.4	ASSD: 3.6
<a href="#">Yu et al (2023b)</a>	2023	ATLAS v1.2	$0.559 \pm 0.180$	$0.576 \pm 0.162$	-	-	F1: $0.545 \pm 0.162$
<a href="#">Yu et al (2023a)</a>	2023	ATLAS v1.2	$0.571 \pm 0.195$	$0.597 \pm 0.158$	-	-	F1: $0.562 \pm 0.192$
<a href="#">Huo et al (2022)</a>	2022	ATLAS v2.0 training dataset	$0.646 \pm 0.270$	-	-	$21.51 \pm 26.82$	VD: $5729 \pm 11565$ SLC: $3.382 \pm 6.786$
<a href="#">Chalcroft et al (2023)</a>	2023	ISLES 2022	0.693	-	-	-	F1: 0.657
		ATLAS v2.0	0.678*	-	-	-	F1: 0.474*
<a href="#">Liang et al (2021)</a>	2021	AISD	0.5784	0.5880	0.6597	-	F1: 0.6218
<a href="#">Ni et al (2022)</a>	2022	AISD	0.5245	-	-	39.18	-

selected Transformer-based methods for stroke segmentation across various datasets. [de Vries et al \(2021\)](#) proposed a hybrid Transformer-CNN pipeline for ischemic stroke infarct segmentation from CT perfusion scans. They considered the axial slices of 3D images as temporal information and incorporated the flipped and registered form of each slice as an additional channel in the input to exploit the brain’s bilateral quasi-symmetry property. The spatio-temporal data were fed into a Transformer block consisting of the Linformer [Wang et al \(2020\)](#) backbone to generate an attention map representing the probability of infarction. Subsequently, this attention map, along with the source data, was input into a traditional U-Net to produce the final segmentation. The Transformer part was trained using cross-entropy loss, while the U-Net was trained using the generalized Dice loss function.

UCATR [Luo et al \(2021\)](#) was proposed for acute ischemic stroke segmentation from non-contrast CT images. A Transformer-based block succeeded the CNN encoder in the bottleneck, and irrelevant information was filtered by employing Multi-Head Cross-Attention modules in the skip connections. The proposed network was evaluated on a clinical dataset containing information from 11 patients, averaging 95 slices per patient, and it outperformed U-Net, Attention U-Net, and TransUNet by achieving a mean Dice coefficient of 0.7358. UTransNet [Feng et al \(2022\)](#) employed a novel module (CT block) consisting of two convolutional layers and a Transformer module to leverage the advantages of both CNNs and Transformers. To mitigate computational complexity within the self-attention mechanism, the PVT v2 Transformer [Wang et al \(2022c\)](#) was incorporated. In evaluations conducted on the ATLAS dataset, UTransNet achieved superior results compared to other Transformer-based methods, including TransUNet, SwinUNet [Cao et al \(2022\)](#), and UCTransNet [Wang et al \(2022a\)](#).

The Multi-Encoder Transformer (METrans) [Wang et al \(2022b\)](#) was proposed as a novel architecture, incorporating a methodology involving additional encoding modules. These modules served the purpose of extracting abstract features at distinct stages of the primary encoder path. The ensuing step involved the fusion of the multi-scale extracted features. Following each convolutional module within the encoder, Convolutional Block Attention Modules (CBAM) [Woo et al \(2018\)](#) were employed to harness both channel attention and spatial-channel attention. Furthermore, a Transformer-based block was integrated into the bottleneck to facilitate the extraction of global features.

Within the architecture of STHarDNet [Gu et al \(2022\)](#), HarDNet blocks [Chao et al \(2019\)](#) were employed in both the encoder and decoder paths. Notably, the Swin Transformer was incorporated exclusively in the initial layer of the skip connection, while the subsequent layers adhered to a conventional, plain structure. The performance of this network surpassed that of numerous CNN-based and Transformer-based counterparts when evaluated on the ATLAS dataset. LLRHNet [Liu et al \(2022b\)](#) implemented a dual-path approach for feature encoding, wherein the initial path utilized CNN layers to extract local information, and the subsequent path employed a Transformer-based block for encoding global features. The features extracted from these two paths were concatenated, and a final prediction was generated through a CNN decoder. To enhance information transfer from the CNN encoder to the decoder,

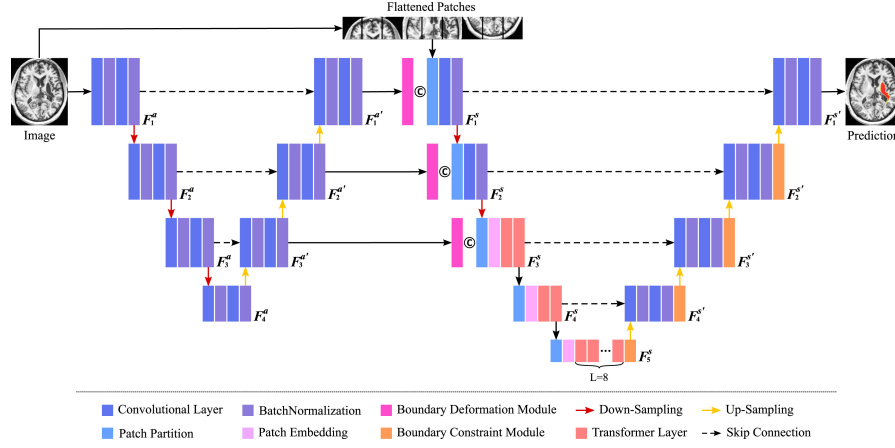
the model integrated multi-level feature fusion skip connections, a departure from conventional skip connection methods. Evaluation of the LLRHNet on a clinical dataset for ischemic stroke segmentation demonstrated its superior performance by achieving a mean Dice coefficient of 0.791.

