# Recent Progress in Transformer-based Medical Image Analysis

Zhaoshan Liu[a], Qiujie Lv[a,b], Ziduo Yang[a,b], Yifan Li[a], Chau Hung Lee[c], Lei Shen[a,*]

[a]*Department of Mechanical Engineering, National University of Singapore, 9 Engineering Drive 1, Singapore, 117575, Singapore*
[b]*School of Intelligent Systems Engineering, Sun Yat-sen University, No.66, Gongchang Road, Guangming District, 518107, China*
[c]*Department of Radiology, Tan Tock Seng Hospital, 11 Jalan Tan Tock Seng, Singapore, 308433, Singapore*

## Abstract

The transformer is primarily used in the field of natural language processing. Recently, it has been adopted and shows promise in the computer vision (CV) field. Medical image analysis (MIA), as a critical branch of CV, also greatly benefits from this state-of-the-art technique. In this review, we first recap the core component of the transformer, the attention mechanism, and the detailed structures of the transformer. After that, we depict the recent progress of the transformer in the field of MIA. We organize the applications in a sequence of different tasks, including classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis. The mainstream classification and segmentation tasks are further divided into eleven medical image modalities. A large number of experiments studied in this review illustrate that the transformer-based method outperforms existing methods through comparisons with multiple evaluation metrics. Finally, we discuss the open challenges and future opportunities in this field. This task-modality review with the latest contents, detailed information, and comprehensive comparison may greatly benefit the broad MIA community.

*Keywords:* Deep Learning, Transformer, Attention Mechanism, Convolutional Neural Network, Medical Image Analysis

## 1. Introduction

Transformer [1] is one of the most widely used models in the natural language processing (NLP) field and has achieved great success in many tasks, such as paraphrase generation [2], text-to-speech synthesis [3], and speech recognition [4]. It is

---

*Corresponding author

*Email addresses:* e0575844@u.nus.edu (Zhaoshan Liu), lvqj5@mail2.sysu.edu.cn (Qiujie Lv), yangzd@mail2.sysu.edu.cn (Ziduo Yang), e0576095@u.nus.edu (Yifan Li), chau_hung_lee@ttsh.com.sg (Chau Hung Lee), mpeshel@nus.edu.sg (Lei Shen)

arXiv:2208.06643v4 [eess.IV] 25 Jul 2023

designed for transduction and sequence modeling with the remarkable capability of modeling long-range dependencies with the data. The convolutional-free transformer is based on the self-attention (attention) mechanism, a successful NLP technique [5–9] that relates different positions of a single sequence to compute the sequence's representation [1]. Unlike the NLP field, the computer vision (CV) field has been dominated by the convolutional neural network (CNN) [10] for a long time [11–15]. Even if, many trials have been carried out to combine CNN and attention in the CV field [16, 17] while none of them overperform CNN. Until 2020, Dosovitskiy et al. [18] proposed a pioneering model and prove that implementing the transformer directly to sequences of image patches works well for image classification. In detail, the proposed method split the input image into multiple patches and embeds each of them linearly. With additional position embeddings added, the resulting vector sequences are fed to the transformer encoder. With the solid foundation they set, the transformer-based method has been widely adopted in the field of CV with superior performance [19–22].

Medical image analysis (MIA) is an essential branch in the CV field. Medical imaging utilizes various modalities to create a visual representation of the inside body [23] and is of great help for further medical diagnosing. There are several kinds of medical imaging modalities, such as magnetic resonance imaging (MRI), computed tomography (CT), ultrasound (US), positron emission tomography (PET), optical coherence tomography (OCT), and digital fundus imaging (DFI). In practice, MIA is usually performed qualitatively by medical personnel. This may result in varying interpretations and degrees of accuracy because of varying degrees of reader experience or varying image quality. Moreover, such image analysis may be time- and labor-expensive. Due to these, the deep learning (DL) method has been widely applied in the field of MIA to reduce inter-reader variation as well as reduce time and manpower costs [24–27]. With the rapid development of the transformer in CV, the transformer-based method has been widely used in MIA either using the transformer solely [28–30] or hybridizing CNN and transformer to capture both local and global information [31–33]. To help researchers catch up with this emerging research field, it is timely and important to have a comprehensive review and perspectives on transformer-based MIA.

In this review, we systematically introduce the transformer technique and its recent progress in the field of MIA, followed by outlooks and perspectives. We first recap the core component of the transformer, the attention mechanism, and the transformer itself. Then, we summarize the transformer-based applications in the sequence of different MIA tasks, including classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis, as shown in Figure 1. For the mainstream classification and segmentation tasks, we further divided their corresponding works into different medical imaging modalities. There are a total of eleven modalities in our review, including MRI, CT, X-ray, microscope, endoscopy, US, dermoscopy, DFI, camera, PET, and OCT, as shown in Figure 2. Finally, the open challenges and future research opportunities of transformer-based MIA tasks
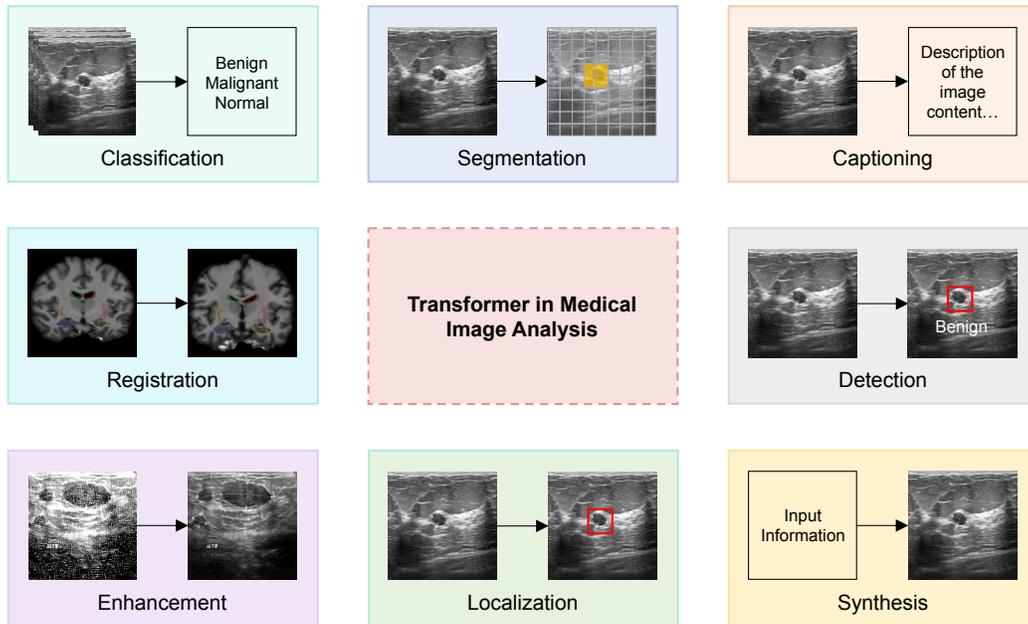
Figure 1: MIA tasks included in this review. The tasks are organized in a sequence of classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis [34, 35].
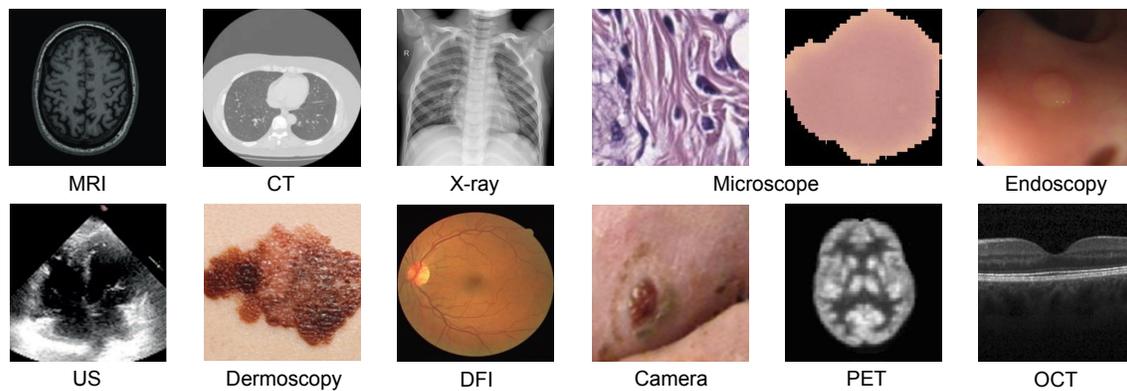


Figure 2: Examples of modalities included in this work. Sequences are MRI [36], CT [37], X-ray [29], microscope [38] (left), [39] (right), endoscopy [40], US [41], dermoscopy [42], DFI [43], camera [44], PET [45], and OCT [46].

are also discussed.

As this review is being written, we become aware of two similar review works [47, 48]. The first review overviews the attention mechanism and other components essential for building the transformer. Then, the transformer architectures designed for MIA applications are discussed. The transformer-based applications are discussed task-wise, in which classification and segmentation applications are further divided into pure and hybrid transformer-based. Extensive discussions are then conducted, including learning manners, model improvements, and performance comparison with CNN. The second review starts by illustrating various medical imaging modalities, followed by an introduction to DL concepts, techniques, and architectures. The different tasks and the transformer-based applications are further discussed. Finally, research trends, current challenges, and future prospects are discussed. Here, we would like to highlight a couple of main differences between our review and others published. First, we cover more than 100 of the latest relevant papers, providing readers with a view of the latest research progress. We also highlight the latest transformer models leveraged in MIA, such as the swin transformer [49], O-Net [50], and transformer-based region-edge aggregation network [51]. The swin transformer is introduced in great detail as it is widely accepted and well-performed across different tasks and modalities. Next, our review contains more details. We include works with more than one modality and summarize them in a one-to-multi manner. We also summarize the objects researched, the datasets used, and the disease corresponding to these datasets when applicable. The summary of contents is coherently performed throughout different tasks. Sufficient details can help new researchers in the field grasp the necessary concepts easier and faster. Finally, instead of giving a quantitative performance evaluation of the transformer-based method, we also provide a comprehensive performance comparison between the transformer-based models and existing state-of-the-art DL methods in MIA. This proves the effectiveness of the transformer-based method. In summary, our task-modality review presents updated contents, detailed information, and comprehensive comparison that will greatly benefit the MIA community.

The rest of this review is organized as follows: In Section 2, we show the methodology performed for our systematic review. In Section 3, we recap the principle of the attention mechanism, the detailed structures of the transformer, and depict how the transformer is adopted into the MIA field. An introduction to different training manners and MIA tasks is also included. Section 4 organizes the transformer-based MIA applications from the perspective of different tasks. To better organize the large number of works related to mainstream classification and segmentation tasks, we further categorize them based on the imaging modalities. The objects in the relevant references are tabulated in detail. The datasets used as well as the diseases corresponding to the datasets are also tabulated when applicable. Moreover, a quantitative performance comparison across the transformer-based method and existing methods are summarized separately. In Section 5, we point out the current challenges and future opportunities of the transformer-based MIA. A concise and com-

prehensive conclusion can be found in Section 6.

## 2. Methodology

With the fast development of the transformer in the CV field, the research on the transformer has become one of the most popular research directions in the MIA field. We search on the Scopus database using "transformer" on the "title" field and "vision" and "medical" on the "title-abs-key" field and the results show that the number of papers in 2022 is more than four times compared to that of in 2021. What is more, the number of publications in 2023 is already more than that of 2021. We thus want to explore several research questions for the transformer-based MIA. First, on which MIA tasks have the transformer-based been successfully applied? Second, does the transformer outperforms previous DL methods, such as the CNN-based method? Finally, what are the current challenges or problems to be solved and the corresponding potential feasible solutions?

In this systematic review, the relevant references are collected by searching within the Web of Science and Scopus databases. The search timeframe is from 2021 to 8 Feb 2023. In the Web of Science database, we include keywords including "transformer", "classification", "segmentation", "captioning", "registration", "detection", "enhancement", "reconstruction", "denoising", "localization", "synthesis", "generation", and "diagnosis" in the "title" field. We also include the keywords "medical" and "vision" in the "topic" field. Regarding the Scopus database, we include keywords in the "TITLE" and "TITLE-ABS-KEY" fields identical to that of the "title" and "topic" fields, respectively. With papers searched, we then exclude the results that are duplicated. Then, we remove the records, in which full text is inaccessible. Conference posters, as well as review, survey, and benchmark papers, are also excluded. Note that the conference papers are not excluded. The papers submitted on arXiv are also reserved. Finally, we check the content of the papers to remove papers that are not vision related. We find that several papers are included even if we include the "vision" keyword. We also exclude the papers using non-medical datasets. Following PRISMA [52], we show the flow diagram for our systematic review in Figure 3.

## 3. Background

### 3.1. Attention Mechanism

The attention mechanism is the core component of the transformer. It differentially weights the significance of each part of the input data and allows the inputs to interact with each other to find to whom they should pay more attention. The attention mechanism expresses the importance of each input (e.g., token) in the current context as the attention score. The outputs are the aggregation of these interactions weighted by the corresponding attention scores. Specifically, with three attention vectors named query, keys, and values, the mechanism maps a query and a set of
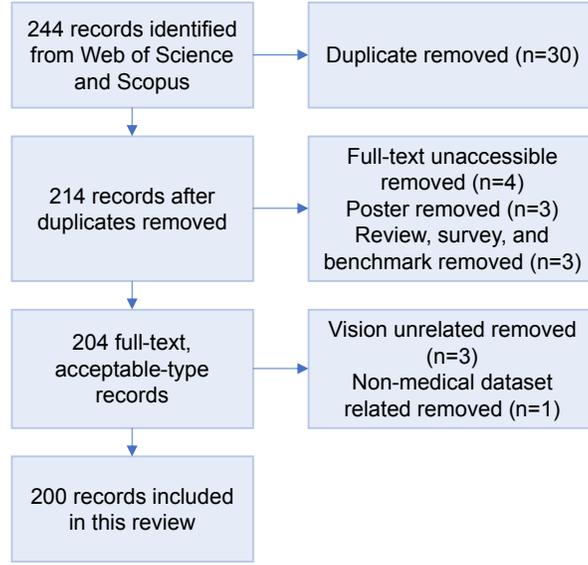
Figure 3: Flow diagram for our systematic review.

key-value pairs to output. The output is computed by summarizing all values according to the weight, which is calculated using a compatibility function of the query with the related key [1]. The attention mechanism can be divided into scaled dot-product attention and multi-head attention.

**Scaled dot-product attention**. Several queries, keys, and values compose the inputs of the scaled dot-product attention. The queries and keys have a dimensionality of $d_k$ and the values have a dimensionality of $d_v$. The dot product of all keys and the query is calculated. The resulting values are divided by a scale factor, $\sqrt{d_k}$, and pass a softmax function. The attention can be calculated through the dot product with the values. Practically, a set of queries, keys, and values are packed into corresponding matrices, *Q*, *K*, *V*, and the outputs matrix can be calculated using the below formula:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

**Multi-head attention**. The scaled dot-product attention is single-head attention. In practice, multi-head attention is more often used as it improves the expressive power of the model and stabilizes training. The multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. Specifically, the queries, keys, and values are projected for *h* times, where *h* is the number of heads. The attention function is performed on each of the projected results concurrently. The output values are concatenated and projected again to obtain the final results [1]. A concise illustration of the multi-head attention is shown in
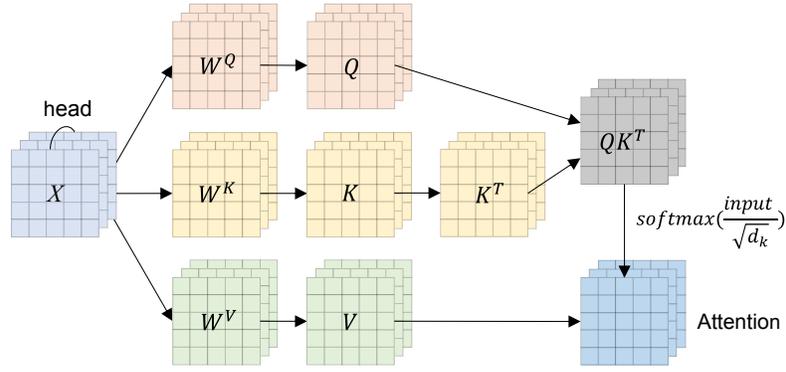
6

Figure 4: An intuitive illustration of the multi-head attention mechanism. The result of each head is calculated respectively, and the values are then concatenated. *X* denotes the input, $W^Q$, $W^K$, and $W^V$ represent the parameter matrices for projections, and $\sqrt{d_k}$ stands the scale factor.

Figure 4, and the multi-head attention can be described using the below equation:

$$\text{MultiHead}(Q, K, V) = \text{Concat}\left(\text{head}_1, \ldots, \text{head}_h\right) W^O$$

$$\text{where head}_i = \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right)$$

where $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$, $W^O \in \mathbb{R}^{h d_v \times d_{\text{model}}}$ are the parameter matrices for projections, and $d_{\text{model}}$ is dimensional keys, values, and queries.

### 3.2. Transformer

Researchers have proposed several variant models based on the attention mechanism, which generally combines the attention mechanism with the recurrent neural network (RNN), such as LSTM. These models are usually limited in training speed due to the sequential structure, and the parallel computing ability is limited. Since the attention model itself can capture global information, a natural question raised is whether we can remove the RNN structure and rely only on the attention model. The answer is yes. The transformer is such a novel model utilized to address sequence-to-sequence tasks, taking a sequence as the input and generating the predicted probabilities as the output. It is mainly composed of the attention mechanism and has an encoder-decoder structure, as shown in Figure 5. Both the encoder and decoder are tandem by several identical blocks. The blocks are composed of several parts, including the masked multi-head attention module, multi-head attention module, layer normalization, and position-wise feed-forward network. The masked multi-head attention module employs the attention mechanism up to the current position and excluded the unpredicted positions till now. Combining the masked multi-head attention and the position offset of the output embeddings, the predictions for the position only lie on the outputs at earlier positions can be ensured [1]. The fully
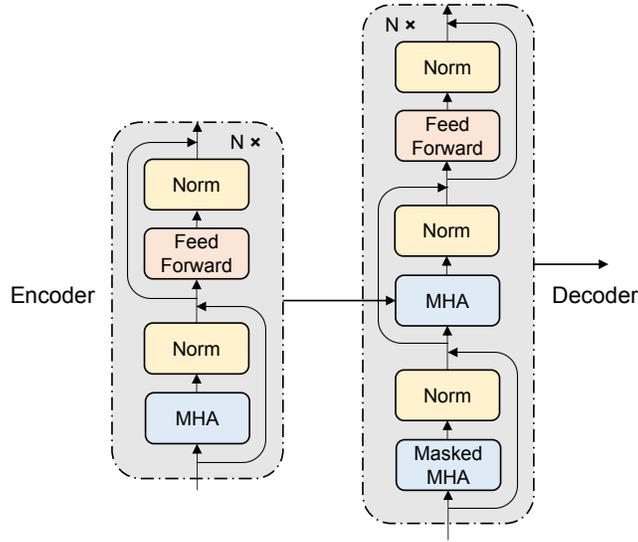
Figure 5: The encoder-decoder structure in transformer [1]. MHA refers to the multi-head attention module and Norm represents layer normalization. The output of the encoder is fed into the MHA in the decoder. **Left**: encoder, and **Right**: decoder.

connected feed-forward network consists of two linear layers and a ReLU activation function in between. The procedure can be expressed using the below equation:

$$\text{FFN}(x) = \max\left(0, xW_1 + b_1\right)W_2 + b_2$$

**Encoder**. The transformer encoder consists of *N* identical blocks. For each block, there are two main parts and both parts employ the residual connection proposed by He et al. [53]. The bottom part is composed of a multi-head attention module and layer normalization. The top part consists of a fully connected feed-forward network and layer normalization.

**Decoder**. The transformer decoder is composed of *N* identical blocks like the encoder and consists of three main parts in each block. The bottom part is composed of a masked multi-head attention module and layer normalization. The middle part consists of a multi-head attention module followed by layer normalization. It is worth noting that the multi-head attention module in the decoder also takes the encoder's output as the input. The top part is composed of a feed-forward network and layer normalization. The residual connection is implemented for all three parts.

### 3.3. Transformer in MIA

To employ the transformer in the field of MIA, the input medical images need to be pre-processed [18] as they are 2D. The pre-processing process can be divided into two main steps, patch generation, and embedding. The patch generation splits the input image into several patches, and the embedding flattens the patches and generates patch embeddings. Positional embeddings and class embedding are
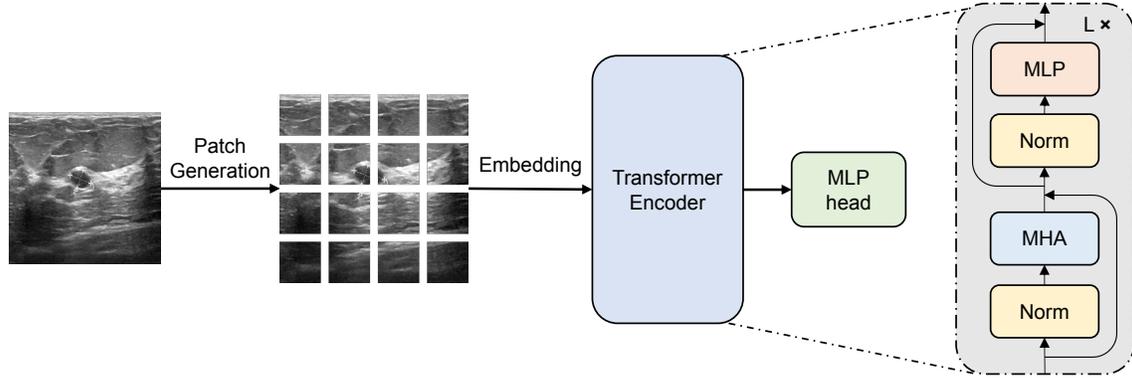
Figure 6: Detailed structure of the transformer used in the MIA field. MHA refers to the multi-head attention module. Norm represents layer normalization and MLP illustrates the multilayer perception module.

also added. The processed images are fed into the transformer encoder for feature extraction. The transformer encoder undergoes slight modification by altering the sequence of layer normalization. The output from the transformer encoder serves as the input for the multilayer perceptron head, resulting in the image class prediction. The detailed structure of the transformer used in MIA is illustrated in Figure 6.

**Patch generation**. The patch generation receives image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ as the input, where $H$ and $W$ are the height and the weight of the image, respectively. The input image is reshaped into a patch sequence $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$, where $P$ is the dimension of each image patch, and $N$ is the number of patches [18]. The relationship between $N$, $H$, $W$, and $P$ can be expressed as:

$$N = HW/P^2$$

**Embedding**. The embedding obtains patches as the input. It flattens and maps them to $D$ dimensions using a linear projection $\mathbf{E}$, where $D$ is the latent vector size of the transformer layers. The resulting patch embeddings are concatenated with the class embedding [54] and then the concatenated embeddings are summed with the position embeddings to retain positional information. The embedding can be described through the below equation [18]:

$$\mathbf{y} = \left[ \mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots ; \mathbf{x}_p^N \mathbf{E} \right] + \mathbf{E}_{\text{pos}} \quad \text{where } \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D}$$

**Transformer encoder**. The transformer encoder takes the resulting embeddings as the input. The overall encoder structure is similar to that in Figure 5 while the layer normalization is moved before the multi-head attention module and multilayer perceptron module (feed-forward network). The multilayer perceptron module is composed of two linear layers and both of the layers use GELU as the activation function.

**Multilayer perceptron head**. The output of the transformer encoder is fed into
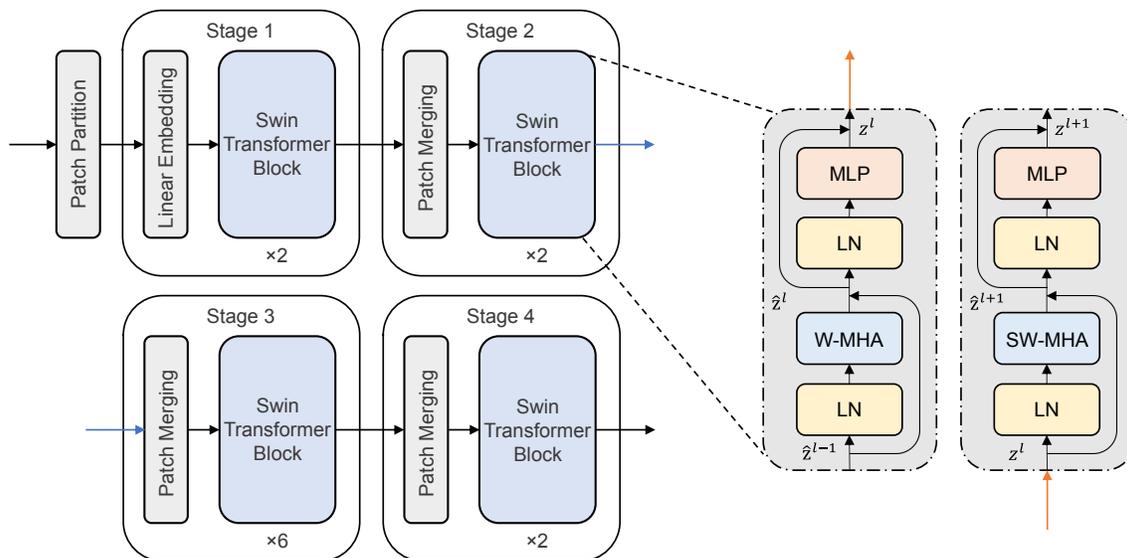
9

Figure 7: The detailed structure of the swin transformer. The swin transformer consists of four stages and two successive swin transformer blocks form a basic unit. MHA, W-MHA, and SW-MHA denote the multi-head attention module, multi-head attention module with regular windowing configurations, and multi-head attention module with shifted windowing configurations, respectively. MLP represents the multilayer perception module, and LN is layer normalization. $\hat{\mathbf{z}}^l$ and $\mathbf{z}^l$ are the output features of the W-MHA and the MLP inside block $l$, respectively. Blue and orange arrows mean the connection of two parts.

the multilayer perceptron head to get the classification result. At the pre-training stage, the multilayer perceptron head is composed of one hidden layer, while at the fine-tuning stage, it consists of a single linear layer.

The transformer faces several challenges in the CV field that stem from the differences between the CV and NLP fields. For one thing, the scale of visual entities can vary largely. For another, the resolution of images can be very high. To solve these problems, the swin transformer has been proposed. The swin transformer builds hierarchical feature maps and calculates the attention only within each local window. Compared with computing attention globally, the computation complexity is reduced largely from quadratic to linear. To provide connections between different windows, a shifted window approach is proposed. The window partitioning is shifted after the attention is calculated within each window. With such designs, the swin transformer can not only reduce the computation complexity to linear from quadratic but also model at various scales with high flexibility. The overall structure of the swin transformer is shown in Figure 7. There are four stages in the swin transformer and two successive swin transformer blocks form a basic unit. The former swin transformer block has a multi-head attention module with regular windowing configurations, while the latter consists of a shifted windowing multi-head attention module. The multilayer

perceptron module consists of two layers with a GELU activation function in between, and layer normalization is implemented in both of the blocks.