In [Wu et al \(2022\)](#) an architecture consisting of three main elements was proposed. First, the Patch Partition Block (PPB) was employed to encode the image as a patch sequence, simultaneously reducing the number of parameters. Second, the Multi-scale Long-Range Interactive and Regional Attention (MLiRA) mechanism served as the encoder, comprising multiple subsampling Transformers (STR) followed by convolutional blocks. Within STR, subsampling Multi-head Interactive Self-Attention mechanisms were utilized to capture dimensional interactive attention. Moreover, STR exhibited flexibility in adjusting input resolution to attain global information at various spatial resolutions. Third, the Feature Interpolation Path (FIP) was utilized as the decoder, facilitating the recovery of encoded features to the original image resolution.

In [Marcus et al \(2023\)](#) a multi-task Transformer-based network for age estimation and segmentation of ischemic lesions from CT images was proposed. Their architecture was rooted in the DETR architecture [Carion et al \(2020\)](#) with certain modifications. The primary components of the proposed network included: 1) a CNN encoder consisting of four ResNeXt [Xie et al \(2017\)](#) blocks to generate an activation map; 2) a Transformer encoder-decoder, commencing with a pyramid pooling module [Zhao et al \(2017\)](#) to augment the receptive field, followed by a Transformer block using the gated positional self-attention mechanism [d’Ascoli et al \(2021\)](#); 3) heads for age estimation and bounding box predictions; and 4) a CNN decoder serving as the segmentation head. Based on evaluations of their proposed pipeline on a large clinical dataset consisting of 776 CT images collected from two medical centers, they reached a mean Dice coefficient of 0.382. For evaluation of the generalizability of the trained network on unseen data, they also utilized the ISLES 2018 dataset as the test dataset and reached a 0.203 mean Dice coefficient.

[Zhang and Chen \(2023\)](#) introduced a pipeline designed to address the efficient processing of 3D image data to preserve volumetric information. Their proposed architecture, named PDSwin, leverages the Swin Transformer with a pyramidal downsampling approach, spatially downsampling 2D slices. Additionally, the authors addressed the shift domain issue arising from diverse image acquisition centers by proposing a cluster-based domain adversarial algorithm. In [Xu and Ding \(2023\)](#) a U-shaped network was employed to segment stroke infarct from CT scans. In the CNN-based encoder path, multiple CBAM blocks were incorporated. The features extracted from two scales of the encoder path were flattened and inputted into a Transformer module, encompassing several deformable Transformer layers utilizing deformable self-attention [Zhu et al \(2020\)](#).

[Soh et al \(2023\)](#) introduced a Hybrid UNet and Transformer (HUT) network, comprising two parallel stages: a UNet and a Transformer block. The input to the Transformer block consisted of intermediate features extracted by the CNN encoder of the UNet. The output of the Transformer block was then fused with the extracted features from the encoder path in the skip connections at two different scales. In SAMIHS [Wang et al \(2023\)](#) a parameter-efficient fine-tuning strategy to the Segment

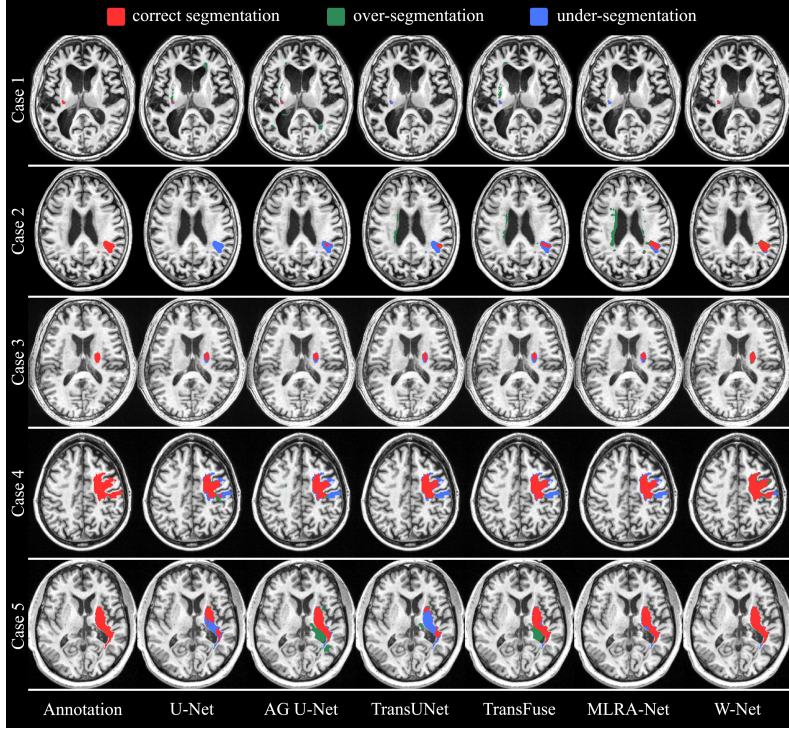


**Fig. 9** W-Net architecture [Wu et al \(2023a\)](#).

Anything Model (SAM) [Kirillov et al \(2023\)](#) model was applied to segment hemorrhagic stroke. To improve segmentation results, they utilized a combination of the binary cross-entropy loss and a boundary-sensitive loss.

Some research endeavours have been pursued to improve the realism of predicted segmentation mask boundaries in stroke nature. One such approach was TransRender [Wu et al \(2023b\)](#), which was proposed to address the issue of overly smooth boundaries. It achieved this by adaptively selecting specific points for computing the boundary features in a point-based rendering manner, intending to enhance the fidelity of the boundary estimation. The hierarchical Transformer-based encoder path captured global information across multiple scales, with additional parallel CNN blocks employed to capture local information. Both local and global features were then provided as input to multiple render modules. These modules, by selecting specific uncertain points and extracting feature representations for those points, facilitated the re-prediction of these uncertain points as boundary points.

Another approach to enhance boundary estimation accuracy was W-Net [Wu et al \(2023a\)](#), which introduced a Boundary Deformation Module (BDM) and a Boundary Constraint Module (BCM) to address fuzzy boundaries. W-Net integrated both a CNN network and a Transformer-based network as backbone networks. Initially, a U-shaped CNN network was employed for coarse segmentation, leveraging the advantages of CNNs to extract local features. In the decoder path, features of different scales were inputted into proposed BDM blocks for further optimization through iterative boundary deformation, correcting the initially predicted boundaries using circular convolutions. The second stage of W-Net consisted of a Transformer-based U-shaped network. The output of the BDM blocks was fused with the encoder features at multiple levels. The decoder utilized BCM blocks to refine the encoded global features by constraining the boundary curves, employing several parallel dilated convolution layers. Refer to Figure 9 for an illustration of the W-Net architecture, and Figures 10 and 11 for qualitative analyses on the ATLAS and ISLES 2022 datasets, respectively.



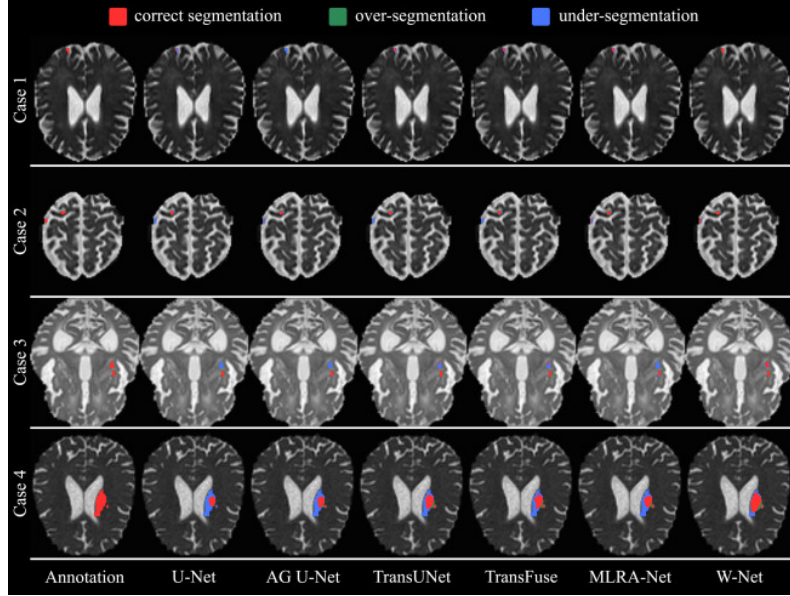
**Fig. 10** Qualitative analysis of the W-Net performance compared to other five benchmarks using ATLAS dataset [Wu et al \(2023a\)](#).

**Table 3** Performance Analysis of Transformer-Based Approaches for Stroke Segmentation.

Reference	Year	Dataset	Performance				
			Dice	Sensitivity/Recall	Precision	HD	Others
<a href="#">de Vries et al (2021)</a>	2021	ISLES 2018 test dataset	0.42	0.53	0.44	-	-
<a href="#">Xu and Ding (2023)</a>	2023	ISLES 2018	0.4667	0.4724	0.5888	-	F1: 0.5242
<a href="#">Wang et al (2022b)</a>	2022	ISLES 2018	0.67	0.641	0.72	-	-
		ATLAS v1.2	0.931	0.91	0.94	-	-
<a href="#">Wu et al (2023b)</a>	2023	ISLES 2022	0.8537	0.8394	0.8648	27.60	F2: 0.8487
		ATLAS v1.2	0.5979	0.6808	0.6391	33.98	F2: 0.5938
<a href="#">Wu et al (2023a)</a>	2023	ISLES 2022	0.8560	0.8539	0.8834	27.34	F2: 0.8529
		ATLAS v1.2	0.6176	0.6868	0.6286	32.47	F2: 0.6460
<a href="#">Feng et al (2022)</a>	2022	ATLAS v1.2	0.8597	-	-	-	F1: 0.8494 IoU: 0.8251
<a href="#">Gu et al (2022)</a>	2022	ATLAS v1.2	0.5547	0.5286	0.6764	-	IoU: 0.4184
<a href="#">Wu et al (2022)</a>	2022	ATLAS v1.2	0.6119	0.6765	0.6330	13.49	F2: 0.6376
<a href="#">Soh et al (2023)</a>	2023	ATLAS v1.2	$0.737 \pm 0.127$	$0.706 \pm 0.153$	$0.825 \pm 0.172$	$10.335 \pm 10.074$	IoU: $0.598 \pm 0.114$
<a href="#">Zhang and Chen (2023)</a>	2023	ATLAS v2.0	0.6273	-	-	-	-
<a href="#">Wang et al (2023)</a>	2023	IHS	0.6977	-	-	3.31	-
		INSTANCE	0.7652	-	-	3.71	-

## 6 Open Challenges and Future Directions

Current solutions for stroke segmentation, whether they employ CNN networks or Transformer-based architectures, have shown less satisfactory results compared to tasks such as tumor segmentation [Ranjbarzadeh et al \(2023\)](#); [Liu et al \(2023b\)](#). The underlying cause of this suboptimal performance is attributed to various factors, including high variability in the location, number, size, and pattern of the infarct.



**Fig. 11** Qualitative analysis of the W-Net performance compared to other five benchmarks using ISLES 2022 dataset [Wu et al \(2023a\)](#).