There are distinct advantages of using the transformer model over CNN in the field of MIA. In the MIA field, many involved images have repetitive patterns or symmetrical patterns, such as the microscope image (left) and the CT image shown in Figure 2. CNN is not sensitive to such symmetrical or repetitive patterns as it only focuses on local features. In contrast, the transformer can capture these global features by exploring the relationships among local regions. Take a tumor image consisting of repetitive normal patterns and a tumor pattern as an example, the calculated attention scores may show high similarities among repeatedly normal regions and low similarities between the tumor region and the normal regions. This can demonstrate that the transformer is sensitive to repetitive patterns.

The scarcity of datasets is a potential issue for the transformer. Training the transformer requires a large amount of data as it lacks some of CNN's inherent inductive biases, such as locality and translation equivariance. It is reported that training the transformer with less than 100 million images usually obtains a suboptimal solution compared to CNN, while the performance of the transformer continues increasing with the increase of the dataset size. However, the process of creating a large dataset of medical images is substantially different from nature images due to various reasons. For instance, it heavily relies on costly equipment to capture medical images, which subsequently requires human experts for annotation. Additionally, medical datasets cannot always be made publicly available due to patient privacy concerns. Therefore, collecting a dataset with more than 100 million images is a significant challenge in the MIA field, and sometimes there are only thousands or even hundreds of images in a dataset [55, 56]. In the case of common data shortage, data augmentation methods are widely implemented in MIA, such as traditional image transformations like flip [57], or newer image synthesis methods [58–60]. Besides, transfer learning is another widely used technique in the MIA field [58, 61, 62]. Transfer learning allows the model to be pre-trained on a larger dataset, such as ImageNet, and the learned knowledge from the bigger dataset can be utilized when training the smaller ones.

### 3.4. Training Manner

Similar to other DL models, the transformer is trained in different training manners. This can include supervised (full-supervised), unsupervised, semi-supervised, self-supervised, weakly supervised, et al. The supervised learning manner is the most widely implemented learning manner and data used for supervised learning contain full labels for model evaluation. Opposite to supervised learning, unsupervised learning only receives the input data and learns intrinsic data properties by discovering underlying structures, patterns, or relationships in the data. Weakly-supervised learning is a learning method between supervised learning and unsupervised learning, in which the learning algorithm is trained using incomplete, imprecise, or noisy labeled data. Semi-supervised learning is a combination of supervised

11

learning and unsupervised, which uses a large amount of unlabeled data and a small amount of labeled data for training to improve model performance. Self-supervised learning enables the model to learn richer and more effective feature representations by designing specific pretext tasks and using the input data to generate supervisory signals. The two most widely implemented self-supervised learning algorithms are DINO [63] and BYOL [64], where DINO utilizes the knowledge distillation technique and BYOL minimizes the difference in latent representations between augmented image pairs. The BYOL is composed of two networks, which are the online network and the target network. The target network is utilized to provide targets for online network training. The introduction of learning methods outside of supervised learning can reduce the requirement of the label amount to a large extent and has been widely implemented in the MIA field.

### 3.5. MIA Task

There are numerous tasks in the field of MIA. Here, we include classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis in our review. Classification is the process of categorizing images into distinct classes. Segmentation partitions images into various objects or subgroups and can be regarded as pixel-level classification. Captioning generates descriptive language using visual information. Registration involves transforming multiple sets of data obtained from different sensors, viewpoints, etc. [65] into a unified coordinate system. Object detection predicts the boundary and the classification result across different objects. It is worth noting that one type of object detection [66] also performs pixel-level classification, while it is not widely implemented in the MIA field partly due to computation resource considerations. Localization is a similar task to object detection while it predicts the object boundary solely. Enhancement works to enhance patterns and remove noise artifacts, which primarily includes reconstruction and denoising. Reconstruction enhances image quality by addressing potential low signal-to-noise ratio, contrast-to-noise ratio, and artifacts [67], while denoising enhances visual images through noise removal. Synthesis creates desired images, which is the opposite of classification. A concise graphical illustration for each task can be found in Figure 1.

## 4. Applications

We discuss the applications of the transformer in classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis tasks. Given the large number of papers involved in mainstream classification and segmentation tasks, we further categorize them by modalities. There are eleven modalities included in this work, which are MRI, CT, X-ray, microscope, endoscopy, US, dermoscopy, DFI, camera, PET, and OCT. It is worth noting that endoscopy includes colonoscopy, laryngoscopy, etc, and the whole slide image (WSI) is included in the microscope as it is also referred to as virtual microscope [68]. Besides, magnetic
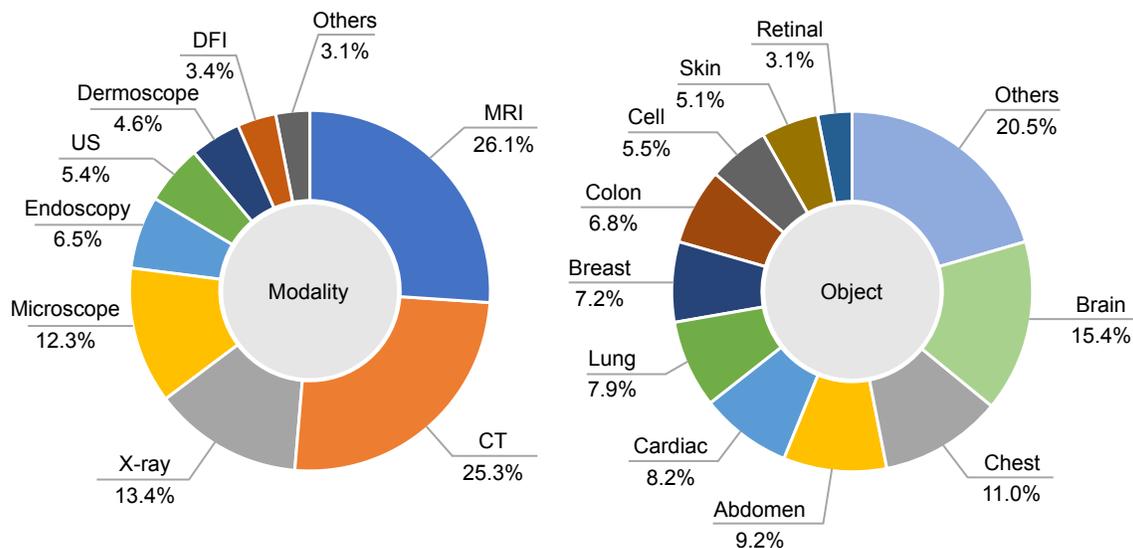
Figure 8: The proportion that each modality and object. **Left**: modality, and **Right**: object. There are around 30 objects in "Others", and each one has a small proportion.

resonance angiography is included in the MRI. The objects contained in the selected references are tabulated in detail. The proportion of each modality and object in all cited references are summarized in Figure 8. In order to facilitate statistical analysis, we present complex objects in a "grouped" manner, whereby several small objects are grouped into the same category. For instance, the colon covers the rectum. Additionally, the related datasets corresponding to the works are summarized in the tables, and the diseases with respect to the datasets are also tabulated when applicable. Diseases are also organized in a "grouped" way when appropriate. For example, lung disease can cover tuberculosis, COVID-19, pneumonia, pneumothorax, etc. The main narrative sequence within each section is that we usually start from the works using existing transformer models followed by the works with newly developed models. It is worth noting that we discuss the qualitative results like segmented masks within the applications for the sake of coherence, while the quantitative performance comparison is illustrated separately in Section 4.4. This can make the performance comparison between transformer-based models and existing models more intuitive.

### 4.1. Classification

Since classification is one of the most widely-studied applications in MIA, we organize this sub-section in the sequence of different modalities, including X-ray, microscope, CT, MRI, DFI, dermoscopy, endoscopy, US, PET, camera, and OCT. In the case of papers containing more than one modality, we put them into the "multiple" category following the OCT. Table 1 summarizes the classification applications with the transformer. In the classification application, most of the works combine

the transformer with CNN to capture both local and global information for further robustness and performance improvement.

**X-ray**. Many of the researchers use existing transformer models to classify medical X-ray images [28, 29, 69, 77, 79, 81, 82]. The implemented transformer models include the transformer, the swin transformer, and the DeiT [180]. Besides using existing transformer models, some researcher aims at combining transformer with other CNN models. Within these, several papers tandems the CNN and transformer. For instance, Duong and co-workers [32] constructed a ViT-Eff method, where input images are fed into EfficientNet [181] and the extracted feature maps are then projected into the transformer, followed by the developed classification head. Similar works include the PneuNet proposed by Wang et al. [87], in which ResNet and transformer are tandem, and the Chest L-Transformer [31] that tandems ResNeXt [182] and transformer. Jalalifar et al. [81] built their method on a DeiT structure and let it benefits from the teacher-student scheme. The DenseNet [183] is set as the teacher while the adapted transformer is chosen as the student. Leveraging existing transformer models or combining the transformer with existing CNN models is intuitive and effective. However, the investigation of the model architecture is lacking, and potential performance improvement may be realized with an in-depth structure design.

Instead of combining existing models directly, several researchers build their models from scratch. Jiang et al. [99] designed an MXT method consisting of five stages. The first four stages are composed of several downsample spatial reduction transformer blocks and a multi-layer overlap patch embedding block. The last stage is composed of two class token transformer blocks and a multi-label attention block. Qi and co-workers [102] proposed a multi-feature fusion transformer where the cross-attention mechanism is deployed to learn information from both original images and corresponding enhanced local phase images. Jiang and Chen [33] developed an MP-ViT model, where images are fed to the patch fuser after enhancement and layer normalization. The obtained fusion features are then trained together with smoothed labels to obtain final prediction results. Park et al. [107] developed a federated split transformer with the FESTA learning process. In FESTA, the server first initializes the weights of the transformer as well as task-specific heads and tails for each task. Then, it distributes the heads and tails weights to each client. For each round, each client (e.g., hospital) carries out the forward propagation on their head and conveys the intermediate feature to the server. Finally, the server aggregates and averages the weights of local heads and tails and distributes the updated global weights back to the clients.

**Microscope**. Some researchers utilize the existing transformer models to classify medical microscope images [110, 112]. For instance, Zeid and co-workers [112] proposed to use the transformer and the compact convolutional transformer [184], in which the convolutional tokenizer is implemented instead of the patch-based tokenizer. The convolutional tokenization is composed of a convolutional layer, a ReLU activation function, followed by max pooling and reshaping, and can benefit the

Table 1: Transformer-based classification applications. The mark "-" shows the corresponding information is not publicly available.

| Method | Year | Modality | Object | Dataset |
|---|---|---|---|---|
| transformer [69] | 2022 | X-ray | lung | lung disease [70, 71] |
| swin transformer [28] | 2022 | X-ray | chest | lung disease [72–75] |
| transformer [29] | 2021 | X-ray | chest | lung disease [71, 76] |
| transformer, DeiT [77] | 2022 | X-ray | chest | [78] |
| transformer [79] | 2023 | X-ray | breast | breast cancer [80] |
| DeiT [81] | 2022 | X-ray | chest | - |
| DeiT [82] | 2022 | X-ray | breast | breast cancer [83] |
| ViT-Eff [32] | 2021 | X-ray | chest | lung disease [17, 84–86] |
| PneuNet [87] | 2023 | X-ray | lung | lung disease [88–97] |
| Chest L-Transformer [31] | 2022 | X-ray | chest | lung disease [98] |
| MXT [99] | 2022 | X-ray | chest | chest disease [100], [101] |
| multi-feature fusion transformer [102] | 2022 | X-ray | chest | lung disease [74, 89, 103–106] |
| MP-ViT [33] | 2022 | X-ray | chest | lung disease [46] |
| federated split transformer [107] | 2021 | X-ray | lung | lung disease [103, 108, 109] |
| transformer [110] | 2022 | Microscope | prostate | prostate cancer [111] |
| transformer, compact convolutional transformer [112] | 2021 | Microscope | colon | colorectal cancer [38] |
| explainable transformer-based [39] | 2022 | Microscope | cell | malaria parasite [113, 114] |
| ensembled swin transformer [115] | 2022 | Microscope | breast | breast tumor [116] |
| IMGL-VTNet [117] | 2022 | Microscope | gastric | gastric intestinal metaplasia [118] |
| AMIL-Trans [119] | 2022 | Microscope | breast | breast cancer [120] |
| Self-ViT-MIL [121] | 2022 | Microscope | breast | breast cancer [120] |
| TransPath [122] | 2021 | Microscope | breast, colon | breast cancer [120], colorectal cancer [123], polyps [124] |
| Fourier ViT [125] | 2022 | Microscope | breast | breast cancer [126] |
| RAMST [127] | 2022 | Microscope | gastrointestinal | - |
| CWC-Transformer [128] | 2023 | Microscope | breast, lung | breast cancer [120], breast cancer[129] |
| transformer [130] | 2022 | CT | lung | lung disease [37, 131] |
| transformer [132] | 2022 | CT | lung | lung disease [133] |
| transformer [134] | 2021 | CT | lung | - |
| transformer [135] | 2022 | CT | artery | - |
| transformer [136] | 2022 | CT | lung | lung disease [133] |
| transformer [30] | 2021 | CT | lung | lung disease [137, 138] |
| multi-view convolutional transformer [139] | 2022 | CT | lung | - |
| DenseTransformer [140] | 2022 | CT | lung | lung disease [141] |
| transformer-based factorized encoder [142] | 2022 | CT | lung | lung disease [143] |
| multi-granularity dilated transformer [144] | 2023 | CT | lung | lung disease [145] |
| transformer [146] | 2022 | MRI | pancreas | intraductal papillary mucosal neoplasms [147] |
| TransMed [148] | 2021 | MRI | head, neck, knee | anterior cruciate ligament, meniscal tears [149] |
| double-scale GAN [36] | 2021 | MRI | brain | [150] |
| MEST [151] | 2022 | MRI | brain | Parkinson's disease [152] |
| MIL-VT [153] | 2021 | DFI | retinal | retinal disease [154, 155] |
| VTGAN [156] | 2021 | DFI | retinal | retinal disease [157] |
| MVT-based framework [42] | 2022 | Dermoscopy | skin | pigmented skin lesion [158] |
| O-Net [50] | 2022 | Dermoscopy | skin | melanoma [159] |
| transformer [160] | 2022 | Endoscopy | gastrointestinal | gastrointestinal disease [161] |
| transformer [162] | 2022 | Endoscopy | colon | - |
| transformer [163] | 2022 | US | breast | breast disease [35, 164] |
| multi-scale feature fusion transformer [165] | 2022 | US | breast | - |
| Advit [166] | 2022 | PET | brain | Alzheimer's Disease [167] |
| multi-model transformer [44] | 2021 | Camera | toe | toe disease [168] |
| ViT-P [169] | 2021 | OCT | genitourinary | genitourinary syndrome [46] |
| SSBTN [170] | 2022 | X-ray, Microscope | breast, small intestine | breast cancer [171], Crohn's disease [172] |
| symmetric dual transformer [173] | 2022 | X-ray, CT | chest | lung disease [104, 133] |
| grouped bottleneck transformer [174] | 2022 | CT, MRI, Microscope | tooth, abdomen, chest, brain, synapse | [175, 176] |
| FPViT [177] | 2022 | Microscope, X-ray, Dermoscopy, US, CT, DFI | colon, chest, skin, chest, breast, abdomen, retinal | [175] |
| SEViT [178] | 2022 | X-ray, DFI | chest, retinal | retinal disease [154], [179] |

model from the inductive bias. Similar works include the explainable transformer-based model [39], in which the compact convolutional transformer and a gradient-weighted class activation map technique are implemented to show the attention paid to different parts by generating a heatmap. Some works made modifications based on the existing transformer model, such as the ensembled swin transformer proposed in 2022 [115]. The ensembled swin transformer averages the predicted vectors of all individual ones. Tandeming the CNN and transformer is also commonly used. For instance, the IMGL-VTNet [117] tandems the ResNet and the deformable transformer encoder. The deformable transformer encoder is composed of a multi-scale deformable attention module and a feed-forward network. In the multi-scale deformable attention module, the multi-scale deformable attention function is leveraged to produce the feature map via weighted average. Zhang et al. [119] proposed an AMIL-Trans network composed of two stages. In the first stage, features are captured using ResNet as well as the efficient channel attention module [185]. In the second stage, the transformer encoder for discriminant instance features takes the features as the input and outputs the prediction.

Instead of using the existing models or making minor modifications based on them, the remaining works constructed the novel models more deeply. Gul et al. [121] implemented a Self-ViT-MIL method, in which the transformer is first trained in a self-supervised manner using the DINO training approach. The multiple-instance learning aggregator is then trained with frozen transformer weights. Wang and co-workers [122] developed a TransPath model consisting of a CNN encoder, a transformer encoder, and a token-aggregating and excitation module. The proposed self-supervised model is trained using the BYOL. Duan et al. [125] constructed a Fourier ViT model consisting of two branches. The one branch is composed of two transformer encoders and their output information is exchanged with cross attention. Another branch normalizes the tokens and performs the 2D discrete Fourier transform. Elementwise multiplication is then implemented followed by the 2D inverse Fourier transform. The outputs of two encoders and the Fourier branch are concatenated before passing the fully connected layer. Lv and co-workers [127] constructed a RAMST, which can be further divided into the region-level RAMST and the WSI-level RAMST. Both region-level and WSI-level RAMST are composed of CNN and transformer while the WSI-level RAMST consists of an additional CNN branch. A novel feature weight uniform sampling method is also developed and implemented in both RAMSTs for patch subset sampling to preserve representative region features. In 2023, the [128] CWC-Transformer is proposed by Wang et al. to solve the problem of feature extraction and spatial information loss effectively. In the compression stage, a feature compression method is implemented to extract discriminative features and reduce data bias. During the learning phase, the strengths of CNN and the transformer are extended to enhance the interrelationship between local and global information.

**CT**. A majority of researchers utilize developed transformer models to classify medical CT images [30, 130, 132, 134–136]. A typical example is the medical diag-

nostic platform developed by Li et al. [134]. The platform is based on the transformer and can gain more medical information from the traditional image recognition model by distilling technology. There are a few works that contain novel transformer-based models. Xiong et al. [139] developed a multi-view convolutional transformer composed of four stages, which are view generation, visual backbone, feature decorrelation, and classifier. The visual backbone introduces the non-local self-attention into the last layer of ResNet, and the feature decorrelation learns a set of sample weights for eliminating the dependence between features. Mei [140] married CNN and the transformer and constructed the DenseTransformer. The CNN and transformer are combined in three ways, including CNN and transformer in parallel, transformer in front of CNN in series, and CNN in front of transformer in series. Huang and co-workers [142] proposed a transformer-based factorization encoder consisting of two transformer encoders. The former encoder enables the intra-slice interaction via encoding feature maps from the same slice, and the latter encoder investigates the inter-slice interaction via encoding feature maps from different slices. The multi-granularity dilated transformer [144] developed in 2022 leveraging the local focus scheme for guiding the deformable dilated transformer. The local focus scheme aims at discriminative local features more via modeling channel-wise grouped topology, and the deformable dilated transformer incorporates diverse contextual information.

**MRI**. Salanitri et al. proposed to leverage the transformer to diagnose intraductal papillary mucosal neoplasms [146]. Dai and co-workers [148] developed a TransMed method by connecting the ResNet and the transformer in series. The double-scale generative adversarial network (GAN) method proposed by Hu et al. [36] is composed of a generator and two discriminators. The local CNN-based discriminator guides the generator to capture structural representation with inductive bias, while the transformer-based global discriminator directs the generator to extract comprehensive features via leveraging long-range dependencies. The MEST framework [151] developed in 2022 uses pre-trained VGGNet [186] and attention mechanism to learn multi-plane dynamic images. Time-series information is used to construct dynamic functional connection images. Spatial-temporal connectivity transformer is utilized to solve spatiotemporal redundancy and dependencies, and ensemble learning is also employed to integrate multimodality data.

**DFI**. Yu and co-workers [153] developed a MIL-VT network. The MIL-VT uses the transformer as the backbone and a multiple-instance learning head is proposed to exploit the feature representations captured by individual patches better. The cross-entropy loss between the multilayer perceptron head and the label and between the multiple-instance learning head and label are computed. The VTGAN model [156] constructed by Kamran et al. is composed of a coarse and a fine generator as well as two transformers as discriminators. The generators synthesize images according to the input and synthesized images are fed to the transformer encodes for classification. The two encoders also determine whether the synthesized images are from input or artificially generated.

**Dermoscopy**. In 2022, Aladhadh et al. [42] developed an MVT-based framework
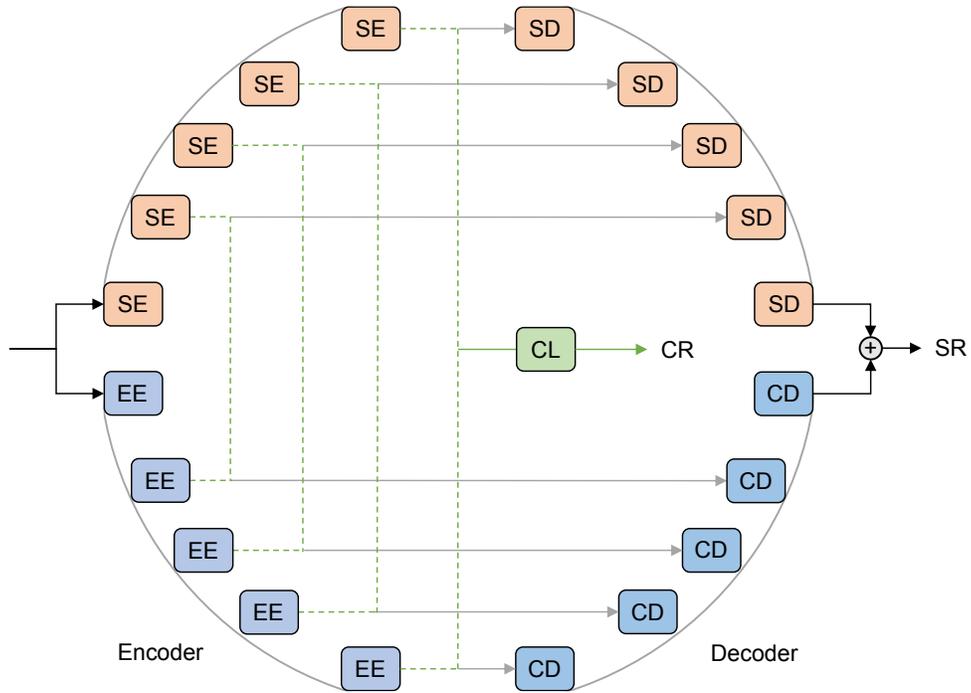
Figure 9: Structure of the O-Net. O-Net has a novel architecture design and serves as a universal model for both classification and segmentation. EE, SE, CD, and SD represent the EfficientNet encoder block, swin transformer encoder block, CNN decoder block, and swin transformer decoder block, respectively. The gray solid lines show the skip connection. The green dotted lines depict fusion and unification. CL means the classification layer. CR represents the classification result and the SR illustrates the segmentation result.

based on the transformer. Different data augmentation methods, including image flip, scaling, rotation, and contrast. Wang and co-workers designed a novel architecture O-Net, which serves as a universal model for both classification and segmentation. We show the complicated structure of the O-Net in Figure 9. The O-Net has a circle shape and is composed of four core blocks: EfficientNet encoder block, swin transformer encoder block, CNN decoder block, and swin transformer decoder block. Several EfficientNet encoder blocks and swin transformer encoder blocks composed the encoder, while the remaining two blocks compose the decoder. The encoder and decoder are connected by several skip connections.

**Endoscopy**. Hosain et al. [160] proposed to compare the classification performance between the transformer and CNN. Tamhane and co-workers [162] developed a landmark detection pipeline composed of three stages, including preprocessing, feature extracting, and classifying. The features are extracted using the transformer. CNN models including ResNet and ConvNeXt [187] are also implemented for performance comparison.

**US**. Gheflati and Rivaz [163] utilize the transformer and compare its performance with several CNN models. Li et al. [165] developed a multi-scale feature fusion transformer by combining CNN and transformer. Short-distance feature interaction block is designed for the two feature maps within the CNN block, while a long-distance feature interaction block is developed for the feature maps between stages. Cross-SE block is introduced in the transformer block, which is mainly composed of the global average pooling layer and fully connected layer.

**PET**. Xing and co-workers [166] proposed an Advit model composed of two branches processing different modalities of PET. For each branch, a 3D-to-2D operation is implemented to project the 3D PET images into 2D fusion images using the proposed CNN module. The fused 2D images are then forwarded to the transformer. The output of two transformers is finally concatenated.

**Camera**. Qayyum et al. [44] developed a multi-model transformer consisting of two separate pre-trained transformers. The outputs of both transformers are concatenated using pair-wise feature concatenation. The pair-wise feature concatenation is composed of two branches. In the first branch, the output of the second transformer is concatenated after the output of the first transformer. Regarding the second branch, the sequence is inverted. Outputs of two branches are then concatenated, in which the output of the first branch is placed in the front.

**OCT**. Wang and co-workers [169] developed a ViT-P architecture, consisting of a slim model and several transformer encoders in series. The model is mainly composed of four stages, in which each stage consists of the multi-branch convolutional and channel attention mechanism [188]. Besides, DCGAN [189] and proposed B-DCGAN are implemented to perform data augmentation. Though their method shows promising results, we have to point out that the images generated by GAN are not quantitatively evaluated using metrics like IS [190] and FID [191] thus causing the performance difficult to represent intuitively.

**Multiple**. In 2022, Gong et al. [170] developed an SSBTN composed of three

modules, which are the pretext channel module, the transformer-based transfer module, and the downstream channel module. The transformer-based transfer module uses bi-channel transformer encoders and the loss between two channels is calculated. Rahhal and co-workers [173] proposed a symmetric dual transformer consisting of two transformers. Original images are fed to one of the transformers with a class classifier while augmented images are sent to another transformer with a distill classifier. The outputs of the two transformers are then passed through a weighted fusion layer. Gao et al. [174] built a grouped bottleneck transformer. The grouped bottleneck transformer block is composed of two branches, consisting of convolution operation only and both convolution operation and multi-head self-attention mechanism, respectively. The FPViT network [177] developed by Liu et al. extracts feature from three different layers from ResNet and input them into the transformer heads at different scales. The activation vectors are obtained through three transformers and the ResNet head, and prediction results are made through concatenated vectors. Almalik and co-workers [178] proposed a SEViT model, which trains separate multilayer perceptron layers using extracted patch tokens. A self-ensemble of different multilayer perceptron layers, together with the transformer classifier, enhances the robustness of the transformer. Besides, the consistency between ensemble predictions is leveraged for detecting adversarial samples.