Furthermore, the intensity differences resulting from the varied imaging vendors and stroke ages pose a significant challenge for automated algorithms. The proposed methods must effectively distinguish between healthy and infarcted regions of the brain while accommodating the diverse variability introduced by different medical imaging systems and inherent stroke features.

An additional crucial feature of proposed methods in stroke segmentation is their generalizability to unseen data from different vendors, making them applicable. Current methods often exhibit a lack of generalization, with testing on unseen data acquired from diverse centers leading to significantly lower performance. It is imperative to investigate and improve the robustness of these methods to handle unseen data effectively. Exploring domain adaptation methods could prove beneficial in achieving this objective.

Another inherent characteristic of stroke infarcts is their ability to affect various parts of the brain at different stages, exhibiting different sizes in different locations. It is crucial for segmentation methods to accurately handle multi-instance infarcts of varying sizes. Despite the high values of the Dice coefficient suggesting good overall performance, instance-wise measurements often yield lower results [Kofler et al \(2023\)](#). This discrepancy is primarily due to the neglect of small infarcts in the presence of larger ones when calculating commonly used metrics. The currently proposed architectures encounter challenges in detecting and segmenting small infarcts, resulting from information loss within the deeper layers of the networks designed to represent abstract features. Thoughtful improvements are necessary to adapt proposed pipelines for the segmentation of small infarcts.



Transformers are widely employed in medical image segmentation due to their ability to represent global information and capture long-range dependencies, providing a robust representation of shape-based features for segmentation. Given the high variability in stroke infarct shape, location, and pattern, the incorporation of texture information derived from CNNs appears to be more beneficial. Hybrid CNN-Transformer architectures address this challenge by combining texture-based and global information. However, a careful selection of how to employ Transformer blocks is necessary to fully leverage their advantages for the integration of texture information alongside global features to improve the effectiveness of the stroke segmentation methodologies.

A notable challenge in the effective application of Transformers stems from their dependence on large datasets, which becomes especially pronounced in the medical field. The limited availability of labeled data, coupled with difficulties in acquiring annotations and privacy considerations, imposes constraints on the accessibility of medical data. Although pre-training on alternative datasets, including natural images, might be considered, it is suboptimal due to the inherent domain shift. Medical images often have high dimensions, resulting in a large number of parameters that demand significant computational resources. The increased parameter count, combined with a limited dataset, can lead to overfitting issues. Various strategies, such as slicing 3D data from different anatomical planes or adopting patch-based data inputting, have been attempted previously. However, these approaches can result in the loss of valuable information. Exploring realistic data augmentation techniques and optimizing data input methods are crucial avenues to investigate in order to mitigate these challenges.

To establish effective pipelines for stroke segmentation, it is crucial to incorporate additional characteristics, such as privacy-preserving algorithms [Sheller et al \(2020\)](#); [Li et al \(2021, 2019\)](#), and embrace explainable artificial intelligence (XAI) [Alicioglu and Sun \(2022\)](#); [Mondal et al \(2021\)](#); [Singh et al \(2020\)](#) as integral components of trustworthy AI [Liu et al \(2022a\)](#). In recent years, researchers have dedicated their efforts to interpret deep learning-based models, which were previously considered black boxes. Understanding the decision-making process of Transformers can enhance predictions and facilitate their use in aiding decision-making for medical diagnoses. Additionally, exploring privacy-preserving algorithms is imperative. This investigation aims to provide a platform that can be utilized in different centers, allowing the sharing of knowledge derived from diverse datasets without compromising the privacy of medical information. This approach aims to continually enhance the performance of the underlying pipeline.

## 7 Discussion and Conclusion

In this paper, we present a comprehensive review of Transformer-based architectures for segmenting stroke infarcts from MRI and CT images. We begin by offering preliminary information on the concepts of self-attention, vision Transformers, and several benchmark Transformer-based networks designed for medical image segmentation. Following this, we delve into details about available datasets for stroke segmentation, encompassing both ischemic and hemorrhagic strokes for both MRI and CT modalities.

Subsequently, we discussed commonly used metrics for evaluating segmentation performance and conducted a literature review on stroke segmentation using deep learning methods. Given that a significant portion of previous research has been conducted using CNNs, we selectively extracted and summarized the key ideas with superior performance. Specifically focusing on Transformer-based architectures, we conducted a review of 15 papers, all employing hybrid CNN-Transformer architectures. For each paper, we offered a high-level abstraction of the core techniques utilized in these networks. Additionally, we presented comparison tables for quantitative evaluations of performance, considering both CNN-based and Transformer-based networks. Finally, we outlined the unsolved challenges associated with stroke segmentation and suggested potential avenues for future research directions.

## Acknowledgment

The findings presented in this paper have emerged from a project funded by the Qatar Japan Research Collaboration Research Program under grant number M-QJRC-2023-313. The authors extend their sincere gratitude to Marubeni and Qatar University for their consistent and generous support.

## References

- Abramova V, Clerigues A, Quiles A, et al (2021) Hemorrhagic stroke lesion segmentation using a 3d u-net with squeeze-and-excitation blocks. *Computerized Medical Imaging and Graphics* 90:101908
- Abulnaga SM, Rubin J (2019) Ischemic stroke lesion segmentation in ct perfusion scans using pyramid pooling and focal loss. In: *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Workshop, BrainLes 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Revised Selected Papers, Part I 4*, Springer, pp 352–363
- Ali A, Touvron H, Caron M, et al (2021) Xcit: Cross-covariance image transformers. *Advances in neural information processing systems* 34:20014–20027
- Alicioglu G, Sun B (2022) A survey of visual analytics for explainable artificial intelligence methods. *Computers & Graphics* 102:502–520
- Alom MZ, Hasan M, Yakopcic C, et al (2018) Recurrent residual convolutional neural network based on u-net (r2u-net) for medical image segmentation. *arXiv preprint arXiv:180206955*
- Aoki J, Kimura K, Iguchi Y, et al (2010) Flair can estimate the onset time in acute ischemic stroke patients. *Journal of the neurological sciences* 293(1-2):39–44
- Azad R, Arimond R, Aghdam EK, et al (2023) Dae-former: Dual attention-guided efficient transformer for medical image segmentation. In: *International Workshop on PRedictive Intelligence In MEdicine*, Springer, pp 83–95



- Badrinarayanan V, Kendall A, Cipolla R (2017) Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39(12):2481–2495
- Bal A, Banerjee M, Chaki R, et al (2023) A robust ischemic stroke lesion segmentation technique using two-pathway 3d deep neural network in mr images. *Multimedia Tools and Applications* pp 1–40
- Balakrishnan G, Zhao A, Sabuncu MR, et al (2019) Voxelmorph: a learning framework for deformable medical image registration. *IEEE transactions on medical imaging* 38(8):1788–1800
- Basak H, Rana A (2020) F-unet: A modified u-net architecture for segmentation of stroke lesion. In: *International Conference on Computer Vision and Image Processing*, Springer, pp 32–43
- Cao H, Wang Y, Chen J, et al (2022) Swin-unet: Unet-like pure transformer for medical image segmentation. In: *European conference on computer vision*, Springer, pp 205–218
- Carion N, Massa F, Synnaeve G, et al (2020) End-to-end object detection with transformers. In: *European conference on computer vision*, Springer, pp 213–229
- Cereda CW, Christensen S, Campbell BC, et al (2016) A benchmarking tool to evaluate computer tomography perfusion infarct core predictions against a dwi standard. *Journal of Cerebral Blood Flow & Metabolism* 36(10):1780–1789
- Chalcroft L, Pereira RL, Brudfors M, et al (2023) Large-kernel attention for efficient and robust brain lesion segmentation. *arXiv preprint arXiv:230807251*
- Chalela JA, Kidwell CS, Nentwich LM, et al (2007) Magnetic resonance imaging and computed tomography in emergency assessment of patients with suspected acute stroke: a prospective comparison. *The Lancet* 369(9558):293–298
- Chao P, Kao CY, Ruan YS, et al (2019) Hardnet: A low memory traffic network. In: *Proceedings of the IEEE/CVF international conference on computer vision*, pp 3552–3561
- Chen J, Lu Y, Yu Q, et al (2021) Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:210204306*
- Chen LC, Papandreou G, Kokkinos I, et al (2017) 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
- Chen LC, Zhu Y, Papandreou G, et al (2018) Encoder-decoder with atrous separable convolution for semantic image segmentation. In: *Proceedings of the European*

- conference on computer vision (ECCV), pp 801–818
- Clerigues A, Valverde S, Bernal J, et al (2019) Acute ischemic stroke lesion core segmentation in ct perfusion images using fully convolutional neural networks. *Computers in biology and medicine* 115:103487
- Clèrigues A, Valverde S, Bernal J, et al (2020) Acute and sub-acute stroke lesion segmentation from multimodal mri. *Computer methods and programs in biomedicine* 194:105521
- Dimyan MA, Cohen LG (2011) Neuroplasticity in the context of motor rehabilitation after stroke. *Nature Reviews Neurology* 7(2):76–85
- Dosovitskiy A, Beyer L, Kolesnikov A, et al (2020) An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*
- d’Ascoli S, Touvron H, Leavitt ML, et al (2021) Convit: Improving vision transformers with soft convolutional inductive biases. In: *International Conference on Machine Learning*, PMLR, pp 2286–2296
- Feigin VL, Brainin M, Norrving B, et al (2022) World stroke organization (wso): global stroke fact sheet 2022. *International Journal of Stroke* 17(1):18–29
- Feng P, Ni B, Cai X, et al (2022) Utransnet: Transformer within u-net for stroke lesion segmentation. In: *2022 IEEE 25th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, IEEE, pp 359–364
- Fiebach J, Schellinger P, Jansen O, et al (2002) Ct and diffusion-weighted mr imaging in randomized order: diffusion-weighted imaging results in higher accuracy and lower interrater variability in the diagnosis of hyperacute ischemic stroke. *Stroke* 33(9):2206–2210
- Flossmann E, Redgrave JN, Briley D, et al (2008) Reliability of clinical diagnosis of the symptomatic vascular territory in patients with recent transient ischemic attack or minor stroke. *Stroke* 39(9):2457–2460
- Goldberger AL, Amaral LA, Glass L, et al (2000) Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation* 101(23):e215–e220
- Goldstein LB, Simel DL (2005) Is this patient having a stroke? *Jama* 293(19):2391–2402
- Gómez S, Mantilla D, Garzón G, et al (2023) Apis: A paired ct-mri dataset for ischemic stroke segmentation challenge. *arXiv preprint arXiv:2309.15243*
- Grysiewicz RA, Thomas K, Pandey DK (2008) Epidemiology of ischemic and hemorrhagic stroke: incidence, prevalence, mortality, and risk factors. *Neurologic clinics*