### 4.2. Segmentation

Segmentation-related works are also grouped by different modalities in the sequence of MRI, CT, endoscopy, X-ray, US, microscope, DFI, camera, and dermoscopy. For papers containing more than one modality, we place them into the "multiple" category following the dermoscopy. The transformer-based segmentation works are organized in Table 2, Table 3, and Table 4 according to different modalities, respectively. Most of the segmentation works combine the transformer with the U-Net [192] or its variants like TransUNet [193]. The U-Net is composed of a contracting path (encoder) following the structure of CNN, and an expansive path (decoder). The bottleneck layer is implemented between the encoder and decoder. The encoder is composed of several CNN blocks with the ReLU activation function. The output of each CNN block passes the max pooling for downsampling. With the downsampling, the number of feature channels is doubled. As for the decoder, it is composed of several upsampling, which halves the passed feature channels and several CNN blocks with the ReLU activation function. The features in the CNN blocks are concatenated with the feature obtained from the encoder at different scales. Specifically, the features obtained from a certain encoder scale are first cropped because border pixels are lost during convolution operation. The cropped features are then concatenated along the channel dimension.

**MRI**. The U-Net and its variants are widely used to segment medical MRI images. Within these works, most of the works aim at modifying the encoder, decoder, or the bottleneck layer [194, 196, 199, 203, 205, 206, 210, 212–219, 221, 223, 224, 226, 228]. Many works utilize several U-Nets when constructing their model [229–231,

Table 2: Transformer-based segmentation applications for CT and MRI modalities.

| Method | Year | Modality | Object | Dataset |
|---|---|---|---|---|
| UTNet [194] | 2021 | MRI | cardiac | [195] |
| MRA-TUNet [196] | 2022 | MRI | cardiac | cardiac disease [197, 198] |
| TransConver [199] | 2022 | MRI | brain | brain tumor [200–202] |
| UTransNet [203] | 2022 | MRI | brain | stroke [204] |
| TransBTS [205] | 2021 | MRI | brain | brain tumor [200–202] |
| METrans [206] | 2022 | MRI | brain | stroke [207–209] |
| SwinBTS [210] | 2022 | MRI | brain | brain tumor [200, 202, 211] |
| BTSwin-Unet [212] | 2022 | MRI | brain | brain tumor [200, 201] |
| AST-Net [213] | 2022 | MRI | brain | brain tumor [202] |
| BiTr-Unet [214] | 2022 | MRI | brain | brain tumor [202] |
| Swin UNETR [215] | 2022 | MRI | brain | brain tumor [202] |
| CST-UNET [216] | 2022 | MRI | brain | brain tumor [202] |
| VT-UNet [217] | 2022 | MRI | brain | brain tumor [202] |
| CSU-Net [218] | 2022 | MRI | brain | brain tumor [202] |
| OSTransnet [219] | 2022 | MRI | bone | osteosarcoma [220] |
| 3D PSwinBTS [221] | 2022 | MRI | brain | brain tumor [202, 222] |
| TSEUnet [223] | 2022 | MRI | brain | brain tumor [202] |
| RMTF-Net [224] | 2022 | MRI | brain | brain tumor [202, 225] |
| AMTNet [226] | 2023 | MRI | prostate, brain | brain tumor [202], [227] |
| transformer-based GAN [228] | 2022 | MRI | brain | brain tumor [202] |
| DUconViT [229] | 2022 | MRI | bone | - |
| CTCL [230] | 2022 | MRI | cardiac | cardiac disease [197] |
| symmetrical supervision transformer [231] | 2022 | MRI | abdomen, cardiac | [232, 233] |
| transformer-enhanced U-Net [234] | 2021 | MRI | cardiac | [235] |
| TransUNet-based [236] | 2022 | MRI | brain, cardiac | stroke[209] |
| dual-teacher [237] | 2022 | MRI | cardiac | cardiac disease [197] |
| mmFormer [238] | 2022 | MRI | brain | brain tumor [202] |
| NVTrans-UNet [239] | 2023 | MRI | cardiac | [240] |
| 3D transformer [241] | 2022 | MRI | brain | Alzheimer's [242] |
| iSegFormer [243] | 2022 | MRI | cartilage | [244] |
| transformer-based region-edge aggregation network [51] | 2022 | MRI | cardiac, knee | cardiac disease [198] |
| CESS-ViT [245] | 2022 | MRI | cardiac | cardiac disease [197] |
| uncertainty-aware transformer [246] | 2022 | MRI | cardiac | cardiac disease [197] |
| HybridCTrm [247] | 2021 | MRI | brain | [248], neurodevelopmental disorders [249] |
| statistical features-based [250] | 2022 | MRI | brain, cardiac | [222], brain tumor [200–202] |
| feature fusion-based [251] | 2023 | MRI | brain | brain tumor [200–202] |
| UNTER [252] | 2022 | CT | liver | liver tumor [253, 254] |
| CoTr [255] | 2021 | CT | abdomen | colorectal cancer, ventral hernia [256] |
| ITUnet [257] | 2022 | CT | head, neck | - |
| TFCNs [258] | 2022 | CT | abdomen, chest | colorectal cancer, ventral hernia [256], lung disease [91, 259] |
| TSE DeepLab [260] | 2022 | CT | sinus, patellar | - |
| transformer-UNet [261] | 2021 | CT | lung | [227] |
| AFTer-UNet [262] | 2022 | CT | abdomen, chest | colorectal cancer, ventral hernia [256], organs at risk [263, 264] |
| HT-Net [265] | 2022 | CT | lung, kidney, bladder | kidney tumor [266], lung lesion [267], bladder cancer [268] |
| UCATR [269] | 2021 | CT | brain | - |
| MMViT-Seg [270] | 2023 | CT | lung | lung disease [271, 272] |
| CCAT-net [273] | 2022 | CT | chest | lung disease [37] |
| CAC-EMVT [274] | 2021 | CT | chest | - |
| MSHT [275] | 2021 | CT | liver, kidney | kidney tumor [266], liver tumor [253] |
| RCSHT [276] | 2022 | CT | chest | - |
| design-flexible transformer [277] | 2022 | CT | liver, spine | liver tumor [222], [278] |
| MAPTransNet [279] | 2022 | CT | lung | lung tumor [106, 280] |
| CTUNet [281] | 2022 | CT | pancreas | [282] |

234, 236, 237]. For instance, Xiao et al. [237] developed a semi-supervised dual-teacher architecture, which uses simultaneous dual-teacher to guide the student. The two teachers use U-Net and Swin-UNet [283] as the backbone and the student uses the U-Net as the backbone. The U-Net teacher and U-Net student, and two teachers inside are screened for uncertainty assessment during training. Beyond single modal, there are also several works aiming at multimodal tasks [238, 239], such as the multimodal network NVTrans-UNet [239] developed by Li et al. The input of the NVTrans-UNet is composed of three main parts, including the encoder module, bottleneck layer, and decoder module. The NVTrans-UNet utilizes the neighborhood transformer to localize the receptive field of each token to its nearest neighboring pixel. The multi-modal gated fusion strategy is implemented to adjust the contribution of feature mapping from each modal. Atrous spatial pyramid pooling is also used in the bottleneck layer for expanding the receptive field, reducing parameters, and enhancing extraction ability.

Besides the U-Net-based work, there are several models developed in other ways. Several works made minor modifications based on the existing model or connect two models in series. For example, Karimi et al. [241] developed a 3D transformer, in which the residual connection in the transformer encoder block is removed. Liu and co-workers [243] proposed an iSegFormer, where the swin transformer and lightweight multilayer perceptron decoder are combined in series. Some researchers design their models from scratch [51, 245–247, 250, 251]. For instance, Chen et al. [51] proposed a transformer-based region-edge aggregation network, where the multi-level region and edge features are aggregated by multiple transformer-based inference modules to form multi-level complementary features. These complementary features are utilized to guide the decoding of the corresponding level region and edge features. Sun and co-workers [247] developed a multimodal HybridCTrm network, which is composed of two paths. The first path takes the MRI-T1 and MRI-T2 images together, followed by the parallel CNN and transformer, while the second path takes the MRI-T1 and MRI-T2 images separately. In 2023, a novel method based on deep semantics and edge information fusion is developed [251]. The proposed method is composed of a semantic segmentation module, an edge detection module, as well as a feature fusion module. The segmentation module utilizes the swin transformer as the backbone with shifted patch tokenization strategy. The CNN-based detection module consists of the proposed edge spatial attention block for feature enhancement. Semantic and edge features from two modules are fused by the feature fusion module.

**CT**. Instead of using existing models to segment CT images, such as the framework [252] using UNTER [284] as the backbone, major works use U-Net-based novel networks with different aspects of modifications [255, 257, 258, 260–262, 265, 269, 270, 273–277, 279, 281]. For instance, Xie et al. [255] designed a framework CoTr with three parts: A CNN encoder, a DeTrans encoder, and a Decoder. The DeTrans encoder connects the CNN encoder and the decoder at different scales. The output of the CNN encoder is flattened before feeding to the DeTrans

Table 3: Transformer-based segmentation applications for endoscopy, X-ray, US, microscope, DFI, camera, and dermoscopy modalities.

| Method | Year | Modality | Object | Dataset |
|---|---|---|---|---|
| RANT [287] | 2022 | Endoscopy | throat | [288] |
| BiDFNet [289] | 2022 | Endoscopy | colon | polyp [40, 290–293] |
| Patcher [294] | 2022 | Endoscopy | colon | polyp [290] |
| Polyp2Seg [295] | 2022 | Endoscopy | colon | polyp [40, 290–293] |
| TransHarDNet [296] | 2022 | Endoscopy | colon | polyp [40, 290–293] |
| FCBFormer [297] | 2022 | Endoscopy | colon | polyp [290, 298] |
| U-Net [299] | 2021 | X-ray | breast | - |
| temporary transformer [300] | 2022 | X-ray | catheter | - |
| APSegmenter [301] | 2022 | X-ray | spine | spinal curvature [302] |
| Chest L-Transformer [31] | 2022 | X-ray | chest | lung disease [98] |
| federated split transformer [107] | 2021 | X-ray | lung | lung disease [303] |
| TransBridge [304] | 2021 | US | cardiac | cardiac disease [41] |
| TFNet [305] | 2022 | US | breast, thyroid | breast disease [35], thyroid disorder [306] |
| CSwin-PNet [307] | 2023 | US | breast | breast disease [35, 164] |
| RSTUnet-CR [308] | 2022 | US | breast | - |
| dilated transformer [309] | 2022 | US | breast | breast tumor [310] |
| Swin-PANet [311] | 2022 | Microscope | colon, cell | colon cancer [312], [313] |
| multiple-instance transformer [314] | 2022 | Microscope | colon | colon cancer [315] |
| SMESwin Unet [316] | 2022 | Microscope | colon, cell | colon cancer [312], [313], [317] |
| PCAT-UNet [43] | 2022 | DFI | retinal | retinal disease [318, 319], [320] |
| Polarformer [321] | 2022 | DFI | retinal | retinal disease [322–324] |
| GT-DLA-dsHFF [325] | 2022 | DFI | retinal | retinal disease [318, 320, 326], [327] |
| versatile transformer [328] | 2022 | Camera | skin | - |
| semi-supervised transformer [329] | 2022 | Dermoscopy | skin | melanoma [159, 330, 331] |

encoder and the output of the DeTrans encoder is reshaped and then send to the decoder. The DeTrans only pays attention to a small set of key positions thus the complexity is reduced largely. Kan and co-workers [257] developed an ITUnet, in which the feature map of CNN and transformer are added in the downsampling stage. In the upsampling stage, the segmentation predictions for each feature map obtained by the up block are generated and utilized to calculate the loss. Li et al. [258] constructed a TFCNs model, in which the encoder is constructed by introducing the transformer into the FC-DenseNet [285]. The RL-Transformer layer is added at the end of the encoder and the convolutional linear attention block is introduced in the skip connection to filter non-semantic features by including spatial and channel attention. In 2023, Yang and co-workers [260] developed a TSE DeepLab framework, which leverages atrous convolution in DeepLabv3 [286] as the backbone to extract features. The captured features are then converted into visual tokens and then fed to the transformer. Squeeze and excitation components are also introduced after the transformer for channel importance sorting.

**Endoscopy**. Pan et al. [287] proposed a RANT framework, in which the transformer and CNN are combined in series. Features are cascaded using reverse attention and receptive field block module. The segmentation results are optimized using convolutional conditional random fields. Tang and co-workers [289] proposed a bi-decoder BiDFNet works in both fine-to-coarse and coarse-to-fine ways. The

BiDFNet is composed of an encoder based on PVTv2 [332] as well as two decoders connected in series. The adaptive fusion module and the residual connection module are implemented in the decoders, and the adaptive fusion module aggregates the features from different scales effectively. Ou et al. [294] developed a Patcher method, where the encoder utilizes a cascade of Patcher blocks for expert features capture at different scales. The Patcher block first segment the input to large patches with overlapping contexts and then further divided them into small patches. The divided small patches are next fed to sequential transformer blocks for feature extraction and the large patches are finally reassembled. The mixture-of-experts-based decoder utilizes a gating network to filter a set of suitable expert features for the prediction. Mandujano and co-workers [295] constructed a Polyp2Seg network, which uses PVTv2 to extract a set of multi-scale features. The extracted multi-scale features are then compressed and fed into several feature aggregation modules. A multi-context attention module is implemented to characterize low-level polyp cues and the final predicted results are obtained by several auxiliary outputs. The TransHarDNet network [333] designed in 2022 combines the transformer and HarDNet blocks [334]. HarDNet Blocks are leveraged to extract spatial and depth information, while the transformer captures global semantic context information. Several cascaded partial decoders are implemented to fuse the feature maps and the skip connection with the receptive field block is implemented between the HarDNet blocks and partial decoders. Sanderson et al. [297] designed an FCBFormer architecture consisting of two branches. The transformer branch extracts the most important features based on the transformer, while the fully convolutional branch is implemented as a supplementary. The output of the two branches is then concatenated and passes the prediction head.

**X-ray**. Saidnassim and co-workers [299] proposed to use the BYOL algorithm for U-Net-based breast image segmentation. Zhang et al. [300] proposed a temporary transformer network, which takes both the current and previous frames as the input to obtain temporary information. The current frame is fed into the CNN and transformer, while the previous frame is fed into the CNN only. In 2022, Zhang and co-workers [301] developed an APSegmenter method in which the transformer-based Segmenter [20] is utilized to obtain semantic segmentation results. The proposed adaptive post-processing module is utilized to optimize the results, which takes the vertebral block boundary in the adhesion region as the input and outputs the vertebral mass without adhesion. The Chest L-Transformer and federated split transformer discussed in Table 1 are also leveraged for medical image segmentation.

**US**. Deng et al. [304] proposed a TransBridge model, in which both the encoder and decoder are based on CNN, while the transformer encoder is used to skip-connect them at different scales. Within the transformer encoder, an embedding layer is implemented by using shuffled group convolution and dense patch division. Wang and co-workers [305] constructed a TFNet model, where the channel attention mechanism is introduced for solving the channel modeling defect. A loss function based on KL distance is also proposed to modify the predicted results by
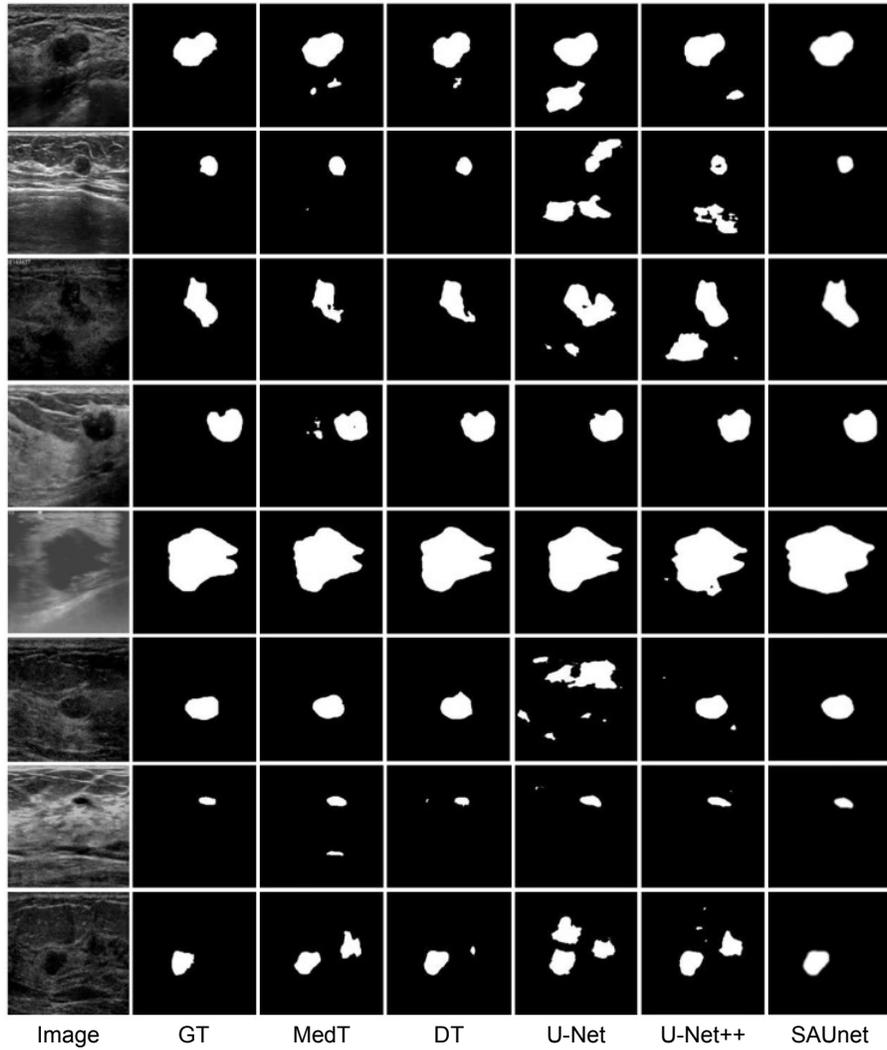
Figure 10: Segmentation results on the BUSIS dataset [310] using MedT [335], dilated transformer (DT) [309], U-Net [192], U-Net++ [336], and SAUnet [337]. GT represents the ground truth.

25

calculating the variance between the results of the main classifier and the auxiliary classifier. In 2023, Yang et al. [307] built a CSwin-PNet model. An interactive channel attention module using channel-wise attention and an SFF module is developed for feature region emphasize and feature supplementary during fusion, respectively. Besides, a boundary detection module is also utilized to extract the boundary information. Zhuang and co-workers [308] designed an RSTUnet-CR model consisting of a shared encoder, a segmentation decoder, and a consistency regularization decoder where long-distance dependence is established using the residual swin transformer block. The dilated transformer model proposed in 2022 [309] uses a dilation convolution block to connect the encoder and decoder. The encoder contains the multi-head attention mechanism and the decoder is mainly composed of deconvolution. As shown in Figure 10, the dilated transformer performs better compared with other state-of-the-art methods. Among all methods, the dilated transformer and the SAUnet [337] perform better due to low false positives, meaning that the boundaries can be distinguished precisely. Within the two methods, the dilated transformer can capture information in more detail. Take the fourth row as an example, the dilated transformer can distinguish the invaginated part at the top tumor better. Though the dilated transformer outperforms other models, we also observed that it can sometimes produce unideal results. Several examples would be the first, seventh, and eighth rows, in which an isolated extra object is mistakenly created. This means the ability to differentiate textures with subtle differences still has a margin to be improved.

**Microscope**. In 2022, Liao et al. [311] proposed a Swin-PANet model following the coarse-to-fine as well as dual supervision strategy. The developed Swin-PANet consists of a prior attention network and a hybrid transformer network. The swin transformer-assisted prior attention network carries out intermediate supervision learning, while the hybrid transformer network with enhanced attention blocks implements direct learning. Besides, the skip connection is employed to connect the encoder and decoder of the hybrid transformer network. Qian and co-workers [314] developed a multiple-instance transformer where the transformer is incorporated into the multiple-instance learning framework. The self-attention establishes the relationship among different instances. Deep supervision is implemented to overcome the annotation limitation existing in weakly-supervised methods. Wang et al. [316] developed a SMESwin Unet model based on their proposed MCCT. The MCCT is designed to fuse multi-scale semantic features and attention maps based on the channel-wise cross-fusion transformer [338]. Superpixel is introduced by dividing the pixel-level feature into district-level and external attention is leveraged to introduce the correlations among all data samples.

**DFI**. Chen and co-workers [43] designed a PCAT-UNet model containing two main components named patches convolution attention transformer block and feature grouping attention module. Both encoder and decoder are composed of several patches convolution attention transformer blocks and the outputs of the feature grouping attention modules are fed into the patches convolution attention transformer

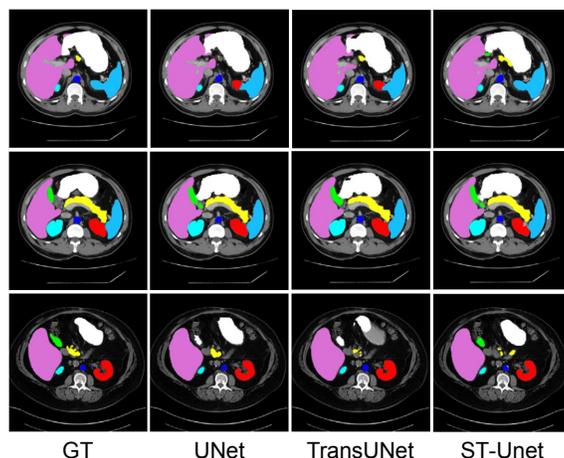GT      UNet     TransUNet    ST-Unet

Figure 11: Segmentation results on the Synapse dataset [256] using UNet, TransUNet [193], and ST-Unet [376].

blocks. The segmentation map of each layer is predicted using the fused enhanced feature map. Feng et al. [321] proposed a Polarformer network mainly composed of the learnable polar transformation module as well as the CNN-transformer module. The polar transformation module carries out a differentiable log-polar transform, while the CNN-transformer module captures features and consolidates global attention. A segmentation head is implemented to output the confidence scores and transmute the predictions back to the Cartesian coordinate system. Li and co-workers [325] developed a GT-DLA-dsHFF model, in which a global transformer and dual local attention network are introduced for global information integration and local vessel information extraction, respectively. Besides, a deep-shallow hierarchical feature algorithm is used to fuse features.

**Camera**. The versatile transformer [328] developed by Junayed et al. is composed of the dual encoder, the feature versatile block, and efficient decoder architecture with skip connections. The dual encoder is based on CNN and transformer to extract features, and the feature versatile block is implemented to distribute and integrate obtained features between the encoder and decoder. A squeeze and excitation block component is also introduced in the decoder to capture channel-wise dependencies as well as the significant feature correlations.

**Dermoscopy**. Alahmadid and co-workers [329] proposed a transformer consisting of a supervised stream and an unsupervised stream. The supervised stream combines CNN and transformer and the output features of CNN and transformer are fused. Specifically, the transformer output is reshaped into the same spatial dimension as the CNN, and then two features are concatenated. The fused features are then fed to the decoder module for semantic segmentation learning. The unsupervised stream is composed of a supplementary decoding head and utilizes the unsupervised technique for encoder module enrichment. A surrogate task is designed on top of the CNN and transformer representations.

27

Table 4: Transformer-based segmentation applications for multiple modalities.

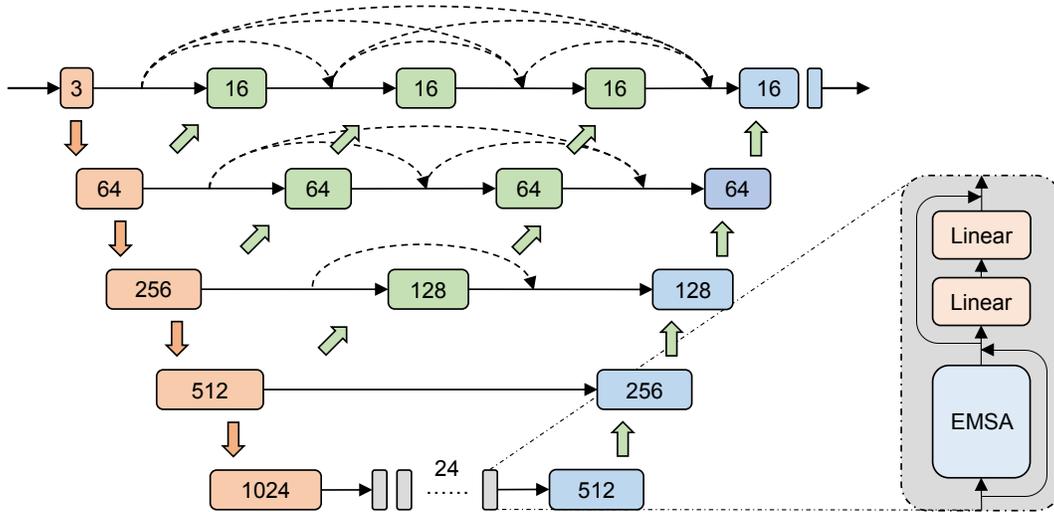| Method | Year | Modality | Object | Dataset |
|---|---|---|---|---|
| MISSFormer [339] | 2022 | CT, MRI, DFI | abdomen, cardiac, retinal | colorectal cancer, ventral hernia [256], cardiac disease [340], retinal disease [318] |
| Dual encoder transformer-CNN [341] | 2022 | CT, MRI | abdomen, cardiac | colorectal cancer, ventral hernia [256], cardiac disease [340] |
| ConTrans [342] | 2022 | Endoscopy, Dermoscopy, CT, Microscope | cell, skin, chest, colon | polyp [40, 290, 291, 293], melanoma [159, 330], pigmented skin lesion [158], lung disease [343], colon cancer [312], [344], cancer [345] |
| ScaleFormer [346] | 2022 | CT, MRI, Microscope | abdomen, cell, cardiac | [313, 347], cardiac disease [197] |
| EMSViT [348] | 2022 | MRI, CT | abdomen, brain | colorectal cancer, ventral hernia [256], brain tumor [201, 202] |
| TransCUNet [349] | 2022 | Microscope, Endoscopy, Dermoscopy | colon, cell, skin | colon cancer [312], polyp [291], [313, 350], melanoma [159] |
| CATS [351] | 2022 | CT, MRI | abdomen, brain, prostate | brain tumor [202], vestibular schwannomas [352], [353] |
| D-former [354] | 2022 | CT, MRI | abdomen, cardiac | brain tumor [202], cardiac disease [197] |
| APT-Net [355] | 2022 | Dermoscopy, Endoscopy, Microscope | skin, colon | melanoma [159, 331], polyp [40, 290–293], colon cancer [312] |
| TransNorm [356] | 2022 | CT, Dermoscopy, Microscope | abdomen, skin, cell | colorectal cancer, ventral hernia [256], melanoma [159, 330, 331], myeloma [357] |
| SwinPA-Net [358] | 2022 | Colonoscopy, Microscope, Camera | colon, cell | polyp [40, 290–293], [344] |
| GPA-TUNet [359] | 2022 | CT, MRI | abdomen, cardiac | colorectal cancer, ventral hernia [256], cardiac disease [197] |
| ConvWin-UNet [360] | 2023 | Microscope, CT | kidney, abdomen | colorectal cancer, ventral hernia [256], [361] |
| PCT [362] | 2023 | US, Microscope, Dermoscopy | parotid, skin, cell | [313], melanoma [159] |
| DS-TransUNet [363] | 2022 | Endoscopy, Dermoscopy, Microscope | colon, skin, cell | colorectal cancer [40], colon cancer [312], polyp [290–293], [344], melanoma [330] |
| DSTUNet [364] | 2022 | MRI, CT | abdomen, cardiac | cardiac disease [197], colorectal cancer, ventral hernia [256], cardiac disease [340] |
| MT-UNet [365] | 2022 | CT, MRI | abdomen, cardiac | colorectal cancer, ventral hernia [256], cardiac disease [340] |
| ViTBIS [366] | 2021 | CT, MRI | abdomen, brain | brain tumor [201, 202], colorectal cancer, ventral hernia [256] |
| TDD-UNet [367] | 2022 | CT, X-ray | lung | lung disease [88, 272] |
| SwinE-Net [368] | 2022 | Endoscopy, MRI | colon, brain | polyp [40, 290–293] |
| USegTransformer [369] | 2022 | Dermoscopy, MRI, CT, Microscope | brain, lung, cell, skin, chest | pigmented skin lesion [158], lung lesion [267], brain tumor [225], [344], melanoma [330], lung disease [343] |
| SegTransVAE [370] | 2022 | CT, MRI | kidney, brain | kidney tumor [266], brain tumor [202] |
| MedT [335] | 2021 | US, Microscope | brain, colon, cell | intraventricular hemorrhage [371, 372], colon cancer [312], [313, 350] |
| medical transformer [373] | 2023 | MRI, US, Camera | prostate, cardiac, tongue | [374] |
| CTC-Net [375] | 2023 | CT, MRI | abdomen, cardiac | colorectal cancer, ventral hernia [256], cardiac disease [197] |
| ST-Unet [376] | 2023 | Dermoscopy, CT | skin, abdomen | colorectal cancer, ventral hernia [256], pigmented skin lesion [158], melanoma [330] |
| MS-TransUNet++ [377] | 2022 | MRI, CT | prostate, liver | liver tumor [253], prostate cancer [378] |
| O-Net [50] | 2022 | Dermoscopy, CT | skin, abdomen | melanoma [159], colorectal cancer, ventral hernia [256] |
| TMSS [379] | 2022 | PET, CT | head, neck | tumor [380] |
| TransDeepLab [381] | 2022 | CT, Dermoscopy | abdomen, skin | colorectal cancer, ventral hernia [256], melanoma [159, 330, 331] |
| transformer [382] | 2023 | CT, MRI | abdomen | colorectal cancer, ventral hernia [256], [232] |
| ECT-NAS [383] | 2021 | CT, MRI | abdomen, cardiac | cardiac disease [197], [384], [232] |
| SMIT [385] | 2022 | CT, MRI | abdomen | colorectal cancer, ventral hernia [256] |
| progressive sampling transformer [386] | 2022 | Microscope, Endoscopy | colon, cell | colon cancer [312]; polyp [291], [313, 350] |
| X-Net [387] | 2021 | Microscope, Endoscopy | colon, cell | [344, 388], polyp [290] |

Figure 12: The structure of the MS-TransUNet++ [377]. The representative TransUNet++ implements the transformer blocks into the U-Net intuitively and uses multiple skip-connection to bridge the features at different resolutions. The orange, green, and blue blocks illustrate CNN-based blocks. The gray blocks represent the transformer layers consisting of efficient multi-head self-attention (EMSA). The orange arrows show the downsampling operation, the green arrows represent the upsampling operation, and the dotted lines show the skip connection.

**Multiple**. A large part of the work for coping with several modalities is U-shaped or its variants [50, 335, 339, 341, 342, 346, 348, 349, 351, 354–356, 358–360, 362–370, 373, 375, 377, 379, 387]. Yuan et al. [375] developed a CTC-Net, where two encoders are designed by the swin transformer and residual CNN to capture complementary features. The cross-domain fusion block is used to concatenate these features. The correlation between features from the ResNet and transformer domains is calculated and channel attention is employed to extract dual attention information. A feature complementary module is constructed by incorporating cross-domain fusion, feature correlation, and dual attention. Zhang et al. developed an ST-Unet [376] which leverages the swin transformer to extract features. Features of each encoder stage are then enhanced by the developed CLFE module and concatenated with the current ones, followed by the up-sampling. The CLFE utilizes the self-attention block to learn the global feature information of a certain layer and fuse and learn the information with the tokens of the previous layer to obtain the enhanced multi-layer feature information. Cross-layer features are finally obtained for decoding feature enhancement. The segmented results across ST-Unet and other models can be found in Figure 11. It can be seen that for images with clear visual and semantic relationships, U-Net, TransUNet, and ST-Unet exhibit accurate segmentation. However, ST-Unet outperforms other methods for images with discreet visual relationships due to better global context encoding and semantic discrimination, and other

29

methods can perform over- and under-segmentation. It is worth noting that due to insufficient semantic information and blurred boundaries, sometimes all the above-mentioned methods cannot produce outstanding results. An example of this would be the pancreas in the first and third rows. Wang and co-workers [377] designed a representative architecture MS-TransUNet++. The MS-TransUNet++ implements the transformer blocks into the U-Net intuitively and uses multiple skip-connection to bridge the features at different resolutions. We show the structure of the MS-TransUNet++ in Figure 12 for a more intuitive explanation. The MS-TransUNet++ has a U-shape and several constructed transformer layers are inserted between the encoder and decoder. Skip connections on different feature scales are implemented densely across different CNN blocks. The O-Net introduced in Table 1 can also segment images using a separate output head, as shown in Figure 9. The end-to-end multimodal TMSS network [379] proposed in 2022 is composed of a transformer encoder and CNN decoder. The transformer encoder takes the projected features from multimodal images and electronic health records.

There are a few works that are not based on U-Net [381–383, 385, 386]. For instance, Jiang et al. [385] proposed an SMIT method to perform self-supervised learning for the transformer. The proposed method combines a dense pixel-wise regression pretext task with masked patch token distillation. Two transformers are utilized in the proposed method, serving as student and teacher, respectively. The parameters of teacher network parameters are updated using an exponential moving average with momentum. In 2022, Jiang and co-workers [386] introduced a progressive sampling transformer, in which a gated position-sensitive axial attention mechanism is introduced in the attention module. Iterative sampling for sampling position updating is also added to ensure the attention stays on the region to be segmented.

### 4.3. Miscellaneous

Miscellaneous works are discussed in a sequence of captioning, registration, detection, enhancement, localization, and synthesis. It is worth noting that we do not further divide the included works by different imaging modalities due to the small number of works. The summary of the transformer-based miscellaneous applications can be found in Table 5.

**Captioning**. Most captioning works modify different parts of the transform from NLP [389, 391, 392, 394, 396, 397, 399]. For instance, Hou et al. [392] proposed a RATCHET network, in which the image is fed into a DenseNet-based encoder and its output passes the masked multi-head attention module. The text tokens are fed for embedding and then pass the transformer decoder. Their method performs well and is capable of generating correct keywords for the given images, as shown in Figure 13. Li and co-workers [397] proposed a CGT network, which can restore a sub-graph from clinical relation. The restored triples are injected into the visual features as prior knowledge to drive the decoding procedure. Then, the visible matrix is utilized to limit the impact of knowledge during encoding. Reports are predicted by the encoded cross-modal features via a transformer decoder. In 2022,

Table 5: Transformer-based miscellaneous applications.

| Task | Method | Year | Modality | Object | Dataset |
|---|---|---|---|---|---|
| Captioning | transformer [389] | 2022 | X-ray | chest | [390] |
| Captioning | CEDT [391] | 2022 | X-ray | chest | chest disease [86], [390] |
| Captioning | RATCHET [392] | 2021 | X-ray | chest | [393] |
| Captioning | TranSQ [394] | 2022 | X-ray | chest | [390, 395] |
| Captioning | multicriteria supervised transformer [396] | 2022 | X-ray | chest | [390, 395] |
| Captioning | CGT [397] | 2022 | DFI | retinal | [398] |
| Captioning | KdTNet [399] | 2022 | Endoscopy, X-ray | gastrointestinal, chest | [390] |
| Captioning | SGT [400] | 2022 | Endoscopy | kidney, small intestine | [401] |
| Captioning | MCGN [333] | 2022 | X-ray | chest | [395] |
| Captioning | Eddie-Transformer [402] | 2022 | X-ray | chest | chest disease [86], [390], lung disease [403] |
| Registration | TD-Net [34] | 2022 | MRI | brain | Alzheimer's [404] |
| Registration | SymTrans [405] | 2022 | MRI | brain | Alzheimer's [404] |
| Registration | TransMorph [406] | 2022 | MRI, CT | brain, chest, abdomen, pelvis | [150, 407] |
| Registration | FTNet [408] | 2022 | MRI | brain | Alzheimer's [404], [409] |
| Registration | Swin-VoxelMorph [410] | 2022 | MRI | brain | Alzheimer's [411], Parkinson's disease [152] |
| Registration | Xmorpher [412] | 2022 | CT | cardiac | [413, 414] |
| Registration | C2FViT [415] | 2022 | MRI | brain | Alzheimer's [404], [416] |
| Detection | swin transformer [417] | 2023 | X-ray | breast | [418] |
| Detection | DETR-based [419] | 2023 | Microscope | cell | [420] |
| Detection | NucDETR [421] | 2022 | Microscope | cell | [422], cancer [423] |
| Detection | lightweight transformer [424] | 2022 | X-ray | breast | breast cancer [425] |
| Detection | MS Transformer [426] | 2022 | CT | lung, bone, kidney, lymph | pulmonary nodules, bone lesions, kidney lesions, lymph node enlargement [427] |
| Detection | SFOD-Trans [428] | 2022 | CT | vein | [429] |
| Detection | federated split transformer [107] | 2021 | X-ray | lung | lung disease [95] |
| Enhancement | SSTrans-3D [430] | 2023 | CT | brain | - |
| Enhancement | 3D CVT-GAN [45] | 2022 | PET | brain | - |
| Enhancement | GVTrans [431] | 2021 | MRI | brain | [150, 432] |
| Enhancement | RSTUnet-CR [308] | 2022 | US | breast | - |
| Enhancement | TED-Net [433] | 2021 | CT | liver | metastatic lesion [434] |
| Enhancement | SIST [435] | 2022 | CT | head, chest, abdomen, spine, lung | acute cognitive or motor deficit, high-risk for pulmonary nodules, metastatic liver lesions [436] |
| Enhancement | Eformer [437] | 2022 | CT | liver | metastatic lesion [434] |
| Localization | transformer graph network [438] | 2022 | CT | artery | coronary plaque [439] |
| Synthesis | ResViT [440] | 2021 | MRI, CT | brain, pelvis | [150], brain tumor [202], [441] |

| Image | True text | Predicted text |
|---|---|---|
| | Stable right greater than left upper lobe fibrotic changes. New opacity in the left mid-to-low lung raises concern for infectious process versus possibly asymmetric edema. Recommend follow up to resolution. | Diffuse bilateral parenchymal opacities, similar compared to the prior exam, with new focal opacity in the left upper lung field. Findings could reflect multifocal infection, though a component of pulmonary edema is also possible. |
| | Cardiomegaly and pulmonary edema which may have progressed since prior although some changes may be accounted for by lower lung volumes on the current exam. Left basilar opacity, potentially atelectasis noting that infection would also be possible. | 1. Low lung volumes with bibasilar atelectasis. 2. Severe cardiomegaly. |
| | Known lung metastases are again noted though better assessed on prior CT. No definite signs of superimposed acute process. | No acute cardiopulmonary process. |
| | In comparison with the study of _ _ _, there is little change in the substantial enlargement of the cardiomediastinal silhouette and moderate pulmonary edema with bilateral pleural effusions. Monitoring and support devices remain in place. | 'As compared to the previous radiograph, there is no relevant change. Moderate cardiomegaly with bilateral pleural effusions and subsequent areas of atelectasis. The monitoring and support devices are in constant position. No new parenchyma opacities.' |

Figure 13: Reports generated by RATCHET [392] using the X-ray dataset [393]. The same color shows the corresponding descriptions.

KdTNet [399] is developed, in which the visual grid and graph convolutional modules are designed to extract fine-grained visual features. The transformer decoder is implemented to generate the hidden semantic states. A BERT-based auxiliary language module is used to obtain the context language features from the pre-defined medical term knowledge. Besides, a multimodal information fusion module is constructed to calculate the contribution of linguistic and visual features adaptively. Several works construct the models using smaller units. For example, Lin et al. [400] designed an SGT network, in which relation-driven attention is proposed to facilitate the interaction described in the report. Instead of directly leveraging the inputs traditionally, relation-driven attention utilizes diverse sampled interactive relationships as augmented memory. Besides, an ingenious approach is also developed to homogenize the input heterogeneous scene graph, in which graph-induced attention is injected into the encoder for local interactions encoding. Wang and co-workers [333] proposed an MCGN method, in which a memory-augmented sparse attention block with bilinear pooling is developed for extracting higher-order interactions. The Eddie-Transformer developed by Nguyen et al. [402] decouples the latent visual features into semantic disease embeddings and disease states using the proposed state-aware mechanism. The learned diseases and corresponding states are entangled into explicit and precise disease representations.

**Registration**. U-shaped networks are used in most of the works for medical image registration [34, 405, 406, 408, 410, 412]. For example, Shi and co-workers

[412] developed a new transformer architecture XMorpher with dual parallel feature extraction networks. The XMorpher exchanges information through cross-attention to discover multi-level semantic correspondence. At the same time, respective features are captured gradually for registration. The cross-attention transformer blocks can find the correspondence automatically and prompt the feature fusion. There is a work that uses a different way to design the model. Mok et al. [415] proposed a C2FViT method, which naturally leverages the global connectivity and locality of the convolutional transformer and the multi-resolution strategy to learn the global affine registration. The proposed C2FViT is divided into three stages to the affine registration in a coarse-to-fine manner. The three stages have an identical architecture, including a convolutional patch embedding layer and several transformer encoder blocks. For any transformer encoder block, it is composed of an alternating multi-head self-attention module with a convolutional feed-forward layer.

**Detection**. Several works made minor modifications based on existing models for medical image detection [417, 419]. One example would be the modified DETR-based [17] proposed by Leng et al. [419]. The authors introduce the PVT [332] and deformable attention module into the DETR. Connecting existing models in series is also widely used [421, 424, 426]. For instance, Zhang and co-workers [424] proposed a lightweight transformer for tumor detection [425], in which images are fed into a ResNet to generate feature maps. The proposed method employs attention to the outputs of ResNet to improve the hidden representations. The outputs are then fed to FPN [442], where the multi-scale pyramidal hierarchy is utilized to construct feature pyramids. A semi-supervised method is also introduced in the detection task. Liu et al. [428] proposed a semi-supervised framework SFOD-Trans, consisting of two parallel branches. The two branches in the SFOD-Trans are utilized to train supervised and unsupervised loss, respectively. The combination of the two branches results in a semi-supervised loss. Besides, a new fusion module named normalized ROI fusion (NRF) is designed for fusing the hepatic portal vein information captured from labeled and unlabeled images. The NRF extracts the ROI of the object region by calculating the geometric gravity center of the bounding box using real and artificial labels. The obtained two ROIs are fused using the MixUp [443] method. The federated split transformer discussed in Table 1 and Table 3 can also be used for image detection.

**Enhancement**. Reconstruction is one of the most widely researched enhancement tasks. Xie and co-workers [430] proposed a network SSTrans-3D, which reconstructs the volume using a slice-by-slice scheme. The structures of the encoder and decoder are the same as the ones in the transformer, while the normalization layers are removed. In 2022, Zeng et al. [45] developed a 3D convolutional visual transformer-GAN model 3D CVT-GAN. A hierarchical generator is designed where multiple 3D CVT blocks are used as the encoder and TCVT blocks are implemented as the decoder. The CVT block is based on convolutional embedding and the transformer, while the TCVT block is based on transpose convolutional embedding as well as the transformer. The discriminator is also based on the 3D CVT block. Korkmaz

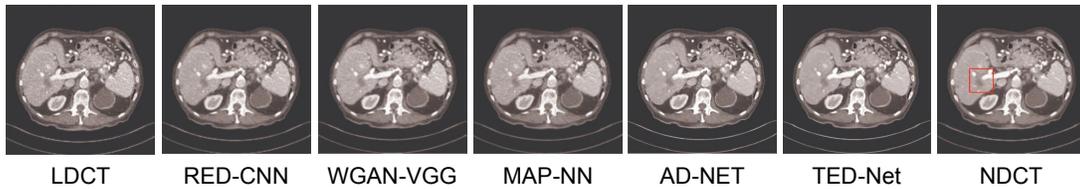| LDCT | RED-CNN | WGAN-VGG | MAP-NN | AD-NET | TED-Net | NDCT |

Figure 14: The denoising results of different methods on the CT dataset [434]. Included methods are RED-CNN [444], WGAN-VGG [445], MAP-NN [446], AD-NET [447], and TED-Net. LDCT represents low-dose CT, and NDCT illustrates normal-dose CT.

and co-workers [431] proposed a novel GVTrans deep generative network. It realizes scan-specific reconstruction by embedding visual converters into the generative network. The multi-layer architecture increases image resolution progressively. Up-sampled feature maps are fed into a cross-attention transformer module within each layer and generated images are masked with the same sampling pattern as in the undersampled acquisition. Besides, the parameters of the model are optimized for consistency. The RSTUnet-CR discussed in Table 2 is also capable of restructuring images through a consistency regularization decoder. Besides reconstruction, denoising is also widely investigated. Wang et al. [433] proposed a symmetric TED-Net network, consisting of an encoder-decoder structure and both the encoder and decoder contain several transformer blocks. The input to the encoder is tokenized and the decoder outputs the detokenized result. A transformer block is employed to link the encoder and decoder and the input is removed from the output to calculate the final result. The denoising results of TED-Net and other methods can be found in Figure 14, and it is easy to find that the TED-Net is capable of keeping high-level smoothness and details when removing the artifact or noise, while other methods left more blotchy noise. Yang and co-workers [435] developed a SIST method, in which denoising is performed in the sinogram and image domains using the internal structure in the sinogram domain. In detail, the CT imaging mechanism and statistical characteristics of sinogram are studied for inner-structure loss design to restore high-quality CT images. A sinogram transformer module is also proposed, in which interrelations between projections of different view angles are exploited for sinogram denoising. Moreover, an image reconstruction module is developed to denoise complementarily in both the sinogram and the image domain. The Eformer [437] developed by Luthra et al. is built based on transformer blocks with non-overlapping window-based self-attention. The learnable Sobel-Feldman operators are incorporated to enhance edges and concatenated in the intermediate layers.

**Localization**. Viti et al. [438] developed a transformer graph network, which exploited the self-attention mechanism of the spatial transformer to embed the contextual features of the coronary tree. Specifically, the local features are extracted by CNN and the positional encodings are then embedded into the extracted features. Positional encodings are locally calculated by utilizing the directed tree structure. A

simple signed hop count from the center node is utilized, that is, +1 for distal and -1 for proximal. The resulting features are merged within the self-attention block of the spatial transformer. Besides, 150 coronary CT angiography scans are collected retrospectively.

**Synthesis**. The ResViT model proposed by Dalmaz et al. [440] is composed of an encoder, an information bottleneck, and a decoder. The generator of the ResViT utilizes a central bottleneck with aggregated residual transformer (ART) blocks with transformer modules. The ART block is composed of three parts, which are the transformer encoder-based part, the channel compression part, as well as the residual CNN part. For the given input feature maps, it first passes the transformer encoder-based part, in which the residual connection is implemented. The concatenated features are then fed to the channel compression part with two CNN branches followed by the sum operation. Output feature maps are obtained by feeding the output of the channel compression part to the residual CNN part. As the name implies, the residual connection is also implemented in this block.

### 4.4. Quantitative Evaluation

Table 6 shows the performance of the representative transformer-based models and the performance comparison with other state-of-the-art methods. For classification [29, 115, 121, 162, 165, 173], optimal accuracy, F1 score, area under the curve, precision, recall, balanced accuracy, and Matthew's correlation coefficient are observed. For accuracy, we observe an accuracy of up to 99.6% with the leadership of up to 5.6%. Regarding the F1 score, the highest F1 score of 99.5% and the highest leadership of 3.1% are reached. As for area under the curve, precision, recall, balanced accuracy, and Matthew's correlation coefficient, the highest results are 99.4%, 99.5%, 98.8%, 99.4%, and 98.9%, respectively. For segmentation [194, 269, 273, 274, 328, 329], superior dice similarity coefficient, sensitivity, specificity, pixel accuracy, and intersection over union are achieved. The highest dice similarity coefficient is 90.6% and the highest improvement reaches 4.4%. As for sensitivity, the highest values and leadership are 94.8% and 9.6%, respectively. Regarding the specificity, pixel accuracy, and intersection over union, the best results are 97.7%, 95.9%, and 80.7% respectively. For captioning [399], superior bilingual evaluation understudy-4, consensus-based image description evaluation, and recall-oriented understudy for gisting evaluation-longest common subsequence are reached, equaling 0.58, 0.69, and 0.75, respectively. As for registration [34] and enhancement [433, 435], better dice similarity coefficient, structural similarity index, root mean squared error, and peak signal-to-noise ratio of 74.3%, 0.92, 8.77, and 41.80 are observed.

As can be seen, the reviewed transformer-based method outperforms most existing methods on different MIA tasks. These methods include both the CNN-based methods [53, 183, 448–452] and the transformer-based methods [49, 193, 447, 453–455]. Overperforming CNN-based models can prove the inherent advantage of the transformer-based method while outperforming existing transform-based methods il-

35

lustrates the effectiveness of their proposed improvements. The outstanding results prove the versatility and adaptability of the transformer-based method in the field of MIA. Though the transformer-based method mostly outperforms existing methods to a large extent, we need to point out that there exist some cases that it is not well-at. For example, the transformer-based method can sometimes get a relatively low specificity. To sum up, with overall satisfactory performance, the development of MIA can be significantly boosted by the transformer-based method.

## 5. Challenges and Perspectives

Despite the significant progress and successful deployment of transformer-based methods as a major game changer in the CV area of MIA, future challenges still exist. We summarize several main challenges and give corresponding perspectives on how to solve or improve them. We organize the main challenges together with the corresponding perspectives into three parts, which are feature integration and computing cost reduction, data augmentation and dataset collection, and learning manner and modality-object distribution.

**Feature integration and computing cost reduction**. In order to improve the model performance by capturing both local and global features, most current works only simply hybridize CNN and transformer, such as inserting the transformer encoder block into a CNN model. However, the integration of local and global features in this way may not be firm enough. To integrate CNN and the transformer closer, two-fold ways can be implemented by benefiting the transformer from inductive biases, which are inherent in CNN. On the one hand, inductive bias in CNN can be brought back to the transformer [457–459]. On the other hand, the transformer can learn with CNN simultaneously under the mutual learning framework [460]. High computing cost is always an inevitable problem for the transformer due to the quadratic computational complexity of the input size, especially when the image resolution is high. However, seldom works mentioned or try to solve this problem. To improve the training efficiency of the transformer, more attention computing methods, such as shifted window attention [49], efficient attention [461], and multi-head linear self-attention [462] can be taken into consideration. Besides, projection parameters in the transformer can be shared at different levels. The FLOPs and the number of parameters of the model can be calculated for quantitative model complexity evaluation and further comparison.