26(4):871–895

- Gu Y, Piao Z, Yoo SJ (2022) Sthardnet: Swin transformer with hardnet for mri segmentation. *Applied Sciences* 12(1):468
- Guo MH, Xu TX, Liu JJ, et al (2022) Attention mechanisms in computer vision: A survey. *Computational visual media* 8(3):331–368
- Guo MH, Lu CZ, Liu ZN, et al (2023) Visual attention network. *Computational Visual Media* 9(4):733–752
- Hakim A, Christensen S, Winzeck S, et al (2021) Predicting infarct core from computed tomography perfusion in acute ischemia with machine learning: Lessons from the isles challenge. *Stroke* 52(7):2328–2337
- Hatamizadeh A, Nath V, Tang Y, et al (2021) Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In: *International MICCAI Brainlesion Workshop*, Springer, pp 272–284
- Hatamizadeh A, Tang Y, Nath V, et al (2022) Unetr: Transformers for 3d medical image segmentation. In: *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp 574–584
- Hernandez Petzsche MR, de la Rosa E, Hanning U, et al (2022) Isles 2022: A multi-center magnetic resonance imaging stroke lesion segmentation dataset. *Scientific data* 9(1):762
- Hssayeni MD, Croock MS, Salman AD, et al (2020) Intracranial hemorrhage segmentation using a deep convolutional model. *Data* 5(1):14
- Hu J, Shen L, Sun G (2018) Squeeze-and-excitation networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 7132–7141
- Hu X, Luo W, Hu J, et al (2020) Brain segnet: 3d local refinement network for brain lesion segmentation. *BMC medical imaging* 20:1–10
- Huang G, Liu Z, Van Der Maaten L, et al (2017) Densely connected convolutional networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 4700–4708
- Hui H, Zhang X, Li F, et al (2020) A partitioning-stacking prediction fusion network based on an improved attention u-net for stroke lesion segmentation. *IEEE Access* 8:47419–47432
- Huo J, Chen L, Liu Y, et al (2022) Mapping: Model average with post-processing for stroke lesion segmentation. *arXiv preprint arXiv:221115486*

- Hwang DY, Silva GS, Furie KL, et al (2012) Comparative sensitivity of computed tomography vs. magnetic resonance imaging for detecting acute posterior fossa infarct. *The Journal of emergency medicine* 42(5):559–565
- Isensee F, Jaeger PF, Kohl SA, et al (2021) nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* 18(2):203–211
- Islam M, Ren H (2018) Class balanced pixnet for neurological image segmentation. In: *Proceedings of the 2018 6th International Conference on Bioinformatics and Computational Biology*, pp 83–87
- Jia X, Bartlett J, Zhang T, et al (2022) U-net vs transformer: Is u-net outdated in medical image registration? In: *International Workshop on Machine Learning in Medical Imaging*, Springer, pp 151–160
- Kadry S, Damaševičius R, Taniar D, et al (2021) U-net supported segmentation of ischemic-stroke-lesion from brain mri slices. In: *2021 Seventh International conference on Bio Signals, Images, and Instrumentation (ICBSII)*, IEEE, pp 1–5
- Kamnitsas K, Ferrante E, Parisot S, et al (2016) Deepmedic for brain tumor segmentation. In: *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: Second International Workshop, BrainLes 2016, with the Challenges on BRATS, ISLES and mTOP 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 17, 2016, Revised Selected Papers 2*, Springer, pp 138–149
- Karimi D, Vasylechko SD, Gholipour A (2021) Convolution-free medical image segmentation using transformers. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I* 24, Springer, pp 78–88
- Khezipour S, Seyedarabi H, Razavi SN, et al (2022) Automatic segmentation of the brain stroke lesions from mr flair scans using improved u-net framework. *Biomedical Signal Processing and Control* 78:103978
- Kirillov A, Mintun E, Ravi N, et al (2023) Segment anything. *arXiv preprint arXiv:230402643*
- Kofler F, Möller H, Buchner JA, et al (2023) Panoptica–instance-wise evaluation of 3d semantic and instance segmentation maps. *arXiv preprint arXiv:231202608*
- Kumar A, Upadhyay N, Ghosal P, et al (2020) Csnnet: A new deepnet framework for ischemic stroke lesion segmentation. *Computer Methods and Programs in Biomedicine* 193:105524
- Li J, Chen J, Tang Y, et al (2023a) Transforming medical imaging with transformers? a comparative review of key properties, current progresses, and future perspectives. *Medical image analysis* p 102762

- Li X, Chen H, Qi X, et al (2018) H-denseunet: hybrid densely connected unet for liver and tumor segmentation from ct volumes. *IEEE transactions on medical imaging* 37(12):2663–2674
- Li X, Huang K, Yang W, et al (2019) On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:190702189*
- Li X, Jiang M, Zhang X, et al (2021) Fedbn: Federated learning on non-iid features via local batch normalization. *arXiv preprint arXiv:210207623*
- Li X, Luo G, Wang K, et al (2023b) The state-of-the-art 3d anisotropic intracranial hemorrhage segmentation on non-contrast head ct: The instance challenge. *arXiv preprint arXiv:230103281*
- Li Y, Cai W, Gao Y, et al (2022) More than encoder: Introducing transformer decoder to upsample. In: *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, IEEE, pp 1597–1602
- Liang K, Han K, Li X, et al (2021) Symmetry-enhanced attention network for acute ischemic infarct segmentation with non-contrast ct images. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VII 24*, Springer, pp 432–441
- Liew SL, Lo BP, Donnelly MR, et al (2022) A large, curated, open-source stroke neuroimaging dataset to improve lesion segmentation algorithms. *Scientific data* 9(1):320
- Lin TY, Goyal P, Girshick R, et al (2017) Focal loss for dense object detection. In: *Proceedings of the IEEE international conference on computer vision*, pp 2980–2988
- Liu CF, Leigh R, Johnson B, et al (2023a) A large public dataset of annotated clinical mris and metadata of patients with acute stroke. *Scientific Data* 10(1):548
- Liu H, Wang Y, Fan W, et al (2022a) Trustworthy ai: A computational perspective. *ACM Transactions on Intelligent Systems and Technology* 14(1):1–59
- Liu L, Wu FX, Wang J (2019a) Efficient multi-kernel dcnn with pixel dropout for stroke mri segmentation. *Neurocomputing* 350:117–127
- Liu L, Kurgan L, Wu FX, et al (2020) Attention convolutional neural network for accurate segmentation and quantification of lesions in ischemic stroke disease. *Medical Image Analysis* 65:101791
- Liu L, Wang Y, Chang J, et al (2022b) Llrhnet: Multiple lesions segmentation using local-long range features. *Frontiers in Neuroinformatics* 16:859973