**Data augmentation and dataset collection**. In the field of MIA, data shortage always hamper the model performance. The data augmentation technique is an important research direction to address this problem. However, as far as we have seen, many of the transformer-related works have not gone deep into it. Most of the works only use traditional data augmentation techniques such as rotation, crop, and flip. So far, only seldom works utilize advanced data augmentation methods, such as the GAN-based method to synthesize images. Even though, the implemented basic GAN cannot be considered advanced as the quality and resolution of the images

Table 6: Quantitative performance of the representative transformer-based method. For classification, ACC, F1, AUC, PRE, REC, BA, and MCC represent the accuracy, F1 score, area under the curve, precision, recall, balanced accuracy, and Matthew's correlation coefficient, respectively. For segmentation, DSC, SEN, SPE, PA, and IoU stand for dice similarity coefficient, sensitivity, specificity, pixel accuracy, and intersection over union, respectively. For miscellaneous tasks, BLEU-4, CIDEr, ROUGE-L, DSC, SSIM, and RMSE, PSNR mean bilingual evaluation understudy-4, consensus-based image description evaluation, recall-oriented understudy for gisting evaluation-longest common subsequence, dice similarity coefficient, structural similarity index, root mean squared error, and peak signal-to-noise ratio, respectively. We select one of the representative datasets when multiple datasets are used.

| Task | Method | Baseline | Performance (Baseline) |
| --- | --- | --- | --- |
| Classification | transformer [29] | DenseNet | ACC: 97.6% (92.0%), F1: 94.6% (91.5%), PRE: 95.3% (91.0%), REC: 93.8% (92.2%) |
| Classification | transformer [162] | ResNet | ACC: 81.8% (73.1%) |
| Classification | multi-scale feature fusion transformer [165] | DenseNet | ACC: 85.3% (84.3%), F1: 74.2% (72.7%), AUC: 92.3% (90.7%), PRE: 80.2% (80.3%), REC: 70.5% (68.7%) |
| Classification | ensembled swin transformer [115] | swin transformer | ACC: 99.6% (99.2%), F1: 99.5% (99.2%), AUC: 99.4% (99.2%), BA: 99.4% (99.1%), MCC: 98.9% (98.3%) |
| Classification | Self-ViT-MIL [121] | DSMIL [448] | ACC: 91.5% (91.5%), AUC: 94.3% (93.6%) |
| Classification | symmetric dual transformer [173] | EfficientNet-based [449] | ACC: 99.1% (99.0%), F1: 99.1% (99.0%), PRE: 99.5% (99.2%), REC: 98.8% (98.8%) |
| Segmentation | versatile transformer [328] | FrCN [450] | DSC: 85.3% (82.0%), SEN: 83.9% (80.8%), SPE: 82.4% (83.5%), PA: 88.1% (85.3%), IoU: 80.7% (77.4%) |
| Segmentation | semi-supervised transformer [329] | MSA-UNet [456] | DSC: 90.6% (90.3%), SEN: 94.8% (88.7%), SPE: 97.7% (97.1%), PA: 95.9% (95.8%) |
| Segmentation | UTNet [194] | CBAM [453] | DSC: 88.3% (87.3%) |
| Segmentation | UCATR [269] | TransUNet | DSC: 73.6% (70.6%), SEN: 73.1% (69.4%) |
| Segmentation | CCAT-net [273] | FPN [451] | DSC: 65.1% (60.7%), SEN: 76.0% (66.4%), SPE: 97.7% (95.5%) |
| Segmentation | CAC-EMVT [274] | TransUNet | DSC: 75.4% (73.2%), PA: 94.0% (92.3%), IoU: 80.6% (78.0%) |
| Captioning | KdTNet [399] | PPKED [454] | BLEU-4: 0.58 (0.58), CIDEr: 0.69 (0.68), ROUGE-L: 0.75 (0.74) |
| Registration | TD-Net [34] | SYMNet [455] | DSC: 74.3% (73.7%) |
| Enhancement | TED-Net [433] | AD-Net | SSIM: 0.91 (0.90), RMSE: 8.77 (9.72) |
| Enhancement | SIST [435] | DP-ResNet [452] | SSIM: 0.92 (0.91), PSNR: 41.80 (40.92) |

synthesized by basic GAN are difficult to guarantee. In the case of using low-quality or even repetitive (e.g., model collapse) synthesized images for training, the validity of the model performance is questionable. For instance, a classification model can show very high accuracy on a dataset, but there may exist thousands of repetitive synthesized images that are correctly classified in the dataset. To augment data better, state-of-the-art image synthesis models should be taken into consideration. For instance, GAN that suits small datasets like StyleGAN2-ADA [463], independent spatial and appearance transform models [59], and diffusion probabilistic models like 3D-DDPM [60]. Another problem we observed is that many selected papers only compare the model performance with several classic models, and models designed for MIA by other authors are not included. This is especially common for non-mainstream modalities and objects. One of the main reasons that caused it is the lack of widely accepted benching marking datasets like ImageNet [464]. Thus, the collection and publication of new high-quality medical datasets can benefit this research field a lot. Constructing such datasets can also benefit the development of the transfer learning technique in the MIA field. According to our observation, though transfer learning is widely implemented in the field of MIA, most of them transfer from ImageNet. As natural images and medical images can have different data contributions, transferring from medical datasets may further improve model performance.

**Learning manner and modality-object distribution**. There are several state-of-the-art learning manners, such as weakly-supervised learning, and unsupervised learning, which can reduce the need for data labeling. However, these manners are not widely used in transformer-based MIA works. Regarding the modality-object distribution, most existing works mainly concentrate on several mainstream modalities, as shown in Figure 8. However, there is a lot of untapped research potential outside of these mainstream modalities and objects. In terms of modalities, current research primarily focuses on MRI, CT, X-ray, and microscope imaging. Despite being an essential medical image modality, the US has not been fully investigated. In terms of objects, most of the current works focus on the brain, chest, abdomen, and heart, while other objects such as the retina warrant further investigation.

We believe that future efforts in the transformer-based MIA community are certainly not limited to the three points listed above. More research directions like model interpretability should also be fully investigated. With the joining of more and more artificial intelligence and medical researchers, the transformer-based MIA will be developed at an unprecedented speed from both algorithm and data sides. Besides self-benefiting, the fast development of the transformer-based MIA can also benefit related application domains of MIA a lot. For example, state-of-the-art methods [465–467] in these related application domains such as the optimization algorithm can be combined with the transformer. Specifically, the optimization algorithm can be implemented to search the hyperparameter combination in the transformer model to search for further performance improvement. With the fast development of the transformer-based method, the development of MIA can definitely be accelerated to a large extent. This can help doctors to diagnose more fastly, accurately, and smartly,

promoting early intervention.

## 6. Conclusion

The transformer-based MIA is now developing rapidly. In this review, we summarize and analyze the recent progress on transformer-based MIA. The structure of this review is on the basis of different tasks, including classification, segmentation, captioning, registration, detection, enhancement, localization, and synthesis. The task-modality mode can make the readers access their needs faster and easier. We also compare the performance between the transformer-based method and existing state-of-the-art methods. Furthermore, we point out the current challenges and perspectives in the transformer-based MIA field in three core points from both data and algorithm sides. The main advantages of our task-modality review include updated content, detailed information, and comprehensive comparison. Our systematic review may help new DL researchers as well as medical experts without DL knowledge enter this field more quickly. In other words, the detailed content about the latest progress and the performance comparison can be easily accessed with the task-modality mode. However, it is worth noting that the task-modality organization mode may occasionally overlook the sequential relationship between research works. Specifically, subsequent studies building on prior works may be categorized into different tasks or modalities due to the tasks performed or datasets used. While this is not likely to happen frequently, it can occur in certain cases. We show that future work of the transformer-based MIA can be five-fold. First, the method exploration for feature integration and computing cost reduction can be developed. Second, more effort on data augmentation as well as dataset collection should be paid. Next, more focus on the non-mainstream learning manner, modality, and object is needed. Then, deeper research on the model interpretability can be performed. Finally, the transformer-based method and other related application domain methods can be combined. To sum up, benefiting from the fast development of the transformer-based MIA, the medical diagnosis might become more and more convenient and accurate.

## References

[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Adv Neural Inf Process Syst*, 30:1–11, 2017.

[2] Elozino Egonmwan and Yllias Chali. Transformer and seq2seq model for paraphrase generation. In *Proceedings of the Workshop on Neural Generation and Translation*, pages 249–255, 2019.

[3] Li-Wei Chen and Alexander Rudnicky. Fine-grained style control in transformer-based text-to-speech synthesis. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 7907–7911. IEEE, 2022.

[4] Yangyang Shi, Yongqiang Wang, Chunyang Wu, Ching-Feng Yeh, Julian Chan, Frank Zhang, Duc Le, and Mike Seltzer. Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6783–6787. IEEE, 2021.

[5] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv preprint*, 2017. `https://doi.org/10.48550/arXiv.1703.03130`.

[6] Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. *arXiv preprint*, 2016. `https://doi.org/10.48550/arXiv.1606.01933`.

[7] Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for abstractive summarization. *arXiv preprint*, 2017. `https://doi.org/10.48550/arXiv.1705.04304`.

[8] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *arXiv preprint*, 2016. `https://doi.org/10.48550/arXiv.1601.06733`.

[9] Yousry AbdulAzeem, Waleed M Bahgat, and Mahmoud Badawy. A cnn based framework for classification of alzheimer's disease. *Neural Comput. Appl.*, 33(16):10415–10428, 2021. `https://doi.org/10.1007/s00521-021-05799-w`.

[10] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998. `https://doi.org/10.1109/5.726791`.

[11] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Proceedings of the European conference on computer vision*, pages 491–507. Springer, 2020.

[12] Kai Xu, Longyin Wen, Guorong Li, Liefeng Bo, and Qingming Huang. Spatiotemporal cnn for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1379–1388, 2019.

[13] Yang Lei, Xiuxiu He, Jincao Yao, Tonghe Wang, Lijing Wang, Wei Li, Walter J Curran, Tian Liu, Dong Xu, and Xiaofeng Yang. Breast tumor segmentation in 3d automatic breast ultrasound using mask scoring r-cnn. *Med. Phys.*, 48(1):204–214, 2021. `https://doi.org/10.1002/mp.14569`.

[14] Xingxing Xie, Gong Cheng, Jiabao Wang, Xiwen Yao, and Junwei Han. Oriented r-cnn for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3520–3529, 2021.

[15] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages

14454–14463, 2021.

[16] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.

[17] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proceedings of the European conference on computer vision*, pages 213–229. Springer, 2020.

[18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2010.11929`.

[19] Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 357–366, 2021.

[20] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021.

[21] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2906–2917, 2021.

[22] Zhaoshan Liu and Lei Shen. Cect: Controllable ensemble cnn and transformer for covid-19 image classification by capturing both local and global image features. *arXiv preprint*, 2023. `https://doi.org/10.48550/arXiv.2302.02314`.

[23] Ritam Saha and Mrinal Kanti Bhowmik. Active contour model for medical applications. In *Handbook of Research on Natural Computing for Optimization Problems*, pages 937–959. IGI Global, 2016.

[24] Ziduo Yang, Lu Zhao, Shuyu Wu, and Calvin Yu-Chian Chen. Lung lesion localization of covid-19 from chest ct image: A novel weakly supervised learning method. *IEEE J. Biomed. Health Inform.*, 25(6):1864–1872, 2021. `https://doi.org/10.1109/JBHI.2021.3067465`.

[25] S Poonkodi and M Kanchana. 3d-medtrancsgan: 3d medical image transformation using csgan. *Comput. Biol. Med.*, page 106541, 2023. `https://doi.org/10.1016/j.compbiomed.2023.106541`.

[26] Jialei Chen, Chong Fu, Haoyu Xie, Xu Zheng, Rong Geng, and Chiu-Wing Sham. Uncertainty teacher with dense focal loss for semi-supervised medical image segmentation. *Comput. Biol. Med.*, 149:106034, 2022. `https://doi.org/10.1016/j.compbiomed.2022.106034`.

[27] Jun Li, Junyu Chen, Yucheng Tang, Ce Wang, Bennett A Landman, and S Kevin Zhou. Transforming medical imaging with transformers? a comparative review of key properties, current progresses, and future perspectives.

*Med. Image Anal.*, page 102762, 2023. `https://doi.org/10.1016/j.media.2023.102762`.

[28] Tuan Le Dinh, Suk-Hwan Lee, Seong-Geun Kwon, and Ki-Ryong Kwon. Covid-19 chest x-ray classification and severity assessment using convolutional and transformer neural networks. *Appl. Sci.*, 12(10):4861, 2022. `https://doi.org/10.3390/app12104861`.

[29] Koushik Sivarama Krishnan and Karthik Sivarama Krishnan. Vision transformer based covid-19 detection using chest x-rays. In *Proceedings of the International Conference on Signal Processing, Computing and Control*, pages 644–648. IEEE, 2021. `https://doi.org/10.1109/ISPCC53510.2021.9609375`.

[30] Yanan Wu, Shouliang Qi, Yu Sun, Shuyue Xia, Yudong Yao, and Wei Qian. A vision transformer for emphysema classification using ct images. *Phys. Med. Biol.*, 66(24):245016, 2021. `https://doi.org/10.1088/1361-6560/ac3dc8`.

[31] Hong Gu, Hongyu Wang, Pan Qin, and Jia Wang. Chest l-transformer: Local features with position attention for weakly supervised chest radiograph segmentation and classification. *Front. Med.*, page 1619, 2022. `https://doi.org/10.3389/fmed.2022.923456`.

[32] Linh T Duong, Nhi H Le, Toan B Tran, Vuong M Ngo, and Phuong T Nguyen. Detection of tuberculosis from chest x-ray images: boosting the performance with vision transformer and transfer learning. *Expert Syst. Appl.*, 184:115519, 2021. `https://doi.org/10.1016/j.eswa.2021.115519`.

[33] Zheng Jiang and Liang Chen. Multisemantic level patch merger vision transformer for diagnosis of pneumonia. *Comput. Math. Method Med.*, 2022, 2022. `https://doi.org/10.1155/2022/7852958`.

[34] Lei Song, Guixia Liu, and Mingrui Ma. Td-net: unsupervised medical image registration network based on transformer and cnn. *Appl. Intell.*, pages 1–9, 2022. `https://doi.org/10.1007/s10489-022-03472-w`.

[35] Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data Brief*, 28:104863, 2020. `https://doi.org/10.1016/j.dib.2019.104863`.

[36] Zebin Hu, Hao Liu, Zhendong Li, and Zekuan Yu. Cross-model transformer method for medical image synthesis. *Complexity*, 2021:1–7, 2021. `https://doi.org/10.1155/2021/5624909`.

[37] Sergey P Morozov, AE Andreychenko, NA Pavlov, AV Vladzymyrskyy, NV Ledikhova, VA Gombolevskiy, Ivan A Blokhin, PB Gelezhe, AV Gonchar, and V Yu Chernina. Mosmeddata: Chest ct scans with covid-19 related findings dataset. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2005.06465`.

[38] Jakob Nikolas Kather, Cleo-Aron Weis, Francesco Bianconi, Susanne M Melchers, Lothar R Schad, Timo Gaiser, Alexander Marx, and Frank Gerrit Zöllner. Multi-class texture analysis in colorectal cancer histology. *Sci Rep*, 6(1):1–11, 2016. `https://doi.org/10.1038/srep27988`.

[39] Md Robiul Islam, Md Nahiduzzaman, Md Omaer Faruq Goni, Abu Sayeed, Md Shamim Anower, Mominul Ahsan, and Julfikar Haider. Explainable transformer-based deep learning model for the detection of malaria parasites from blood cell images. *Sensors*, 22(12):4358, 2022. `https://doi.org/10.3390/s22124358`.

[40] Juan Silva, Aymeric Histace, Olivier Romain, Xavier Dray, and Bertrand Granado. Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *Int. J. Comput. Assist. Radiol. Surg.*, 9(2):283–293, 2014. `https://doi.org/10.1007/s11548-013-0926-3`.

[41] David Ouyang, Bryan He, Amirata Ghorbani, Neal Yuan, Joseph Ebinger, Curtis P Langlotz, Paul A Heidenreich, Robert A Harrington, David H Liang, Euan A Ashley, et al. Video-based ai for beat-to-beat assessment of cardiac function. *Nature*, 580(7802):252–256, 2020. `https://doi.org/10.1038/s41586-020-2145-8`.

[42] Suliman Aladhadh, Majed Alsanea, Mohammed Aloraini, Taimoor Khan, Shabana Habib, and Muhammad Islam. An effective skin cancer classification mechanism via medical vision transformer. *Sensors*, 22(11):4008, 2022. `https://doi.org/10.3390/s22114008`.

[43] Danny Chen, Wenzhong Yang, Liejun Wang, Sixiang Tan, Jiangzhaung Lin, and Wenxiu Bu. Pcat-unet: Unet-like network fused convolution and transformer for retinal vessel segmentation. *PLoS One*, 17(1):e0262689, 2022. `https://doi.org/10.1371/journal.pone.0262689`.

[44] Abdul Qayyum, Abdesslam Benzinou, Moona Mazher, and Fabrice Meriaudeau. Efficient multi-model vision transformer based on feature fusion for classification of dfuc2021 challenge. In *Proceedings of the Diabetic Foot Ulcers Grand Challenge*, pages 62–75. Springer, 2021. `https://doi.org/10.1007/978-3-030-94907-5_5`.

[45] Pinxian Zeng, Luping Zhou, Chen Zu, Xinyi Zeng, Zhengyang Jiao, Xi Wu, Jiliu Zhou, Dinggang Shen, and Yan Wang. 3d cvt-gan: A 3d convolutional vision transformer-gan for pet reconstruction. In *Medical Image Computing and Computer Assisted Intervention*, pages 516–526. Springer, 2022. `https://doi.org/10.1007/978-3-031-16446-0_49`.

[46] Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell*, 172(5):1122–1131, 2018. `https://doi.org/10.1016/j.cell.2018.02.010`.

[47] Kelei He, Chen Gan, Zhuoyuan Li, Islem Rekik, Zihao Yin, Wen Ji, Yang Gao, Qian Wang, Junfeng Zhang, and Dinggang Shen. Transformers in medical image analysis: A review. *arXiv preprint*, 2022. `https://doi.org/10.48550/arXiv.2202.12165`.

[48] Arshi Parvaiz, Muhammad Anwaar Khalid, Rukhsana Zafar, Huma Ameer, Muhammad Ali, and Muhammad Moazam Fraz. Vision transformers in med-

ical computer vision–a contemplative retrospection. *arXiv preprint*, 2022. `https://doi.org/10.48550/arXiv.2203.15269`.

[49] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the International Conference on Computer Vision*, page 1001210022, 2021.

[50] Tao Wang, Junlin Lan, Zixin Han, Ziwei Hu, Yuxiu Huang, Yanglin Deng, Hejun Zhang, Jianchao Wang, Musheng Chen, Haiyan Jiang, et al. O-net: A novel framework with deep fusion of cnn and transformer for simultaneous segmentation and classification. *Front. Neurosci.*, 16, 2022. `https://doi.org/10.3389/fnins.2022.876065`.

[51] Shaolong Chen, Lijie Zhong, Changzhen Qiu, Zhiyong Zhang, and Xiaodong Zhang. Transformer-based multilevel region and edge aggregation network for magnetic resonance image segmentation. *Comput. Biol. Med.*, 152:106427, 2023. `https://doi.org/10.1016/j.compbiomed.2022.106427`.

[52] Flow diagram. `http://www.prisma-statement.org/PRISMAStatement/FlowDiagram`. Accessed 9 May 2023.

[53] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the conference on computer vision and pattern recognition*, pages 770–778, 2016.

[54] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint*, 2018. `https://doi.org/10.48550/arXiv.1810.04805`.

[55] Saidi Guo, Lin Xu, Cheng Feng, Huahua Xiong, Zhifan Gao, and Heye Zhang. Multi-level semantic adaptation for few-shot segmentation on cardiac image sequences. *Med. Image Anal.*, 73:102170, 2021. `https://doi.org/10.1016/j.media.2021.102170`.

[56] Hao Tang, Xingwei Liu, Shanlin Sun, Xiangyi Yan, and Xiaohui Xie. Recurrent mask refinement for few-shot medical image segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3918–3928, 2021.

[57] Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Fahmy Aly. Deep learning approaches for data augmentation and classification of breast masses using ultrasound images. *Int. J. Adv. Comput. Sci. Appl*, 10(5):1–11, 2019.

[58] Zhaoshan Liu, Chau Hung Lee, and Lei Shen. Semi-supervised classification of medical ultrasound images based on generative adversarial network. *arXiv preprint*, 2022. `https://doi.org/10.48550/arXiv.2203.06184`.

[59] Amy Zhao, Guha Balakrishnan, Fredo Durand, John V Guttag, and Adrian V Dalca. Data augmentation using learned transformations for one-shot medical image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8543–8553, 2019.

[60] Zolnamar Dorjsembe, Sodtavilan Odonchimed, and Furen Xiao. Three-

dimensional medical image synthesis with denoising diffusion probabilistic models. In *Proceedings of the Medical Imaging with Deep Learning*, pages 1–3, 2022.

[61] Sihong Chen, Kai Ma, and Yefeng Zheng. Med3d: Transfer learning for 3d medical image analysis. *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1904.00625`.

[62] Jianlong Zhou, Zelin Li, Weiming Zhi, Bin Liang, Daniel Moses, and Laughlin Dawes. Using convolutional neural networks and transfer learning for bone age classification. In *2017 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–6. IEEE, 2017. `https://doi.org/10.1109/DICTA.2017.8227503`.

[63] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.

[64] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Adv Neural Inf Process Syst*, 33:21271–21284, 2020.

[65] Lisa Gottesfeld Brown. A survey of image registration techniques. *ACM Comput. Surv.*, 24(4):325–376, 1992. `https://doi.org/10.1145/146370.146374`.

[66] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[67] Sukhjinder Singh, RK Bansal, and Savina Bansal. Medical image enhancement using histogram processing techniques followed by median filter. *Ijipa*, 3(1):1–9, 2012.

[68] Liron Pantanowitz, Paul N Valenstein, Andrew J Evans, Keith J Kaplan, John D Pfeifer, David C Wilbur, Laura C Collins, and Terence J Colgan. Review of the current state of whole slide imaging in pathology. *Journal of pathology informatics*, 2(1):36, 2011. `https://doi.org/10.4103/2153-3539.83746`.

[69] Chiagoziem C Ukwuoma, Zhiguang Qin, Md Belal Bin Heyat, Faijan Akhtar, Abla Smahi, Jehoiada K Jackson, Syed Furqan Qadri, Abdullah Y Muaad, Happy N Monday, and Grace U Nneji. Automated lung-related pneumonia and covid-19 detection based on novel feature extraction framework and vision transformer approaches using chest x-ray images. *Bioengineering*, 9(11):709, 2022. `https://doi.org/10.3390/bioengineering9110709`.

[70] Abeer Badawi and Khalid Elgazzar. Detecting coronavirus from chest x-rays using transfer learning. *Covid*, 1(1):403–415, 2021. `https://doi.org/10.3390/covid1010034`.

[71] Muhammad EH Chowdhury, Tawsifur Rahman, Amith Khandakar, Rashid Mazhar, Muhammad Abdul Kadir, Zaid Bin Mahbub, Khandakar Reajul Islam,

Muhammad Salman Khan, Atif Iqbal, Nasser Al Emadi, et al. Can ai help in screening viral and covid-19 pneumonia? *IEEE Access*, 8:132665–132676, 2020. https://doi.org/10.1109/ACCESS.2020.3010287.

[72] Covidx cxr-2. https://www.kaggle.com/datasets/andyczhao/covidx-cxr2?select=competition_test. Accessed 27 July 2022.

[73] Daniel Kermany, Kang Zhang, Michael Goldbaum, et al. Labeled optical coherence tomography (oct) and chest x-ray images for classification. *Mendeley data*, 2(2), 2018. https://doi.org/10.17632/rscbjbr9sj.2.

[74] Emily B Tsai, Scott Simpson, Matthew P Lungren, Michelle Hershman, Leonid Roshkovan, Errol Colak, Bradley J Erickson, George Shih, Anouk Stein, Jayashree Kalpathy-Cramer, et al. The rsna international covid-19 open radiology database (ricord). *Radiology*, 299(1):E204, 2021. https://doi.org/10.1148/radiol.2021203957.

[75] Joseph Paul Cohen, Beiyi Shen, Almas Abbasi, Mahsa Hoshmand-Kochi, Samantha Glass, Haifang Li, Matthew P Lungren, Akshay Chaudhari, and Tim Q Duong. Radiographic assessment of lung opacity score dataset. *Zenodo*, 4633999, 2021. https://doi.org/10.5281/zenodo.4634000.

[76] Tawsifur Rahman, Amith Khandakar, Yazan Qiblawey, Anas Tahir, Serkan Kiranyaz, Saad Bin Abul Kashem, Mohammad Tariqul Islam, Somaya Al Maadeed, Susu M Zughaier, Muhammad Salman Khan, et al. Exploring the effect of image enhancement techniques on covid-19 detection using chest x-ray images. *Comput. Biol. Med.*, 132:104319, 2021. https://doi.org/10.1016/j.compbiomed.2021.104319.

[77] Finn Behrendt, Debayan Bhattacharya, Julia Krüger, Roland Opfer, and Alexander Schlaefer. Data-efficient vision transformers for multi-label disease classification on chest radiographs. *Current Directions in Biomedical Engineering*, 8(1):34–37, 2022. https://doi.org/10.1515/cdbme-2022-0009.

[78] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019. https://doi.org/10.1609/aaai.v33i01.3301590.

[79] Gelan Ayana, Kokeb Dese, Yisak Dereje, Yonas Kebede, Hika Barki, Dechassa Amdissa, Nahimiya Husen, Fikadu Mulugeta, Bontu Habtamu, and Se-woon Choe. Vision-transformer-based transfer learning for mammogram classification. *Diagnostics*, 13(2):178, 2023. https://doi.org/10.3390/diagnostics13020178.

[80] Dataset of breast mammography images with masses. https://data.mendeley.com/datasets/ywsbh3ndr8/2. Accessed 14 February 2023.

[81] Seyed Ali Jalalifar and Ali Sadeghi-Naini. Data-efficient training of pure vision transformers for the task of chest x-ray abnormality detection using knowledge distillation. In *Proceedings of the Annual International Conference of the IEEE*

*Engineering in Medicine & Biology Society*, pages 1444–1447. IEEE, 2022. https://doi.org/10.1109/EMBC48229.2022.9871372.

[82] Xuxin Chen, Ke Zhang, Neman Abdoli, Patrik W Gilley, Ximin Wang, Hong Liu, Bin Zheng, and Yuchen Qiu. Transformers improve breast cancer diagnosis from unregistered multi-view mammograms. *Diagnostics*, 12(7):1549, 2022. https://doi.org/10.3390/diagnostics12071549.

[83] Bin Zheng, Jules H Sumkin, Margarita L Zuley, Dror Lederman, Xingwei Wang, and David Gur. Computer-aided detection of breast masses depicted on full-field digital mammograms: a performance assessment. *Br. J. Radiol.*, 85(1014):e153–e161, 2012. https://doi.org/10.1259/bjr/51461617.

[84] Stefan Jaeger, Sema Candemir, Sameer Antani, Yì-Xiáng J Wáng, Pu-Xuan Lu, and George Thoma. Two public chest x-ray datasets for computer-aided screening of pulmonary diseases. *Quant. Imaging Med. Surg.*, 4(6):475, 2014. https://doi.org/10.3978/j.issn.2223-4292.2014.11.20.

[85] Joseph Paul Cohen, Paul Morrison, Lan Dao, Karsten Roth, Tim Q Duong, and Marzyeh Ghassemi. Covid-19 image data collection: Prospective predictions are the future. *arXiv preprint*, 2020. https://doi.org/10.48550/arXiv.2006.11988.

[86] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, pages 2097–2106, 2017.

[87] Tianmu Wang, Zhenguo Nie, Ruijing Wang, Qingfeng Xu, Hongshi Huang, Handing Xu, Fugui Xie, and Xin-Jun Liu. Pneunet: deep learning for covid-19 pneumonia diagnosis on chest x-ray image analysis using vision transformer. *Med. Biol. Eng. Comput.*, pages 1–14, 2023. https://doi.org/10.1007/s11517-022-02746-2.

[88] Qata-cov19 dataset. https://www.kaggle.com/datasets/aysendegerli/qatacov19-dataset. Accessed 14 February 2023.

[89] Covid-19 image repository. https://github.com/ml-workgroup/COVID-19-image-repository/tree/master/png. Accessed 14 February 2023.

[90] Eurorad dataset. https://www.eurorad.org/. Accessed 14 February 2023.

[91] Covid chestxray dataset. https://github.com/ieee8023/COVID-chestxray-dataset. Accessed 14 February 2023.

[92] Covid-19 dataset. https://www.sirm.org/category/senza-categoria/COVID-19/. Accessed 14 February 2023.

[93] Covid-19 radiography database. https://www.kaggle.com/datasets/tawsifurrahman/covid19-radiography-database. Accessed 14 February 2023.

[94] Covid-cxnet. https://github.com/armiro/COVID-CXNet. Accessed 14 February 2023.

[95] Rsna pneumonia detection challenge. `https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/data`. Accessed 14 February 2023.

[96] Chest x-ray images (pneumonia). `https://www.kaggle.com/datasets/paultimothymooney/chest-xray-pneumonia`. Accessed 4 August 2022.

[97] Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria de la Iglesia-Vayá. Padchest: A large chest x-ray image dataset with multi-label annotated reports. *Med. Image Anal.*, 66:101797, 2020. `https://doi.org/10.1016/j.media.2020.101797`.

[98] Ross W Filice, Anouk Stein, Carol C Wu, Veronica A Arteaga, Stephen Borstelmann, Ramya Gaddikeri, Maya Galperin-Aizenberg, Ritu R Gill, Myrna C Godoy, Stephen B Hobbs, et al. Crowdsourcing pneumothorax annotations using machine learning annotations on the nih chest x-ray dataset. *J. Digit. Imaging*, 33(2):490–496, 2020. `https://doi.org/10.1007/s10278-019-00299-9`.

[99] Xiaoben Jiang, Yu Zhu, Gan Cai, Bingbing Zheng, and Dawei Yang. Mxt: A new variant of pyramid vision transformer for multi-label chest x-ray image classification. *Cogn. Comput.*, pages 1–16, 2022. `https://doi.org/10.1007/s12559-022-10032-4`.

[100] Chestx-ray14 dataset. `https://nihcc.app.box.com/v/ChestXray-NIHCC`. Accessed 27 July 2022.

[101] Ranzcr clip - catheter and line position challenge. `https://www.kaggle.com/competitions/ranzcr-clip-catheter-line-classification/data`. Accessed 27 July 2022.

[102] Xiao Qi, David J Foran, John L Nosher, and Ilker Hacihaliloglu. Multi-feature vision transformer via self-supervised representation learning for improvement of covid-19 diagnosis. In *Proceedings of the Medical Image Learning with Limited and Noisy Data*, pages 76–85. Springer, 2022. `https://doi.org/10.1007/978-3-031-16760-7_8`.

[103] Maria De La Iglesia Vayá, Jose Manuel Saborit, Joaquim Angel Montell, Antonio Pertusa, Aurelia Bustos, Miguel Cazorla, Joaquin Galant, Xavier Barber, Domingo Orozco-Beltrán, Francisco García-García, et al. Bimcv covid-19+: a large annotated dataset of rx and ct images from covid-19 patients. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2006.01174`.

[104] Linda Wang, Zhong Qiu Lin, and Alexander Wong. Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Sci Rep*, 10(1):1–12, 2020. `https://doi.org/10.1038/s41598-020-76550-z`.

[105] Shivang Desai, Ahmad Baghal, Thidathip Wongsurawat, Piroon Jenjaroenpun, Thomas Powell, Shaymaa Al-Shukri, Kim Gates, Phillip Farmer, Michael Rutherford, Geri Blake, et al. Chest imaging representing a covid-19 positive rural us population. *Sci. Data*, 7(1):414, 2020. `https://doi.org/10.1038/s41597-020-00741-6`.

[106] Kenneth Clark, Bruce Vendt, Kirk Smith, John Freymann, Justin Kirby, Paul Koppel, Stephen Moore, Stanley Phillips, David Maffitt, Michael Pringle, et al.

The cancer imaging archive (tcia): maintaining and operating a public information repository. *J. Digit. Imaging*, 26:1045–1057, 2013. `https://doi.org/10.1007/s10278-013-9622-7`.

[107] Sangjoon Park, Gwanghyun Kim, Jeongsol Kim, Boah Kim, and Jong Chul Ye. Federated split vision transformer for covid-19 cxr diagnosis using task-agnostic training. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2111.01338`.

[108] Alberto Signoroni, Mattia Savardi, Sergio Benini, Nicola Adami, Riccardo Leonardi, Paolo Gibellini, Filippo Vaccher, Marco Ravanelli, Andrea Borghesi, Roberto Maroldi, et al. Bs-net: Learning covid-19 pneumonia severity on a large chest x-ray dataset. *Med. Image Anal.*, 71:102046, 2021. `https://doi.org/10.1016/j.media.2021.102046`.

[109] Andrea Borghesi and Roberto Maroldi. Covid-19 outbreak in italy: experimental chest x-ray scoring system for quantifying and monitoring disease progression. *La radiologia medica*, 125(5):509–513, 2020. `https://doi.org/10.1007/s11547-020-01200-3`.

[110] Kobiljon Ikromjanov, Subrata Bhattacharjee, Yeong-Byn Hwang, Rashadul Islam Sumon, Hee-Cheol Kim, and Heung-Kook Choi. Whole slide image analysis and detection of prostate cancer using vision transformers. In *Proceedings of the International Conference on Artificial Intelligence in Information and Communication*, pages 399–402. IEEE, 2022. `https://doi.org/10.1109/ICAIIC54071.2022.9722635`.

[111] Prostate cancer grade assessment (panda) challenge. `https://www.kaggle.com/c/prostate-cancer-grade-assessment`. Accessed 3 August 2022.

[112] Magdy Abd-Elghany Zeid, Khaled El-Bahnasy, and SE Abo-Youssef. Multi-class colorectal cancer histology images classification using vision transformers. In *Proceedings of the International Conference on Intelligent Computing and Information Systems*, pages 224–230. IEEE, 2021. `https://doi.org/10.1109/ICICIS52592.2021.9694125`.

[113] Sivaramakrishnan Rajaraman, Sameer K Antani, Mahdieh Poostchi, Kamolrat Silamut, Md A Hossain, Richard J Maude, Stefan Jaeger, and George R Thoma. Pre-trained convolutional neural networks as feature extractors toward improved malaria parasite detection in thin blood smear images. *PeerJ*, 6:e4568, 2018. `https://doi.org/10.7717/peerj.4568`.

[114] KM Faizullah Fuhad, Jannat Ferdousey Tuba, Md Rabiul Ali Sarker, Sifat Momen, Nabeel Mohammed, and Tanzilur Rahman. Deep learning based automatic malaria parasite detection from blood smear and its smartphone based application. *Diagnostics*, 10(5):329, 2020. `https://doi.org/10.3390/diagnostics10050329`.

[115] Sudhakar Tummala, Jungeun Kim, and Seifedine Kadry. Breast-net: Multi-class classification of breast cancer from histopathological images using ensemble of swin transformers. *Mathematics*, 10(21):4109, 2022. `https://doi.org/10.3390/math10214109`.

[116] Fabio A Spanhol, Luiz S Oliveira, Caroline Petitjean, and Laurent Heutte. A dataset for breast cancer histopathological image classification. *IEEE Trans. Biomed. Eng.*, 63(7):1455–1462, 2015. `https://doi.org/10.1109/TBME.2015.2496264`.

[117] Panagiotis Barmpoutis, Jing Yuan, William Waddingham, Christopher Ross, Kayhanian Hamzeh, Tania Stathaki, Daniel C Alexander, and Marnix Jansen. Multi-scale deformable transformer for the classification of gastric glands: The imgl dataset. In *Proceedings of the Cancer Prevention Through Early Detection*, pages 24–33. Springer, 2022. `https://doi.org/10.1007/978-3-031-17979-2_3`.

[118] Imgl dataset. `https://zenodo.org/record/6908133#.Y-svhnYzZPZ`. Accessed 14 February 2023.

[119] Jianxin Zhang, Cunqiao Hou, Wen Zhu, Mingli Zhang, Ying Zou, Lizhi Zhang, and Qiang Zhang. Attention multiple instance learning with transformer aggregation for breast cancer whole slide image classification. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1804–1809. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9994848`.

[120] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson, Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama*, 318(22):2199–2210, 2017. `https://doi.org/10.1001/jama.2017.14585`.

[121] Ahmet Gokberk Gul, Oezdemir Cetin, Christoph Reich, Nadine Flinner, Tim Prangemeier, and Heinz Koeppl. Histopathological image classification based on self-supervised vision transformer and weak labels. In *Proceedings of the Medical Imaging 2022: Digital and Computational Pathology*, volume 12039, pages 366–373. SPIE, 2022. `https://doi.org/10.1117/12.2624609`.

[122] Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Junzhou Huang, Wei Yang, and Xiao Han. Transpath: Transformer-based self-supervised learning for histopathological image classification. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 186–195. Springer, 2021. `https://doi.org/10.1007/978-3-030-87237-3_18`.

[123] Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A Valous, Dyke Ferber, et al. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS Med.*, 16(1):e1002730, 2019. `https://doi.org/10.1371/journal.pmed.1002730`.

[124] Jerry Wei, Arief Suriawinata, Bing Ren, Xiaoying Liu, Mikhail Lisovsky, Louis Vaickus, Charles Brown, Michael Baker, Naofumi Tomita, Lorenzo Torresani,

et al. A petri dish for histopathology image analysis. In *Proceedings of the International Conference on Artificial Intelligence in Medicine*, pages 11–24. Springer, 2021. `https://doi.org/10.1007/978-3-030-77211-6_2`.

[125] Hufei Duan, Yiqing Liu, Hui Yan, Qiming He, Yonghong He, and Tian Guan. Fourier vit: A multi-scale vision transformer with fourier transform for histopathological image classification. In *Proceedings of the International Conference on Automation, Control and Robotics Engineering*, pages 189–193. IEEE, 2022. `https://doi.org/10.1109/CACRE54574.2022.9834158`.

[126] Mohamed Amgad, Habiba Elfandy, Hagar Hussein, Lamees A Atteya, Mai AT Elsebaie, Lamia S Abo Elnasr, Rokia A Sakr, Hazem SE Salem, Ahmed F Ismail, Anas M Saad, et al. Structured crowdsourcing enables convolutional segmentation of histology images. *Bioinformatics*, 35(18):3461–3467, 2019. `https://doi.org/10.1093/bioinformatics/btz083`.

[127] Zhilong Lv, Rui Yan, Yuexiao Lin, Ying Wang, and Fa Zhang. Joint region-attention and multi-scale transformer for microsatellite instability detection from whole slide images in gastrointestinal cancer. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 293–302. Springer, 2022. `https://doi.org/10.1007/978-3-031-16434-7_29`.

[128] Yaowei Wang, Jing Guo, Yun Yang, Yan Kang, Yuelong Xia, Zhenhui Li, Yongchun Duan, and Kelong Wang. Cwc-transformer: a visual transformer approach for compressed whole slide image classification. *Neural Comput. Appl.*, pages 1–13, 2023. `https://doi.org/10.1007/s00521-022-07857-3`.

[129] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16519–16529, 2021.

[130] Lulu Gai, Wei Chen, Rui Gao, Yan-Wei Chen, and Xu Qiao. Using vision transformers in 3-d medical image classifications. In *Proceedings of the IEEE International Conference on Image Processing*, pages 696–700. IEEE, 2022. `https://doi.org/10.1109/ICIP46576.2022.9897966`.

[131] Mohammad Rahimzadeh, Abolfazl Attar, and Seyed Mohammad Sakhaei. A fully automated deep learning-based network for detecting covid-19 from a new and large lung ct scan dataset. *Biomed. Signal Process. Control*, 68:102588, 2021. `https://doi.org/10.1016/j.bspc.2021.102588`.

[132] Maisarah Mohd Sufian, Ervin Gubin Moung, Jamal Ahmad Dargham, Farashazillah Yahya, and Sigeru Omatu. Pre-trained deep learning models for covid19 classification: Cnns vs. vision transformer. In *Proceedings of the IEEE International Conference on Artificial Intelligence in Engineering and Technology*, pages 1–6. IEEE, 2022. `https://doi.org/10.1109/IICAIET55139.2022.9936852`.

[133] Eduardo Soares, Plamen Angelov, Sarah Biaso, Michele Higa Froes, and Daniel Kanda Abe. Sars-cov-2 ct-scan dataset: A large dataset of real patients ct scans for sars-cov-2 identification. *MedRxiv*, 2020. `https://doi.`

org/10.1101/2020.04.24.20078584.

[134] Jingxing Li, Zhanglei Yang, and Yifan Yu. A medical ai diagnosis platform based on vision transformer for coronavirus. In *Proceedings of the IEEE International Conference on Computer Science, Electronic Information Engineering and Intelligent Control Technology*, pages 246–252. IEEE, 2021. https://doi.org/10.1109/CEI52496.2021.9574576.

[135] Anish Salvi, Raj Shah, Luke Higgins, and Prahlad G Menon. Vision transformers for ai-driven classification of peripheral artery disease from maximum intensity projections of runoff ct angiograms. In *Proceedings of the International Conference on Bioinformatics and Biomedicine*, pages 3870–3872. IEEE, 2022. https://doi.org/10.1109/BIBM55620.2022.9995337.

[136] Pranab Sahoo, Sriparna Saha, Samrat Mondal, and Suraj Gowda. Vision transformer based covid-19 detection using chest ct-scan images. In *Proceedings of the IEEE-EMBS International Conference on Biomedical and Health Informatics*, pages 01–04. IEEE, 2022. https://doi.org/10.1109/BHI56158.2022.9926823.

[137] Lauge Sorensen, Saher B Shaker, and Marleen De Bruijne. Computed tomography emphysema database. https://lauge-soerensen.github.io/emphysema-database/. Accessed 27 July 2022.

[138] Lauge Sorensen, Saher B Shaker, and Marleen De Bruijne. Quantitative analysis of pulmonary emphysema using local binary patterns. *IEEE Trans. Med. Imaging*, 29(2):559–569, 2010. https://doi.org/10.1109/TMI.2009.2038575.

[139] Yuxuan Xiong, Bo Du, Yongchao Xu, Jiajun Deng, Yunlang She, and Chang Chen. Pulmonary nodule classification with multi-view convolutional vision transformer. In *Proceedings of the International Joint Conference on Neural Networks*, pages 1–7. IEEE, 2022. https://doi.org/10.1109/IJCNN55064.2022.9892716.

[140] Jie Mei. Marrying convolution and transformer for covid-19 diagnosis based on ct scans. In *Proceedings of the International Joint Conference on Neural Networks*, pages 1–7. IEEE, 2022. https://doi.org/10.1109/IJCNN55064.2022.9892015.

[141] Jinyu Zhao, Yichen Zhang, Xuehai He, and Pengtao Xie. Covid-ct-dataset: a ct scan dataset about covid-19. *arXiv preprint*, 2020. https://doi.org/10.48550/arXiv.2003.13865.

[142] Yingying Huang, Yang Si, Bingliang Hu, Yan Zhang, Shuang Wu, Dongsheng Wu, and Quan Wang. Transformer-based factorized encoder for classification of pneumoconiosis on 3d ct images. *Comput. Biol. Med.*, 150:106137, 2022. https://doi.org/10.1016/j.compbiomed.2022.106137.

[143] Parnian Afshar, Shahin Heidarian, Nastaran Enshaei, Farnoosh Naderkhani, Moezedin Javad Rafiee, Anastasia Oikonomou, Faranak Babaki Fard, Kaveh Samimi, Konstantinos N Plataniotis, and Arash Mohammadi. Covid-ct-md, covid-19 computed tomography scan dataset applicable in machine learning

and deep learning. *Sci. Data*, 8(1):121, 2021. `https://doi.org/10.1038/s41597-021-00900-3`.

[144] Kunlun Wu, Bo Peng, and Donghai Zhai. Multi-granularity dilated transformer for lung nodule classification via local focus scheme. *Appl. Sci.*, 13(1):377, 2022. `https://doi.org/10.3390/app13010377`.

[145] Samuel G Armato III, Geoffrey McLennan, Luc Bidaut, Michael F McNitt-Gray, Charles R Meyer, Anthony P Reeves, Binsheng Zhao, Denise R Aberle, Claudia I Henschke, Eric A Hoffman, et al. The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans. *Med. Phys.*, 38(2):915–931, 2011. `https://doi.org/10.1118/1.3528204`.

[146] F Proietto Salanitri, Giovanni Bellitto, Simone Palazzo, Ismail Irmakci, M Wallace, C Bolan, Megan Engels, Sanne Hoogenboom, Marco Aldinucci, Ulas Bagci, et al. Neural transformers for intraductal papillary mucosal neoplasms (ipmn) classification in mri images. In *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine & Biology Society*, pages 475–479. IEEE, 2022. `https://doi.org/10.1109/EMBC48229.2022.9871547`.

[147] Rodney LaLonde, Irene Tanner, Katerina Nikiforaki, Georgios Z Papadakis, Pujan Kandel, Candice W Bolan, Michael B Wallace, and Ulas Bagci. Inn: inflated neural networks for ipmn diagnosis. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 101–109. Springer, 2019. `https://doi.org/10.1007/978-3-030-32254-0_12`.

[148] Yin Dai, Yifan Gao, and Fayu Liu. Transmed: Transformers advance multimodal medical image classification. *Diagnostics*, 11(8):1384, 2021. `https://doi.org/10.3390/diagnostics11081384`.

[149] Nicholas Bien, Pranav Rajpurkar, Robyn L Ball, Jeremy Irvin, Allison Park, Erik Jones, Michael Bereket, Bhavik N Patel, Kristen W Yeom, Katie Shpanskaya, et al. Deep-learning-assisted diagnosis for knee magnetic resonance imaging: development and retrospective validation of mrnet. *PLos Med.*, 15(11):e1002699, 2018. `https://doi.org/10.1371/journal.pmed.1002699`.

[150] Ixi dataset. `https://brain-development.org/ixi-dataset/`. Accessed 30 July 2022.

[151] Jin Liu, Hao Du, Qian Bi, Haiyan Liao, and Yi Pan. Mest: Multi-plane embedding and spatial-temporal transformer for parkinson's disease diagnosis. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1072–1077. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9995498`.

[152] Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey, Karl Kieburtz, Emily Flagg, Sohini Chowdhury, et al. The parkinson progression marker initiative (ppmi). *Prog. Neurobiol.*, 95(4):629–635, 2011. `https://doi.org/10.1016/j.pneurobio.2011.09.005`.

[153] Shuang Yu, Kai Ma, Qi Bi, Cheng Bian, Munan Ning, Nanjun He, Yuexiang

Li, Hanruo Liu, and Yefeng Zheng. Mil-vt: Multiple instance learning enhanced vision transformer for fundus image classification. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 45–54. Springer, 2021. `https://doi.org/10.1007/978-3-030-87237-3_5`.

[154] Aptos 2019 blindness detection. `https://www.kaggle.com/c/aptos2019-blindness-detection`. Accessed 27 July 2022.

[155] Retinal image analysis for multi-disease detection challenge. `https://riadd.grand-challenge.org/`. Accessed 27 July 2022.

[156] Sharif Amit Kamran, Khondker Fariha Hossain, Alireza Tavakkoli, Stewart Lee Zuckerbrod, and Salah A Baker. Vtgan: Semi-supervised retinal image synthesis and disease prediction using vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3235–3245, 2021.

[157] Shirin Hajeb Mohammad Alipour, Hossein Rabbani, and Mohammad Reza Akhlaghi. Diabetic retinopathy grading by digital curvelet transform. *Comput. Math. Method Med.*, 2012, 2012.

[158] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Sci. Data*, 5(1):1–9, 2018. `https://doi.org/10.1038/sdata.2018.161`.

[159] Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In *Proceedings of the International Symposium on Biomedical Imaging*, pages 168–172. IEEE, 2018.

[160] AKM Salman Hosain, Mynul Islam, Md Humaion Kabir Mehedi, Irteza Enan Kabir, and Zarin Tasnim Khan. Gastrointestinal disorder detection with a transformer based approach. In *Proceedings of the Annual Information Technology, Electronics and Mobile Communication Conference*, pages 0280–0285. IEEE, 2022. `https://doi.org/10.1109/IEMCON56893.2022.9946531`.

[161] Gastrointestinal tract or colon diseases image dataset. `https://www.kaggle.com/datasets/francismon/curated-colon-dataset-for-deep-learning`. Accessed 10 February 2023.

[162] Aniruddha Tamhane, Tse'ela Mida, Erez Posner, and Moshe Bouhnik. Colonoscopy landmark detection using vision transformers. In *Proceedings of the Imaging Systems for GI Endoscopy, and Graphs in Biomedical Image Analysis*, pages 24–34. Springer, 2022. `https://doi.org/10.1007/978-3-031-21083-9_3`.

[163] Behnaz Gheflati and Hassan Rivaz. Vision transformers for classification of breast ultrasound images. In *Proceedings of the Annual International Confer-*

*ence of the IEEE Engineering in Medicine & Biology Society*, pages 480–483. IEEE, 2022. `https://doi.org/10.1109/EMBC48229.2022.9871809`.

[164] Moi Hoon Yap, Gerard Pons, Joan Marti, Sergi Ganau, Melcior Sentis, Reyer Zwiggelaar, Adrian K Davison, and Robert Marti. Automated breast ultrasound lesions detection using convolutional neural networks. *IEEE J. Biomed. Health Inform.*, 22(4):1218–1226, 2017. `https://doi.org/10.1109/JBHI.2017.2731873`.

[165] Lele Li, Ziling Wu, Juan Liu, Lang Wang, Yu Jin, Peng Jiang, Jing Feng, and Meng Wu. Cross-attention based multi-scale feature fusion vision transformer for breast ultrasound image classification. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1616–1619. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9994966`.

[166] Xin Xing, Gongbo Liang, Yu Zhang, Subash Khanal, Ai-Ling Lin, and Nathan Jacobs. Advit: Vision transformer on multi-modality pet images for alzheimer disease diagnosis. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–4. IEEE, 2022. `https://doi.org/10.1109/ISBI52829.2022.9761584`.

[167] Xin Xing, Gongbo Liang, Hunter Blanton, Muhammad Usman Rafique, Chris Wang, Ai-Ling Lin, and Nathan Jacobs. Dynamic image for 3d mri image alzheimer's disease classification. In *Proceedings of the Computer Vision*, pages 355–364. Springer, 2021. `https://doi.org/10.1007/978-3-030-66415-2_23`.

[168] Moi Hoon Yap, Bill Cassidy, Joseph M Pappachan, Claire O'Shea, David Gillespie, and Neil D Reeves. Analysis towards classification of infection and ischaemia of diabetic foot ulcers. In *Proceedings of the EMBS International Conference on Biomedical and Health Informatics*, pages 1–4. IEEE, 2021. `https://doi.org/10.1109/BHI50953.2021.9508563`.

[169] Haoran Wang, Yanju Ji, Kaiwen Song, Mingyang Sun, Peitong Lv, and Tianyu Zhang. Vit-p: Classification of genitourinary syndrome of menopause from oct images based on vision transformer models. *IEEE Trans. Instrum. Meas.*, 70:1–14, 2021. `https://doi.org/10.1109/TIM.2021.3122121`.

[170] Ronglin Gong, Xiangmin Han, Jun Wang, Shihui Ying, and Jun Shi. Self-supervised bi-channel transformer networks for computer-aided diagnosis. *IEEE J. Biomed. Health Inform.*, 2022. `https://doi.org/10.1109/JBHI.2022.3153902`.

[171] Inês C Moreira, Igor Amaral, Inês Domingues, António Cardoso, Maria Joao Cardoso, and Jaime S Cardoso. Inbreast: toward a full-field digital mammographic database. *Acad. Radiol.*, 19(2):236–248, 2012. `https://doi.org/10.1016/j.acra.2011.09.014`.

[172] Rémi Vallée, Astrid De Maissin, Antoine Coutrot, Harold Mouchère, Arnaud Bourreille, and Nicolas Normand. Crohnipi: An endoscopic image database for the evaluation of automatic crohn's disease lesions recognition algorithms. In *Proceedings of the Medical Imaging 2020: Biomedical Applications in Molecu-*

*lar, Structural, and Functional Imaging*, volume 11317, pages 440–446. SPIE, 2020.

[173] Mohamad Mahmoud Al Rahhal, Yakoub Bazi, Rami M Jomaa, Ahmad AlShibli, Naif Alajlan, Mohamed Lamine Mekhalfi, and Farid Melgani. Covid-19 detection in ct/x-ray imagery using vision transformers. *J. Pers. Med.*, 12(2):310, 2022. `https://doi.org/10.3390/jpm12020310`.

[174] Shen Gao, Xuguang Li, Xin Li, Zhen Li, and Yongqiang Deng. Transformer based tooth classification from cone-beam computed tomography for dental charting. *Comput. Biol. Med.*, 148:105880, 2022. `https://doi.org/10.1016/j.compbiomed.2022.105880`.

[175] Jiancheng Yang, Rui Shi, and Bingbing Ni. Medmnist classification decathlon: A lightweight automl benchmark for medical image analysis. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 191–195. IEEE, 2021. `https://doi.org/10.1109/ISBI48211.2021.9434062`.

[176] Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Sci. Data*, 10(1):41, 2023. `https://doi.org/10.1038/s41597-022-01721-8`.

[177] Jinwei Liu, Yan Li, Guitao Cao, Yong Liu, and Wenming Cao. Feature pyramid vision transformer for medmnist classification decathlon. In *Proceedings of the International Joint Conference on Neural Networks*, pages 1–8. IEEE, 2022. `https://doi.org/10.1109/IJCNN55064.2022.9892282`.

[178] Faris Almalik, Mohammad Yaqub, and Karthik Nandakumar. Self-ensembling vision transformer (sevit) for robust medical image classification. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 376–386. Springer, 2022. `https://doi.org/10.1007/978-3-031-16437-8_36`.

[179] Tawsifur Rahman, Amith Khandakar, Muhammad Abdul Kadir, Khandaker Rejaul Islam, Khandakar F Islam, Rashid Mazhar, Tahir Hamid, Mohammad Tariqul Islam, Saad Kashem, Zaid Bin Mahbub, et al. Reliable tuberculosis detection using chest x-ray with deep learning, segmentation and visualization. *IEEE Access*, 8:191586–191601, 2020. `https://doi.org/10.1109/ACCESS.2020.3031384`.

[180] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *Proceedings of the International conference on machine learning*, pages 10347–10357. PMLR, 2021.

[181] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the International conference on machine learning*, pages 6105–6114. PMLR, 2019.

[182] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages

1492–1500, 2017.

[183] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[184] Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey Shi. Escaping the big data paradigm with compact transformers. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2104.05704`.

[185] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11534–11542, 2020.

[186] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint*, 2014. `https://doi.org/10.48550/arXiv.1409.1556`.

[187] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022.

[188] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.

[189] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint*, 2015. `https://doi.org/10.48550/arXiv.1511.06434`.

[190] Tim Salimans, Ian J Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Proceedings of the Conference on Neural Information Processing Systems*, pages 1–10, 2016.

[191] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the Conference on Neural Information Processing Systems*, pages 1–38, 2017.