- Liu R, Pu W, Zou Y, et al (2022c) Pool-unet: Ischemic stroke segmentation from ct perfusion scans using poolformer unet. In: 2022 6th Asian Conference on Artificial Intelligence Technology (ACAIT), IEEE, pp 1–6
- Liu X, Yang H, Qi K, et al (2019b) Msdf-net: Multi-scale deep fusion network for stroke lesion segmentation. *IEEE Access* 7:178486–178495
- Liu Z, Lin Y, Cao Y, et al (2021) Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 10012–10022
- Liu Z, Tong L, Chen L, et al (2023b) Deep learning based brain tumor segmentation: a survey. *Complex & intelligent systems* 9(1):1001–1026
- Lucas C, Kemmling A, Mamlouk AM, et al (2018) Multi-scale neural network for automatic segmentation of ischemic strokes on acute perfusion images. In: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018), IEEE, pp 1118–1121
- Luo C, Zhang J, Chen X, et al (2021) Ucatr: Based on cnn and transformer encoding and cross-attention decoding for lesion segmentation of acute ischemic stroke in non-contrast computed tomography images. In: 2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC), IEEE, pp 3565–3568
- Maier O, Menze BH, Von der Gablentz J, et al (2017) Isles 2015-a public evaluation benchmark for ischemic stroke lesion segmentation from multispectral mri. *Medical image analysis* 35:250–269
- Marcus A, Bentley P, Rueckert D (2023) Concurrent ischemic lesion age estimation and segmentation of ct brain using a transformer-based network. *IEEE Transactions on Medical Imaging*
- Meyer MJ, Pereira S, McClure A, et al (2015) A systematic review of studies reporting multivariable models to predict functional outcomes after post-stroke inpatient rehabilitation. *Disability and rehabilitation* 37(15):1316–1323
- Milletari F, Navab N, Ahmadi SA (2016) V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: 2016 fourth international conference on 3D vision (3DV), Ieee, pp 565–571
- Mondal AK, Bhattacharjee A, Singla P, et al (2021) xvitcos: explainable vision transformer based covid-19 screening using radiography. *IEEE Journal of Translational Engineering in Health and Medicine* 10:1–10
- Ni H, Xue Y, Wong K, et al (2022) Asymmetry disentanglement network for interpretable acute ischemic stroke infarct segmentation in non-contrast ct scans. In:

- International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp 416–426
- O’Shea K, Nash R (2015) An introduction to convolutional neural networks. arXiv preprint arXiv:151108458
- Pereira S, Pinto A, Amorim J, et al (2019) Adaptive feature recombination and recalibration for semantic segmentation with fully convolutional networks. *IEEE transactions on medical imaging* 38(12):2914–2925
- Praveen G, Agrawal A, Sundaram P, et al (2018) Ischemic stroke lesion segmentation using stacked sparse autoencoder. *Computers in biology and medicine* 99:38–52
- Qi K, Yang H, Li C, et al (2019) X-net: Brain stroke lesion segmentation based on depthwise separable convolution and long-range dependencies. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III* 22, Springer, pp 247–255
- Ranjbarzadeh R, Caputo A, Tirkolaee EB, et al (2023) Brain tumor segmentation of mri images: A comprehensive review on the application of artificial intelligence tools. *Computers in biology and medicine* 152:106405
- Rao Y, Zhao W, Liu B, et al (2021) Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems* 34:13937–13949
- Ronneberger O, Fischer P, Brox T (2015) U-net: Convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III* 18, Springer, pp 234–241
- Rubin J, Abulnaga SM (2019) Ct-to-mr conditional generative adversarial networks for ischemic stroke lesion segmentation. In: *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, IEEE, pp 1–7
- Sheller MJ, Edwards B, Reina GA, et al (2020) Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data. *Scientific reports* 10(1):12598
- Shen Z, Zhang M, Zhao H, et al (2021) Efficient attention: Attention with linear complexities. In: *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp 3531–3539
- Simonsen CZ, Madsen MH, Schmitz ML, et al (2015) Sensitivity of diffusion-and perfusion-weighted imaging for diagnosing acute ischemic stroke is 97.5%. *Stroke* 46(1):98–101

- Singh A, Sengupta S, Lakshminarayanan V (2020) Explainable deep learning models in medical image analysis. *Journal of imaging* 6(6):52
- Soh WK, Yuen HY, Rajapakse JC (2023) Hut: Hybrid unet transformer for brain lesion and tumour segmentation. *Heliyon*
- Sudre CH, Li W, Vercauteren T, et al (2017) Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In: *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*, Springer, pp 240–248
- Tomita N, Jiang S, Maeder ME, et al (2020) Automatic post-stroke lesion segmentation on mr images using 3d residual convolutional neural network. *NeuroImage: clinical* 27:102276
- Tragakis A, Kaul C, Murray-Smith R, et al (2023) The fully convolutional transformer for medical image segmentation. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp 3660–3669
- Vaswani A, Shazeer N, Parmar N, et al (2017) Attention is all you need. *Advances in neural information processing systems* 30
- de Vries L, Emmer B, Majoie C, et al (2021) Transformers for ischemic stroke infarct core segmentation from spatio-temporal ct perfusion scans. In: *Medical Imaging with Deep Learning*
- de Vries L, Emmer BJ, Majoie CB, et al (2023) Perfu-net: Baseline infarct estimation from ct perfusion source data for acute ischemic stroke. *Medical Image Analysis* 85:102749
- Vupputuri A, Dighade S, Prasanth P, et al (2018) Symmetry determined superpixels for efficient lesion segmentation of ischemic stroke from mri. In: *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE, pp 742–745
- Wang D, Wu Z, Yu H (2021) Ted-net: Convolution-free t2t vision transformer-based encoder-decoder dilation network for low-dose ct denoising. In: *Machine Learning in Medical Imaging: 12th International Workshop, MLMI 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Proceedings 12*, Springer, pp 416–425
- Wang H, Cao P, Wang J, et al (2022a) Uctransnet: rethinking the skip connections in u-net from a channel-wise perspective with transformer. In: *Proceedings of the AAAI conference on artificial intelligence*, pp 2441–2449