[192] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proceedings of the International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. `https://doi.org/10.1007/978-3-319-24574-4_28`.

[193] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2102.04306`.

[194] Yunhe Gao, Mu Zhou, and Dimitris N Metaxas. Utnet: a hybrid transformer architecture for medical image segmentation. In *Proceedings of*

the *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 61–71. Springer, 2021. `https://doi.org/10.1007/978-3-030-87199-4_6`.

[195] Victor M Campello, Polyxeni Gkontra, Cristian Izquierdo, Carlos Martin-Isla, Alireza Sojoudi, Peter M Full, Klaus Maier-Hein, Yao Zhang, Zhiqiang He, Jun Ma, et al. Multi-centre, multi-vendor and multi-disease cardiac segmentation: the m&ms challenge. *IEEE Trans. Med. Imaging*, 40(12):3543–3554, 2021. `https://doi.org/10.1109/TMI.2021.3090082`.

[196] Shaolong Chen, Changzhen Qiu, Weiping Yang, and Zhiyong Zhang. Multiresolution aggregation transformer unet based on multiscale input and coordinate attention for medical image segmentation. *Sensors*, 22(10):3820, 2022. `https://doi.org/10.3390/s22103820`.

[197] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE Trans. Med. Imaging*, 37(11):2514–2525, 2018. `https://doi.org/10.1109/TMI.2018.2837502`.

[198] Zhaohan Xiong, Qing Xia, Zhiqiang Hu, Ning Huang, Cheng Bian, Yefeng Zheng, Sulaiman Vesal, Nishant Ravikumar, Andreas Maier, Xin Yang, et al. A global benchmark of algorithms for segmenting the left atrium from late gadolinium-enhanced cardiac magnetic resonance imaging. *Med. Image Anal.*, 67:101832, 2021. `https://doi.org/10.1016/j.media.2020.101832`.

[199] Junjie Liang, Cihui Yang, Mengjie Zeng, and Xixi Wang. Transconver: transformer and convolution parallel network for developing automatic brain tumor segmentation in mri images. *Quant. Imaging Med. Surg.*, 12(4):2397, 2022. `https://doi.org/10.21037/qims-21-919`.

[200] Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby, John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. *Sci. Data*, 4(1):1–13, 2017. `https://doi.org/10.1038/sdata.2017.117`.

[201] Spyridon Bakas, Mauricio Reyes, Andras Jakab, Stefan Bauer, Markus Rempfler, Alessandro Crimi, Russell Takeshi Shinohara, Christoph Berger, Sung Min Ha, Martin Rozycki, et al. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge. *arXiv preprint*, 2018. `https://doi.org/10.48550/arXiv.1811.02629`.

[202] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE Trans. Med. Imaging*, 34(10):1993–2024, 2014. `https://doi.org/10.1109/TMI.2014.2377694`.

[203] Pan Feng, Bo Ni, Xiantao Cai, and Yutao Xie. Utransnet: Transformer within u-net for stroke lesion segmentation. In *Proceedings of the International Conference on Computer Supported Cooperative Work in Design*, pages 359–364. IEEE, 2022. `https://doi.org/10.1109/CSCWD54268.2022.9776250`.

[204] Anatomical tracings of lesions after stroke. `http://fcon_1000.projects.nitrc.org/indi/retro/atlas.html`. Accessed 4 August 2022.

[205] Wenxuan Wang, Chen Chen, Meng Ding, Hong Yu, Sen Zha, and Jiangyun Li. Transbts: Multimodal brain tumor segmentation using transformer. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 109–119. Springer, 2021. `https://doi.org/10.1007/978-3-030-87193-2_11`.

[206] Jing Wang, Shuyu Wang, and Wei Liang. Metrans: Multi-encoder transformer for ischemic stroke segmentation. *Electron. Lett.*, 58(9):340–342, 2022. `https://doi.org/10.1049/ell2.12444`.

[207] Sook-Lei Liew, Julia M Anglin, Nick W Banks, Matt Sondag, Kaori L Ito, Hosung Kim, Jennifer Chan, Joyce Ito, Connie Jung, Nima Khoshab, et al. A large, open source dataset of stroke anatomical brain images and manual lesion segmentations. *Sci. Data*, 5(1):1–11, 2018. `https://doi.org/10.1038/sdata.2018.11`.

[208] Oskar Maier, Bjoern H Menze, Janina von der Gablentz, Levin Häni, Mattias P Heinrich, Matthias Liebrand, Stefan Winzeck, Abdul Basit, Paul Bentley, Liang Chen, et al. Isles 2015-a public evaluation benchmark for ischemic stroke lesion segmentation from multispectral mri. *Med. Image Anal.*, 35:250–269, 2017. `https://doi.org/10.1016/j.media.2016.07.009`.

[209] Ischemic stroke lesion segmentation. `http://www.isles-challenge.org/ISLES2018/`. Accessed 29 July 2022.

[210] Yun Jiang, Yuan Zhang, Xin Lin, Jinkun Dong, Tongtong Cheng, and Jing Liang. Swinbts: A method for 3d multimodal brain tumor segmentation using swin transformer. *Brain Sci.*, 12(6):797, 2022. `https://doi.org/10.3390/brainsci12060797`.

[211] Segmentation labels and radiomic features for the pre-operative scans of the tcga-gbm collection. `https://wiki.cancerimagingarchive.net/pages/viewpage.action?pageId=24282666`. Accessed 29 July 2022.

[212] Junjie Liang, Cihui Yang, Jingting Zhong, and Xiaoli Ye. Btswin-unet: 3d u-shaped symmetrical swin transformer-based network for brain tumor segmentation with self-supervised pre-training. *Neural Process. Lett.*, pages 1–19, 2022. `https://doi.org/10.1007/s11063-022-10919-1`.

[213] Peixu Wang, Shikun Liu, and Jialin Peng. Ast-net: Lightweight hybrid transformer for multimodal brain tumor segmentation. In *Proceedings of the International Conference on Pattern Recognition*, pages 4623–4629. IEEE, 2022. `https://doi.org/10.1109/ICPR56361.2022.9956705`.

[214] Qiran Jia and Hai Shu. Bitr-unet: a cnn-transformer combined network for mri brain tumor segmentation. In *Proceedings of the Brainlesion: Glioma, Multiple*

*Sclerosis, Stroke and Traumatic Brain Injuries*, pages 3–14. Springer, 2022. `https://doi.org/10.1007/978-3-031-09002-8_1`.

[215] Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R Roth, and Daguang Xu. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In *Proceedings of the Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, pages 272–284. Springer, 2022. `https://doi.org/10.1007/978-3-031-08999-2_22`.

[216] Fazhan Zhu, Jiaxing Lv, Kun Lu, Wenyan Wang, Hongshou Cong, Jun Zhang, Peng Chen, Yuan Zhao, and Ziheng Wu. A 3d medical image segmentation framework fusing convolution and transformer features. In *Proceedings of the Intelligent Computing Theories and Application*, pages 772–786. Springer, 2022. `https://doi.org/10.1007/978-3-031-13870-6_63`.

[217] Himashi Peiris, Munawar Hayat, Zhaolin Chen, Gary Egan, and Mehrtash Harandi. A robust volumetric transformer for accurate 3d tumor segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 162–172. Springer, 2022. `https://doi.org/10.1007/978-3-031-16443-9_16`.

[218] Yu Chen, Ming Yin, Yu Li, and Qian Cai. Csu-net: A cnn-transformer parallel network for multimodal brain tumour segmentation. *Electronics*, 11(14):2226, 2022. `https://doi.org/10.3390/electronics11142226`.

[219] Feng Liu, Jun Zhu, Baolong Lv, Lei Yang, Wenyan Sun, Zhehao Dai, Fangfang Gou, Jia Wu, et al. Auxiliary segmentation method of osteosarcoma mri image based on transformer and u-net. *Comput. Intell. Neurosci.*, 2022, 2022. `https://doi.org/10.1155/2022/9990092`.

[220] Jia Wu, Shun Yang, Fangfang Gou, Zhixun Zhou, Peng Xie, Nuo Xu, and Zhehao Dai. Intelligent segmentation medical assistance system for mri images of osteosarcoma in developing countries. *Comput. Math. Method Med.*, 2022, 2022. `https://doi.org/10.1155/2022/7703583`.

[221] Junjie Liang, Cihui Yang, and Lingguo Zeng. 3d pswinbts: An efficient transformer-based unet using 3d parallel shifted windows for brain tumor segmentation. *Digit. Signal Prog.*, 131:103784, 2022. `https://doi.org/10.1016/j.dsp.2022.103784`.

[222] Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nat. Commun.*, 13(1):4128, 2022. `https://doi.org/10.1038/s41467-022-30695-9`.

[223] Yanmei Chen and Jiajun Wang. Tseunet: A 3d neural network with fused transformer and se-attention for brain tumor segmentation. In *Proceedings of the International Symposium on Computer-Based Medical Systems*, pages 131–136. IEEE, 2022. `https://doi.org/10.1109/CBMS55023.2022.00030`.

[224] Di Gai, Jiqian Zhang, Yusong Xiao, Weidong Min, Yunfei Zhong, and Yuling Zhong. Rmtf-net: Residual mix transformer fusion net for 2d brain tu-

mor segmentation. *Brain Sci.*, 12(9):1145, 2022. `https://doi.org/10.3390/brainsci12091145`.

[225] Brain mri segmentation. `https://www.kaggle.com/datasets/mateuszbuda/lgg-mri-segmentation`. Accessed 29 July 2022.

[226] Shenhai Zheng, Jiaxin Tan, Chuangbo Jiang, and Laquan Li. Automated multi-modal transformer network (amtnet) for 3d medical images segmentation. *Phys. Med. Biol.*, 2022. `https://doi.org/10.1088/1361-6560/aca74c`.

[227] Amber L Simpson, Michela Antonelli, Spyridon Bakas, Michel Bilello, Keyvan Farahani, Bram Van Ginneken, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, et al. A large annotated medical image dataset for the development and evaluation of segmentation algorithms. *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1902.09063`.

[228] Liqun Huang, Long Chen, Baihai Zhang, and Senchun Chai. A transformer-based generative adversarial network for brain tumor segmentation. *arXiv preprint*, 2022. `https://doi.org/10.48550/arXiv.2207.14134`.

[229] Ziqiang Ling, Shun Yang, Fangfang Gou, Zhehao Dai, and Jia Wu. Intelligent assistant diagnosis system of osteosarcoma mri image based on transformer and convolution in developing countries. *IEEE J. Biomed. Health Inform.*, 26(11):5563–5574, 2022. `https://doi.org/10.1109/JBHI.2022.3196043`.

[230] Wei Li and Huihua Yang. Collaborative transformer-cnn learning for semi-supervised medical image segmentation. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1058–1065. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9995501`.

[231] Yao Niu, Zhiming Luo, Sheng Lian, Lei Li, Shaozi Li, and Haixin Song. Symmetrical supervision with transformer for few-shot medical image segmentation. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1683–1687. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9995238`.

[232] A Emre Kavur, N Sinem Gezer, Mustafa Barış, Sinem Aslan, Pierre-Henri Conze, Vladimir Groza, Duc Duy Pham, Soumick Chatterjee, Philipp Ernst, Savaş Özkan, et al. Chaos challenge-combined (ct-mr) healthy abdominal organ segmentation. *Med. Image Anal.*, 69:101950, 2021. `https://doi.org/10.1016/j.media.2020.101950`.

[233] Xiahai Zhuang. Multivariate mixture model for myocardial segmentation combining multi-source images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(12):2933–2946, 2018. `https://doi.org/10.1109/TPAMI.2018.2869576`.

[234] Zheyao Gao and Xiahai Zhuang. Consistency based co-segmentation for multi-view cardiac mri using vision transformer. In *Proceedings of the International Workshop on Statistical Atlases and Computational Models of the Heart*, pages 306–314. Springer, 2021. `https://doi.org/10.1007/978-3-030-93722-5_33`.

[235] Multi-disease, multi-view & multi-center right ventricular segmentation in cardiac mri. `https://www.ub.edu/mnms-2/`. Accessed 4 August 2022.

[236] Abel A Reyes, Sidike Paheding, Makarand Deo, and Michel Audette. Gabor filter-embedded u-net with transformer-based encoding for biomedical image segmentation. In *Proceedings of the Multiscale Multimodal Medical Imaging*, pages 76–88. Springer, 2022. https://doi.org/10.1007/978-3-031-18814-5_8.

[237] Zhiyong Xiao, Yixin Su, Zhaohong Deng, and Weidong Zhang. Efficient combination of cnn and transformer for dual-teacher uncertainty-guided semi-supervised medical image segmentation. *Comput. Meth. Programs Biomed.*, 226:107099, 2022. https://doi.org/10.1016/j.cmpb.2022.107099.

[238] Yao Zhang, Nanjun He, Jiawei Yang, Yuexiang Li, Dong Wei, Yawen Huang, Yang Zhang, Zhiqiang He, and Yefeng Zheng. mmformer: Multimodal medical transformer for incomplete multimodal learning of brain tumor segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 107–117. Springer, 2022. https://doi.org/10.1007/978-3-031-16443-9_11.

[239] Bingjie Li, Tiejun Yang, and Xiang Zhao. Nvtrans-unet: Neighborhood vision transformer based u-net for multi-modal cardiac mr image segmentation. *J. Appl. Clin. Med. Phys*, page e13908, 2023. https://doi.org/10.1002/acm2.13908.

[240] Lei Li, Fuping Wu, Sihan Wang, Xinzhe Luo, Carlos Martin-Isla, Shuwei Zhai, Jianpeng Zhang, Yanfei Liu, Zhen Zhang, Markus J Ankenbrand, et al. Myops: A benchmark of myocardial pathology segmentation combining three-sequence cardiac magnetic resonance images. *arXiv preprint*, 2022. https://doi.org/10.48550/arXiv.2201.03186.

[241] Davood Karimi, Haoran Dou, and Ali Gholipour. Medical image segmentation using transformer networks. *IEEE Access*, 10:29322–29332, 2022. https://doi.org/10.1109/ACCESS.2022.3156894.

[242] Mri hippocampus segmentation. https://www.kaggle.com/datasets/sabermalek/mrihs. Accessed 29 July 2022.

[243] Qin Liu, Zhenlin Xu, Yining Jiao, and Marc Niethammer. isegformer: Interactive segmentation via transformers with application to 3d knee mr images. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 464–474. Springer, 2022. https://doi.org/10.1007/978-3-031-16443-9_45.

[244] Felix Ambellan, Alexander Tack, Moritz Ehlke, and Stefan Zachow. Automated segmentation of knee bone and cartilage combining statistical shape knowledge and convolutional neural networks: Data from the osteoarthritis initiative. *Med. Image Anal.*, 52:109–118, 2019. https://doi.org/10.1016/j.media.2018.11.009.

[245] Ziyang Wang, Nanqing Dong, and Irina Voiculescu. Computationally-efficient vision transformer for medical image semantic segmentation via dual pseudo-label supervision. In *Proceedings of the IEEE International Conference on Image Processing*, pages 1961–1965. IEEE, 2022. https://doi.org/10.1109/

ICIP46576.2022.9897482.

[246] Ziyang Wang, Jian-Qing Zheng, and Irina Voiculescu. An uncertainty-aware transformer for mri cardiac semantic segmentation via mean teachers. In *Proceedings of the Medical Image Understanding and Analysis*, pages 494–507. Springer, 2022. `https://doi.org/10.1007/978-3-031-12053-4_37`.

[247] Qixuan Sun, Nianhua Fang, Zhuo Liu, Liang Zhao, Youpeng Wen, and Hongxiang Lin. Hybridctrm: Bridging cnn and transformer for multimodal brain image segmentation. *J. Healthc. Eng.*, 2021, 2021. `https://doi.org/10.1155/2021/7467261`.

[248] Adriënne M Mendrik, Koen L Vincken, Hugo J Kuijf, Marcel Breeuwer, Willem H Bouvy, Jeroen De Bresser, Amir Alansary, Marleen De Bruijne, Aaron Carass, Ayman El-Baz, et al. Mrbrains challenge: online evaluation framework for brain image segmentation in 3t mri scans. *Comput. Intell. Neurosci.*, 2015, 2015. `https://doi.org/10.1155/2015/813696`.

[249] Li Wang, Dong Nie, Guannan Li, Élodie Puybareau, Jose Dolz, Qian Zhang, Fan Wang, Jing Xia, Zhengwang Wu, Jia-Wei Chen, et al. Benchmark on automatic six-month-old infant brain segmentation algorithms: the iseg-2017 challenge. *IEEE Trans. Med. Imaging*, 38(9):2219–2230, 2019. `https://doi.org/10.1109/TMI.2019.2901712`.

[250] Yang Xu, Xianyu He, Guofeng Xu, Guanqiu Qi, Kun Yu, Li Yin, Pan Yang, Yuehui Yin, and Hao Chen. A medical image segmentation method based on multi-dimensional statistical features. *Frontiers in Neuroscience*, 16, 2022.

[251] Zhiqin Zhu, Xianyu He, Guanqiu Qi, Yuanyuan Li, Baisen Cong, and Yu Liu. Brain tumor segmentation based on the fusion of deep semantics and edge information in multimodal mri. *Information Fusion*, 91:376–387, 2023.

[252] Ramtin Mojtahedi, Mohammad Hamghalam, Richard KG Do, and Amber L Simpson. Towards optimal patch size in vision transformers for tumor segmentation. In *Proceedings of the Multiscale Multimodal Medical Imaging*, pages 110–120. Springer, 2022. `https://doi.org/10.1007/978-3-031-18814-5_11`.

[253] Patrick Bilic, Patrick Ferdinand Christ, Eugene Vorontsov, Grzegorz Chlebus, Hao Chen, Qi Dou, Chi-Wing Fu, Xiao Han, Pheng-Ann Heng, Jürgen Hesser, et al. The liver tumor segmentation benchmark (lits). *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1901.04056`.

[254] Amber L Simpson, Alexandre Doussot, John M Creasy, Lauryn B Adams, Peter J Allen, Ronald P DeMatteo, Mithat Gönen, Nancy E Kemeny, T Peter Kingham, Jinru Shia, et al. Computed tomography image texture: a noninvasive prognostic marker of hepatic recurrence after hepatectomy for metastatic colorectal cancer. *Ann. Surg. Oncol.*, 24:2482–2490, 2017. `https://doi.org/10.1245/s10434-017-5896-1`.

[255] Yutong Xie, Jianpeng Zhang, Chunhua Shen, and Yong Xia. Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation. In *Proceedings of the International conference on medical image computing*

*and computer-assisted intervention*, pages 171–180. Springer, 2021. `https://doi.org/10.1007/978-3-030-87199-4_16`.

[256] Multi-atlas labeling beyond the cranial vault - workshop and challenge. `https://www.synapse.org/#!Synapse:syn3193805/wiki/217789`. Accessed 29 July 2022.

[257] Hongyu Kan, Jun Shi, Minfan Zhao, Zhaohui Wang, Wenting Han, Hong An, Zhaoyang Wang, and Shuo Wang. Itunet: Integration of transformers and unet for organs-at-risk segmentation. In *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine & Biology Society*, pages 2123–2127. IEEE, 2022. `https://doi.org/10.1109/EMBC48229.2022.9871945`.

[258] Zihan Li, Dihan Li, Cangbai Xu, Weice Wang, Qingqi Hong, Qingde Li, and Jie Tian. Tfcns: A cnn-transformer hybrid network for medical image segmentation. In *Proceedings of the Artificial Neural Networks and Machine Learning*, pages 781–792. Springer, 2022. `https://doi.org/10.1007/978-3-031-15937-4_65`.

[259] Covid-19 ct scans. `https://www.kaggle.com/datasets/andrewmvd/covid19-ct-scans`. Accessed 17 February 2023.

[260] Jingdong Yang, Jun Tu, Xiaolin Zhang, Shaoqing Yu, and Xianyou Zheng. Tse deeplab: An efficient visual transformer for medical image segmentation. *Biomed. Signal Process. Control*, 80:104376, 2023. `https://doi.org/10.1016/j.bspc.2022.104376`.

[261] Danfeng Guo and Demetri Terzopoulos. A transformer-based network for anisotropic 3d medical image segmentation. In *Proceedings of the International Conference on Pattern Recognition*, pages 8857–8861. IEEE, 2021. `https://doi.org/10.1109/ICPR48806.2021.9411990`.

[262] Xiangyi Yan, Hao Tang, Shanlin Sun, Haoyu Ma, Deying Kong, and Xiaohui Xie. After-unet: Axial fusion transformer unet for medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3971–3981, 2022.

[263] Xuming Chen, Shanlin Sun, Narisu Bai, Kun Han, Qianqian Liu, Shengyu Yao, Hao Tang, Chupeng Zhang, Zhipeng Lu, Qian Huang, et al. A deep learning-based auto-segmentation system for organs-at-risk on whole-body computed tomography images for radiation therapy. *Radiother. Oncol.*, 160:175–184, 2021. `https://doi.org/10.1016/j.radonc.2021.04.019`.

[264] Zoé Lambert, Caroline Petitjean, Bernard Dubray, and Su Kuan. Segthor: segmentation of thoracic organs at risk in ct images. In *Proceedings of the International Conference on Image Processing Theory, Tools and Applications*, pages 1–6. IEEE, 2020. `https://doi.org/10.1109/IPTA50016.2020.9286453`.

[265] Mingjun Ma, Haiying Xia, Yumei Tan, Haisheng Li, and Shuxiang Song. Htnet: hierarchical context-attention transformer network for medical ct image segmentation. *Appl. Intell.*, pages 1–14, 2022. `https://doi.org/10.1007/s10489-021-03010-0`.

[266] Nicholas Heller, Niranjan Sathianathen, Arveen Kalapara, Edward Walczak, Keenan Moore, Heather Kaluzniak, Joel Rosenberg, Paul Blake, Zachary Rengel, Makinna Oestreich, et al. The kits19 challenge data: 300 kidney tumor cases with clinical context, ct semantic segmentations, and surgical outcomes. *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1904.00445`.

[267] Finding and measuring lungs in ct data. `https://www.kaggle.com/datasets/kmader/finding-lungs-in-ct-data`. Accessed 29 July 2022.

[268] Bladder datasets. `https://cdas.cancer.gov/datasets/plco/18/#:~:text=The%20Bladder%20dataset%20is%20a,participants%20in%20the%20PLCO%20trial`. Accessed 30 July 2022.

[269] Chun Luo, Jing Zhang, Xinglin Chen, Yinhao Tang, Xiechuan Weng, and Fan Xu. 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 *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine & Biology Society*, pages 3565–3568. IEEE, 2021. `https://doi.org/10.1109/EMBC46164.2021.9630336`.

[270] Yuan Yang, Lin Zhang, Lei Ren, and Xiaohan Wang. Mmvit-seg: A lightweight transformer and cnn fusion network for covid-19 segmentation. *Comput. Meth. Programs Biomed.*, page 107348, 2023. `https://doi.org/10.1016/j.cmpb.2023.107348`.

[271] Deng-Ping Fan, Tao Zhou, Ge-Peng Ji, Yi Zhou, Geng Chen, Huazhu Fu, Jianbing Shen, and Ling Shao. Inf-net: Automatic covid-19 lung infection segmentation from ct images. *IEEE Trans. Med. Imaging*, 39(8):2626–2637, 2020. `https://doi.org/10.1109/TMI.2020.2996645`.

[272] Covid-19 ct segmentation dataset. `https://medicalsegmentation.com/covid19/`. Accessed 17 February 2023.

[273] Mingyang Liu, Li Xiao, Huiqin Jiang, and Qing He. Ccat-net: A novel transformer based semi-supervised framework for covid-19 lung lesion segmentation. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–5. IEEE, 2022. `https://doi.org/10.1109/ISBI52829.2022.9761533`.

[274] Yang Ning, Shouyi Zhang, Xiaoming Xi, Jie Guo, Peide Liu, and Caiming Zhang. Cac-emvt: Efficient coronary artery calcium segmentation with multi-scale vision transformers. In *Proceedings of the International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1462–1467. IEEE, 2021. `https://doi.org/10.1109/BIBM52615.2021.9669337`.

[275] Luyao Wang, Xiaoyan Wang, Bangze Zhang, Xiaojie Huang, Cong Bai, Ming Xia, and Peiliang Sun. Multi-scale hierarchical transformer structure for 3d medical image segmentation. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1542–1545. IEEE, 2021. `https://doi.org/10.1109/BIBM52615.2021.9669799`.

[276] Yang Ning, Shouyi Zhang, Wei Zhong, Peide Liu, and Caiming Zhang. A hybrid cross-scale transformer architecture for robust medical image segmenta-

tion. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine*, pages 1677–1682. IEEE, 2022. `https://doi.org/10.1109/BIBM55620.2022.9995702`.

[277] Xin You, Yun Gu, Junjun He, Hui Sun, and Jie Yang. A more design-flexible medical transformer for volumetric image segmentation. In *Proceedings of the Machine Learning in Medical Imaging*, pages 62–71. Springer, 2022. `https://doi.org/10.1007/978-3-031-21014-3_7`.

[278] Anjany Sekuboyina, Malek E Husseini, Amirhossein Bayat, Maximilian Löffler, Hans Liebl, Hongwei Li, Giles Tetteh, Jan Kukačka, Christian Payer, Darko Štern, et al. Verse: A vertebrae labelling and segmentation benchmark for multi-detector ct images. *Med. Image Anal.*, 73:102166, 2021. `https://doi.org/10.1016/j.media.2021.102166`.

[279] Duy-Phuong Dao, Hyung-Jeong Yang, Ngoc-Huynh Ho, Sudarshan Pant, Soo-Hyung Kim, Guee-Sang Lee, In-Jae Oh, and Sae-Ryung Kang. Survival analysis based on lung tumor segmentation using global context-aware transformer in multimodality. In *Proceedings of the International Conference on Pattern Recognition*, pages 5162–5169. IEEE, 2022. `https://doi.org/10.1109/ICPR56361.2022.9956406`.

[280] Hugo JWL Aerts, Emmanuel Rios Velazquez, Ralph TH Leijenaar, Chintan Parmar, Patrick Grossmann, Sara Carvalho, Johan Bussink, René Monshouwer, Benjamin Haibe-Kains, Derek Rietveld, et al. Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach. *Nat. Commun.*, 5(1):4006, 2014. `https://doi.org/10.1038/ncomms5006`.

[281] Lifang Chen and Li Wan. Ctunet: automatic pancreas segmentation using a channel-wise transformer and 3d u-net. *Visual Comput.*, pages 1–15, 2022. `https://doi.org/10.1007/s00371-022-02656-2`.

[282] Holger R Roth, Amal Farag, Le Lu, Evrim B Turkbey, and Ronald M Summers. Deep convolutional networks for pancreas segmentation in ct imaging. In *Proceedings of the Medical Imaging*, volume 9413, pages 378–385. SPIE, 2015. `https://doi.org/10.1117/12.2081420`.