- Wang J, Wang S, Liang W (2022b) Metrans: Multi-encoder transformer for ischemic stroke segmentation. *Electronics Letters* 58(9):340–342
- Wang S, Li BZ, Khabisa M, et al (2020) Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:200604768*
- Wang W, Xie E, Li X, et al (2022c) Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media* 8(3):415–424
- Wang X, Shrivastava A, Gupta A (2017) A-fast-rcnn: Hard positive generation via adversary for object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 2606–2615
- Wang Y, Katsaggelos AK, Wang X, et al (2016) A deep symmetry convnet for stroke lesion segmentation. In: *2016 IEEE International Conference on Image Processing (ICIP)*, IEEE, pp 111–115
- Wang Y, Chen K, Yuan W, et al (2023) Samihs: Adaptation of segment anything model for intracranial hemorrhage segmentation. *arXiv preprint arXiv:231108190*
- Wessels T, Wessels C, Ellsiepen A, et al (2006) Contribution of diffusion-weighted imaging in determination of stroke etiology. *American Journal of Neuroradiology* 27(1):35–39
- Winzeck S, Hakim A, McKinley R, et al (2018) Isles 2016 and 2017-benchmarking ischemic stroke lesion outcome prediction based on multispectral mri. *Frontiers in neurology* 9:679
- Woo S, Park J, Lee JY, et al (2018) Cbam: Convolutional block attention module. In: *Proceedings of the European conference on computer vision (ECCV)*, pp 3–19
- Wu Z, Zhang X, Li F, et al (2022) Multi-scale long-range interactive and regional attention network for stroke lesion segmentation. *Computers and Electrical Engineering* 103:108345
- Wu Z, Zhang X, Li F, et al (2023a) W-net: A boundary-enhanced segmentation network for stroke lesions. *Expert Systems with Applications* p 120637
- Wu Z, Zhang X, Li F, et al (2023b) Transrender: a transformer-based boundary rendering segmentation network for stroke lesions. *Frontiers in Neuroscience* 17
- Xie S, Girshick R, Dollár P, et al (2017) Aggregated residual transformations for deep neural networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 1492–1500
- Xiong Y, Zeng Z, Chakraborty R, et al (2021) Nyströmformer: A nyström-based algorithm for approximating self-attention. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp 14138–14148

- Xu Z, Ding C (2023) Combining convolutional attention mechanism and residual deformable transformer for infarct segmentation from ct scans of acute ischemic stroke patients. *Frontiers in Neurology* 14
- Yang H, Huang W, Qi K, et al (2019) Clci-net: Cross-level fusion and context inference networks for lesion segmentation of chronic stroke. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III* 22, Springer, pp 266–274
- Yu W, Luo M, Zhou P, et al (2022) Metaformer is actually what you need for vision. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp 10819–10829
- Yu W, Huang Z, Zhang J, et al (2023a) San-net: Learning generalization to unseen sites for stroke lesion segmentation with self-adaptive normalization. *Computers in Biology and Medicine* 156:106717
- Yu W, Lei Y, Shan H (2023b) Fan-net: Fourier-based adaptive normalization for cross-domain stroke lesion segmentation. In: *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, pp 1–5
- Zhang H, Chen H (2023) Efficient 3d transformer with cluster-based domain-adversarial learning for 3d medical image segmentation. In: *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, IEEE, pp 1–5
- Zhang L, Song R, Wang Y, et al (2020) Ischemic stroke lesion segmentation using multi-plane information fusion. *IEEE Access* 8:45715–45725
- Zhang R, Zhao L, Lou W, et al (2018a) Automatic segmentation of acute ischemic stroke from dwi using 3-d fully convolutional densenets. *IEEE transactions on medical imaging* 37(9):2149–2160
- Zhang Y, Liu H, Hu Q (2021) Transfuse: Fusing transformers and cnns for medical image segmentation. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I* 24, Springer, pp 14–24
- Zhang Y, Liu S, Li C, et al (2022) Application of deep learning method on ischemic stroke lesion segmentation. *Journal of Shanghai Jiaotong University (Science)* pp 1–13
- Zhang Z, Liu Q, Wang Y (2018b) Road extraction by deep residual u-net. *IEEE Geoscience and Remote Sensing Letters* 15(5):749–753
- Zhao H, Shi J, Qi X, et al (2017) Pyramid scene parsing network. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 2881–2890

- Zhou HY, Guo J, Zhang Y, et al (2021) nnformer: Interleaved transformer for volumetric segmentation. arXiv preprint arXiv:210903201
- Zhou SK, Rueckert D, Fichtinger G (2019a) Handbook of medical image computing and computer assisted intervention. Academic Press
- Zhou Y, Huang W, Dong P, et al (2019b) D-unet: a dimension-fusion u shape network for chronic stroke lesion segmentation. IEEE/ACM transactions on computational biology and bioinformatics 18(3):940–950
- Zhu X, Su W, Lu L, et al (2020) Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:201004159