[283] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2105.05537`.

[284] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 574–584, 2022.

[285] Simon Jégou, Michal Drozdzal, David Vazquez, Adriana Romero, and Yoshua Bengio. The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 11–19, 2017.

[286] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam.

Rethinking atrous convolution for semantic image segmentation. *arXiv preprint*, 2017. `https://doi.org/10.48550/arXiv.1706.05587`.

[287] Xiaoying Pan, Weidong Bai, Minjie Ma, and Shaoqiang Zhang. Rant: A cascade reverse attention segmentation framework with hybrid transformer for laryngeal endoscope images. *Biomed. Signal Process. Control*, 78:103890, 2022. `https://doi.org/10.1016/j.bspc.2022.103890`.

[288] Max-Heinrich Laves, Jens Bicker, Lüder A Kahrs, and Tobias Ortmaier. A dataset of laryngeal endoscopic images with comparative study on convolution neural network-based semantic segmentation. *Proceedings of the International journal of computer assisted radiology and surgery*, 14(3):483–492, 2019. `https://doi.org/10.1007/s11548-018-01910-0`.

[289] Shu Tang, Junlin Qiu, Xianzhong Xie, Haiheng Ran, and Guoli Zhang. Bidfnet: Bi-decoder and feedback network for automatic polyp segmentation with vision transformers. In *Proceedings of the Pattern Recognition and Computer Vision*, pages 16–27. Springer, 2022. `https://doi.org/10.1007/978-3-031-18910-4_2`.

[290] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *Proceedings of the International Conference on Multimedia Modeling*, pages 451–462. Springer, 2020. `https://doi.org/10.1007/978-3-030-37734-2_37`.

[291] Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Comput. Med. Imaging Graph.*, 43:99–111, 2015. `https://doi.org/10.1016/j.compmedimag.2015.02.007`.

[292] Nima Tajbakhsh, Suryakanth R Gurudu, and Jianming Liang. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE Trans. Med. Imaging*, 35(2):630–644, 2015. `https://doi.org/10.1109/TMI.2015.2487997`.

[293] David Vázquez, Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Antonio M López, Adriana Romero, Michal Drozdzal, and Aaron Courville. A benchmark for endoluminal scene segmentation of colonoscopy images. *J. Healthc. Eng.*, 2017, 2017. `https://doi.org/10.1155/2017/4037190`.

[294] Yanglan Ou, Ye Yuan, Xiaolei Huang, Stephen TC Wong, John Volpi, James Z Wang, and Kelvin Wong. Patcher: Patch transformers with mixture of experts for precise medical image segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 475–484. Springer, 2022. `https://doi.org/10.1007/978-3-031-16443-9_46`.

[295] Vittorino Mandujano-Cornejo and Javier A Montoya-Zegarra. Polyp2seg: Improved polyp segmentation with vision transformer. In *Proceedings of the Medical Image Understanding and Analysis*, pages 519–534. Springer, 2022. `https://doi.org/10.1007/978-3-031-12053-4_39`.

[296] Qian Wang, Longyan Li, Bo Ni, Yu Li, Dejin Kong, Chen Wang, and Zan Li. Medical image segmentation using transformer. In *Proceedings of the Artificial Intelligence in China*, pages 92–99. Springer, 2022. `https://doi.org/10.1007/978-981-16-9423-3_12`.

[297] Edward Sanderson and Bogdan J Matuszewski. Fcn-transformer feature fusion for polyp segmentation. In *Proceedings of the Medical Image Understanding and Analysis*, pages 892–907. Springer, 2022. `https://doi.org/10.1007/978-3-031-12053-4_65`.

[298] Dataset of endoscopic colonoscopy frames for polyp detection. `https://www.kaggle.com/datasets/balraj98/cvcclinicdb`. Accessed 18 February 2023.

[299] Nurbek Saidnassim, Beibit Abdikenov, Rauan Kelesbekov, Muhammad Tahir Akhtar, and Prashant Jamwal. Self-supervised visual transformers for breast cancer diagnosis. In *Proceedings of the Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*, pages 423–427. IEEE, 2021.

[300] Guifang Zhang, Hon-Cheng Wong, Cheng Wang, Jianjun Zhu, Ligong Lu, and Gaojun Teng. A temporary transformer network for guide-wire segmentation. In *Proceedings of the International Congress on Image and Signal Processing, BioMedical Engineering and Informatics*, pages 1–5. IEEE, 2021. `https://doi.org/10.1109/CISP-BMEI53629.2021.9624350`.

[301] Lingrong Zhang, Jinglin Yang, Dong Liu, Feng Zhang, Sibo Nie, Yuchen Tan, and Taipeng Guo. Spine x-ray image segmentation based on transformer and adaptive optimized postprocessing. In *Proceedings of the International Conference on Software Engineering and Artificial Intelligence*, pages 88–92. IEEE, 2022. `https://doi.org/10.1109/SEAI55746.2022.9832144`.

[302] Accurate automated spinal curvature estimation. `https://aasce19.github.io/`. Accessed 18 February 2023.

[303] Siim-acr pneumothorax segmentation. `https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation`. Accessed 19 February 2023.

[304] Kaizhong Deng, Yanda Meng, Dongxu Gao, Joshua Bridge, Yaochun Shen, Gregory Lip, Yitian Zhao, and Yalin Zheng. Transbridge: A lightweight transformer for left ventricle segmentation in echocardiography. In *Proceedings of the International Workshop on Advances in Simplifying Medical Ultrasound*, pages 63–72. Springer, 2021. `https://doi.org/10.1007/978-3-030-87583-1_7`.

[305] Tao Wang, Zhihui Lai, and Heng Kong. Tfnet: Transformer fusion network for ultrasound image segmentation. In *Proceedings of the Pattern Recognition*, pages 314–325. Springer, 2022. `https://doi.org/10.1007/978-3-031-02375-0_23`.

[306] Lina Pedraza, Carlos Vargas, Fabián Narváez, Oscar Durán, Emma Muñoz, and Eduardo Romero. An open access thyroid ultrasound image database. In *Proceedings of the International symposium on medical information processing and analysis*, volume 9287, pages 188–193. SPIE, 2015. `https:`

`//doi.org/10.1117/12.2073532`.

[307] Haonan Yang and Dapeng Yang. Cswin-pnet: A cnn-swin transformer combined pyramid network for breast lesion segmentation in ultrasound images. *Expert Syst. Appl.*, 213:119024, 2023. `https://doi.org/10.1016/j.eswa.2022.119024`.

[308] Xianwei Zhuang, Xiner Zhu, Haoji Hu, Jincao Yao, Wei Li, Chen Yang, Liping Wang, Na Feng, and Dong Xu. Residual swin transformer unet with consistency regularization for automatic breast ultrasound tumor segmentation. In *Proceedings of the IEEE International Conference on Image Processing*, pages 3071–3075. IEEE, 2022. `https://doi.org/10.1109/ICIP46576.2022.9897941`.

[309] Xiaoyan Shen, Liangyu Wang, Yu Zhao, Ruibo Liu, Wei Qian, and He Ma. Dilated transformer: residual axial attention for breast ultrasound image segmentation. *Quant. Imaging Med. Surg.*, 12(9):4513, 2022. `https://doi.org/10.21037/qims-22-33`.

[310] Yingtao Zhang, Min Xian, Heng-Da Cheng, Bryar Shareef, Jianrui Ding, Fei Xu, Kuan Huang, Boyu Zhang, Chunping Ning, and Ying Wang. Busis: A benchmark for breast ultrasound image segmentation. In *Healthcare*, volume 10, page 729. MDPI, 2022. `https://doi.org/10.3390/healthcare10040729`.

[311] Zhihao Liao, Neng Fan, and Kai Xu. Swin transformer assisted prior attention network for medical image segmentation. *Appl. Sci.*, 12(9):4735, 2022. `https://doi.org/10.3390/app12094735`.

[312] Korsuk Sirinukunwattana, Josien PW Pluim, Hao Chen, Xiaojuan Qi, Pheng-Ann Heng, Yun Bo Guo, Li Yang Wang, Bogdan J Matuszewski, Elia Bruni, Urko Sanchez, et al. Gland segmentation in colon histology images: The glas challenge contest. *Med. Image Anal.*, 35:489–502, 2017. `https://doi.org/10.1016/j.media.2016.08.008`.

[313] Neeraj Kumar, Ruchika Verma, Sanuj Sharma, Surabhi Bhargava, Abhishek Vahadane, and Amit Sethi. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE Trans. Med. Imaging*, 36(7):1550–1560, 2017. `https://doi.org/10.1109/TMI.2017.2677499`.

[314] Ziniu Qian, Kailu Li, Maode Lai, Eric I-Chao Chang, Bingzheng Wei, Yubo Fan, and Yan Xu. Transformer based multiple instance learning for weakly supervised histopathology image segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 160–170. Springer, 2022. `https://doi.org/10.1007/978-3-031-16434-7_16`.

[315] Zhipeng Jia, Xingyi Huang, I Eric, Chao Chang, and Yan Xu. Constrained deep weak supervision for histopathology image segmentation. *IEEE Trans. Med. Imaging*, 36(11):2376–2388, 2017. `https://doi.org/10.1109/TMI.2017.2724070`.

[316] Ziheng Wang, Xiongkuo Min, Fangyu Shi, Ruinian Jin, Saida S Nawrin, Ichen Yu, and Ryoichi Nagatomi. Smeswin unet: Merging cnn and transformer for

medical image segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 517–526. Springer, 2022. `https://doi.org/10.1007/978-3-031-16443-9_50`.

[317] Navid Alemi Koohbanani, Mostafa Jahanifar, Neda Zamani Tajadin, and Nasir Rajpoot. Nuclick: a deep learning framework for interactive segmentation of microscopic images. *Med. Image Anal.*, 65:101771, 2020. `https://doi.org/10.1016/j.media.2020.101771`.

[318] Joes Staal, Michael D Abràmoff, Meindert Niemeijer, Max A Viergever, and Bram Van Ginneken. Ridge-based vessel segmentation in color images of the retina. *IEEE Trans. Med. Imaging*, 23(4):501–509, 2004. `https://doi.org/10.1109/TMI.2004.825627`.

[319] Christopher G Owen, Alicja R Rudnicka, Robert Mullen, Sarah A Barman, Dorothy Monekosso, Peter H Whincup, Jeffrey Ng, and Carl Paterson. Measuring retinal vessel tortuosity in 10-year-old children: validation of the computer-assisted image analysis of the retina (caiar) program. *Invest. Ophthalmol. Vis. Sci.*, 50(5):2004–2010, 2009. `https://doi.org/10.1167/iovs.08-3018`.

[320] AD Hoover, Valentina Kouznetsova, and Michael Goldbaum. Locating blood vessels in retinal images by piecewise threshold probing of a matched filter response. *IEEE Trans. Med. Imaging*, 19(3):203–210, 2000. `https://doi.org/10.1109/42.845178`.

[321] Yaowei Feng, Zhendong Li, Dong Yang, Hongkai Hu, Hui Guo, and Hao Liu. Polarformer: Optic disc and cup segmentation using a hybrid cnn-transformer and polar transformation. *Appl. Sci.*, 13(1):541, 2022. `https://doi.org/10.3390/app13010541`.

[322] José Ignacio Orlando, Huazhu Fu, João Barbosa Breda, Karel Van Keer, Deepti R Bathula, Andrés Diaz-Pinto, Ruogu Fang, Pheng-Ann Heng, Jeyoung Kim, JoonHo Lee, et al. Refuge challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs. *Med. Image Anal.*, 59:101570, 2020. `https://doi.org/10.1016/j.media.2019.101570`.

[323] Jayanthi Sivaswamy, SR Krishnadas, Gopal Datt Joshi, Madhulika Jain, and A Ujjwaft Syed Tabish. Drishti-gs: Retinal image dataset for optic nerve head (onh) segmentation. In *Proceedings of the international symposium on biomedical imaging*, pages 53–56. IEEE, 2014. `https://doi.org/10.1109/ISBI.2014.6867807`.

[324] Francisco Fumero, Silvia Alayón, José L Sanchez, Jose Sigut, and M Gonzalez-Hernandez. Rim-one: An open retinal image database for optic nerve evaluation. In *Proceedings of the international symposium on computer-based medical systems*, pages 1–6. IEEE, 2011. `https://doi.org/10.1109/CBMS.2011.5999143`.

[325] Yang Li, Yue Zhang, Jing-Yu Liu, Kang Wang, Kai Zhang, Gen-Sheng Zhang, Xiao-Feng Liao, and Guang Yang. Global transformer and dual local attention network via deep-shallow hierarchical feature fusion for retinal vessel seg-

mentation. *IEEE T. Cybern.*, 2022. `https://doi.org/10.1109/TCYB.2022.3194099`.

[326] Venkateswararao Cherukuri, Vijay Kumar Bg, Raja Bala, and Vishal Monga. Deep retinal image segmentation with regularization under geometric priors. *IEEE Trans. Image Process.*, 29:2552–2567, 2019. `https://doi.org/10.1109/TIP.2019.2946078`.

[327] Muhammad Moazam Fraz, Paolo Remagnino, Andreas Hoppe, Bunyarit Uyyanonvara, Alicja R Rudnicka, Christopher G Owen, and Sarah A Barman. An ensemble classification-based approach applied to retinal blood vessel segmentation. *IEEE Trans. Biomed. Eng.*, 59(9):2538–2548, 2012. `https://doi.org/10.1109/TBME.2012.2205687`.

[328] Masum Shah Junayed, Md Baharul Islam, and Nipa Anjum. A transformer-based versatile network for acne vulgaris segmentation. In *Proceedings of the Innovations in Intelligent Systems and Applications Conference*, pages 1–6. IEEE, 2022. `https://doi.org/10.1109/ASYU56188.2022.9925323`.

[329] Mohammad D Alahmadi and Wajdi Alghamdi. Semi-supervised skin lesion segmentation with coupling cnn and transformer features. *IEEE Access*, 10:122560–122569, 2022. `https://doi.org/10.1109/ACCESS.2022.3224005`.

[330] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1902.03368`.

[331] Teresa Mendonça, Pedro M Ferreira, Jorge S Marques, André RS Marcal, and Jorge Rozeira. Ph 2-a dermoscopic image database for research and benchmarking. In *Proceedings of the annual international conference of the IEEE engineering in medicine and biology society*, pages 5437–5440. IEEE, 2013. `https://doi.org/10.1109/EMBC.2013.6610779`.

[332] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021.

[333] Zhanyu Wang, Mingkang Tang, Lei Wang, Xiu Li, and Luping Zhou. A medical semantic-assisted transformer for radiographic report generation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 655–664. Springer, 2022. `https://doi.org/10.1007/978-3-031-16437-8_63`.

[334] Ping Chao, Chao-Yang Kao, Yu-Shan Ruan, Chien-Hsiang Huang, and Youn-Long Lin. Hardnet: A low memory traffic network. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3552–3561, 2019.

[335] Jeya Maria Jose Valanarasu, Poojan Oza, Ilker Hacihaliloglu, and Vishal M

Patel. Medical transformer: Gated axial-attention for medical image segmentation. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 36–46. Springer, 2021. `https://doi.org/10.1007/978-3-030-87193-2_4`.

[336] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. *arXiv preprint*, 2018. `https://doi.org/10.48550/arXiv.1807.10165`.

[337] Aleksandar Vakanski, Min Xian, and Phoebe E Freer. Attention-enriched deep learning model for breast tumor segmentation in ultrasound images. *Ultrasound Med. Biol.*, 46(10):2819–2833, 2020. `https://doi.org/10.1016/j.ultrasmedbio.2020.06.015`.

[338] Haonan Wang, Peng Cao, Jiaqi Wang, and Osmar R Zaiane. Uctransnet: rethinking the skip connections in u-net from a channel-wise perspective with transformer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 2441–2449, 2022.

[339] Xiaohong Huang, Zhifang Deng, Dandan Li, Xueguang Yuan, and Ying Fu. Missformer: An effective transformer for 2d medical image segmentation. *IEEE Trans. Med. Imaging*, 2022. `https://doi.org/10.1109/TMI.2022.3230943`.

[340] Automated cardiac diagnosis challenge. `https://acdc.creatis.insa-lyon.fr/description/databases.html`. Accessed 4 August 2022.

[341] Zhifang Hong, Mingzhi Chen, Weijie Hu, Shiyu Yan, Aiping Qu, Lingna Chen, and Junxi Chen. Dual encoder network with transformer-cnn for multi-organ segmentation. *Med. Biol. Eng. Comput.*, pages 1–11, 2022. `https://doi.org/10.1007/s11517-022-02723-9`.

[342] Ailiang Lin, Jiayu Xu, Jinxing Li, and Guangming Lu. Contrans: Improving transformer with convolutional attention for medical image segmentation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 297–307. Springer, 2022. `https://doi.org/10.1007/978-3-031-16443-9_29`.

[343] Covid-19 ct segmentation dataset. `http://medicalsegmentation.com/covid19/`. Accessed 29 July 2022.

[344] Juan C Caicedo, Allen Goodman, Kyle W Karhohs, Beth A Cimini, Jeanelle Ackerman, Marzieh Haghighi, CherKeng Heng, Tim Becker, Minh Doan, Claire McQuin, et al. Nucleus segmentation across imaging experiments: the 2018 data science bowl. *Nat. Methods*, 16(12):1247–1253, 2019. `https://doi.org/10.1038/s41592-019-0612-7`.

[345] Jevgenij Gamper, Navid Alemi Koohbanani, Ksenija Benet, Ali Khuram, and Nasir Rajpoot. Pannuke: an open pan-cancer histology dataset for nuclei instance segmentation and classification. In *Proceedings of the Digital Pathology*, pages 11–19. Springer, 2019. `https://doi.org/10.1007/978-3-030-23937-4_2`.

[346] Huimin Huang, Shiao Xie, Lanfen Lin, Yutaro Iwamoto, Xianhua Han, Yen-

Wei Chen, and Ruofeng Tong. Scaleformer: Revisiting the transformer-based backbones from a scale-wise perspective for medical image segmentation. *arXiv preprint*, 2022. https://doi.org/10.48550/arXiv.2207.14552.

[347] Bennett Landman, Zhoubing Xu, Juan Eugenio Igelsias, M Styner, TR Langerak, and A Klein. Segmentation outside the cranial vault challenge. *Synapse*, 2015.

[348] Abhinav Sagar. Emsvit: Efficient multi scale vision transformer for biomedical image segmentation. In *Proceedings of the Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, pages 39–51. Springer, 2022. https://doi.org/10.1007/978-3-031-08999-2_3.

[349] Shen Jiang and Jinjiang Li. Transcunet: Unet cross fused transformer for medical image segmentation. *Comput. Biol. Med.*, 150:106207, 2022. https://doi.org/10.1016/j.compbiomed.2022.106207.

[350] Neeraj Kumar, Ruchika Verma, Deepak Anand, Yanning Zhou, Omer Fahri Onder, Efstratios Tsougenis, Hao Chen, Pheng-Ann Heng, Jiahui Li, Zhiqiang Hu, et al. A multi-organ nucleus segmentation challenge. *IEEE Trans. Med. Imaging*, 39(5):1380–1391, 2019. https://doi.org/10.1109/TMI.2019.2947628.

[351] Hao Li, Dewei Hu, Han Liu, Jiacheng Wang, and Ipek Oguz. Cats: Complementary cnn and transformer encoders for segmentation. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–5. IEEE, 2022. https://doi.org/10.1109/ISBI52829.2022.9761596.

[352] Cross-modality domain adaptation for medical image segmentation - 2021. https://crossmoda.grand-challenge.org/. Accessed 19 February 2023.

[353] Generalisable 3d semantic segmentation. http://medicaldecathlon.com/. Accessed 19 February 2023.

[354] Yixuan Wu, Kuanlun Liao, Jintai Chen, Jinhong Wang, Danny Z Chen, Honghao Gao, and Jian Wu. D-former: A u-shaped dilated transformer for 3d medical image segmentation. *Neural Comput. Appl.*, pages 1–14, 2022. https://doi.org/10.1007/s00521-022-07859-1.

[355] Ning Zhang, Long Yu, Dezhi Zhang, Weidong Wu, Shengwei Tian, and Xiaojing Kang. Apt-net: Adaptive encoding and parallel decoding transformer for medical image segmentation. *Comput. Biol. Med.*, 151:106292, 2022. https://doi.org/10.1016/j.compbiomed.2022.106292.

[356] Reza Azad, Mohammad T Al-Antary, Moein Heidari, and Dorit Merhof. Transnorm: Transformer provides a strong spatial normalization mechanism for a deep segmentation model. *IEEE Access*, 10:108205–108215, 2022. https://doi.org/10.1109/ACCESS.2022.3211501.

[357] Anubha Gupta, Pramit Mallick, Ojaswa Sharma, Ritu Gupta, and Rahul Duggal. Pcseg: Color model driven probabilistic multiphase level set based tool for plasma cell segmentation in multiple myeloma. *PloS one*, 13(12):e0207908, 2018. https://doi.org/10.1371/journal.pone.0207908.

[358] Hao Du, Jiazheng Wang, Min Liu, Yaonan Wang, and Erik Meijering. Swinpanet: Swin transformer-based multiscale feature pyramid aggregation network

for medical image segmentation. *IEEE Trans. Neural Netw. Learn. Syst.*, 2022. https://doi.org/10.1109/TNNLS.2022.3204090.

[359] Chaoqun Li, Liejun Wang, and Yongming Li. Transformer and group parallel axial attention co-encoder for medical image segmentation. *Sci Rep*, 12(1):16117, 2022. https://doi.org/10.1038/s41598-022-20440-z.

[360] Xiaomeng Feng, Taiping Wang, Xiaohang Yang, Minfei Zhang, Wanpeng Guo, and Weina Wang. Convwin-unet: Unet-like hierarchical vision transformer combined with convolution for medical image segmentation. *Math. Biosci. Eng.*, 20(1):128–144, 2023. https://doi.org/10.3934/mbe.2023007.

[361] Hubmap - hacking the kidney. https://www.kaggle.com/c/hubmap-kidney-segmentation/data. Accessed 19 February 2023.

[362] Gang Zhang, Chenhong Zheng, Jianfeng He, and Sanli Yi. Pct: Pyramid convolutional transformer for parotid gland tumor segmentation in ultrasound images. *Biomed. Signal Process. Control*, 81:104498, 2023. https://doi.org/10.1016/j.bspc.2022.104498.

[363] Ailiang Lin, Bingzhi Chen, Jiayu Xu, Zheng Zhang, Guangming Lu, and David Zhang. Ds-transunet: Dual swin transformer u-net for medical image segmentation. *IEEE Trans. Instrum. Meas.*, 2022. https://doi.org/10.1109/TIM.2022.3178991.

[364] Zhuotong Cai, Jingmin Xin, Peiwen Shi, Jiayi Wu, and Nanning Zheng. Dstunet: Unet with efficient dense swin transformer pathway for medical image segmentation. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–5. IEEE, 2022. https://doi.org/10.1109/ISBI52829.2022.9761536.

[365] Hongyi Wang, Shiao Xie, Lanfen Lin, Yutaro Iwamoto, Xian-Hua Han, Yen-Wei Chen, and Ruofeng Tong. Mixed transformer u-net for medical image segmentation. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 2390–2394. IEEE, 2022. https://doi.org/10.1109/ICASSP43922.2022.9746172.

[366] Abhinav Sagar. Vitbis: Vision transformer for biomedical image segmentation. In *Clinical Image-Based Procedures, Distributed and Collaborative Learning, Artificial Intelligence for Combating COVID-19 and Secure and Privacy-Preserving Machine Learning*, pages 34–45. Springer, 2021. https://doi.org/10.1007/978-3-030-90874-4_4.

[367] Xuping Huang, Junxi Chen, Mingzhi Chen, Lingna Chen, and Yaping Wan. Tdd-unet: Transformer with double decoder unet for covid-19 lesions segmentation. *Comput. Biol. Med.*, 151:106306, 2022. https://doi.org/10.1016/j.compbiomed.2022.106306.

[368] Kyeong-Beom Park and Jae Yeol Lee. Swine-net: hybrid deep learning approach to novel polyp segmentation using convolutional neural network and swin transformer. *J. Comput. Des. Eng.*, 9(2):616–632, 2022. https://doi.org/10.1093/jcde/qwac018.

[369] Tashvik Dhamija, Anunay Gupta, Shreyansh Gupta, Rahul Katarya, Ghan-

shyam Singh, et al. Semantic segmentation in medical images through transfused convolution and transformer networks. *Appl. Intell.*, pages 1–17, 2022. `https://doi.org/10.1007/s10489-022-03642-w`.

[370] Quan-Dung Pham, Hai Nguyen-Truong, Nam Nguyen Phuong, Khoa NA Nguyen, Chanh DT Nguyen, Trung Bui, and Steven QH Truong. Segtransvae: Hybrid cnn-transformer with regularization for medical image segmentation. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–5. IEEE, 2022. `https://doi.org/10.1109/ISBI52829.2022.9761417`.

[371] Jeya Maria Jose Valanarasu, Rajeev Yasarla, Puyang Wang, Ilker Hacihaliloglu, and Vishal M Patel. Learning to segment brain anatomy from 2d ultrasound with less data. *IEEE J. Sel. Top. Signal Process.*, 14(6):1221–1234, 2020. `https://doi.org/10.1109/JSTSP.2020.3001513`.

[372] Puyang Wang, Nick G Cuccolo, Rachana Tyagi, Ilker Hacihaliloglu, and Vishal M Patel. Automatic real-time cnn-based neonatal brain ventricles segmentation. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 716–719. IEEE, 2018. `https://doi.org/10.1109/ISBI.2018.8363674`.

[373] Zhixian Tang, Jintao Duan, Yanming Sun, Yanan Zeng, Yile Zhang, and Xufeng Yao. A combined deformable model and medical transformer algorithm for medical image segmentation. *Med. Biol. Eng. Comput.*, 61(1):129–137, 2023. `https://doi.org/10.1007/s11517-022-02702-0`.

[374] Tonge image dataset. `https://github.com/BioHit/TongeImageDataset`. Accessed 19 February 2023.

[375] Feiniu Yuan, Zhengxiao Zhang, and Zhijun Fang. An effective cnn and transformer complementary network for medical image segmentation. *Pattern Recognit.*, 136:109228, 2023. `https://doi.org/10.1016/j.patcog.2022.109228`.

[376] Jing Zhang, Qiuge Qin, Qi Ye, and Tong Ruan. St-unet: Swin transformer boosted u-net with cross-layer feature enhancement for medical image segmentation. *Computers in Biology and Medicine*, page 106516, 2023.

[377] Bo Wang, Pengwei Dong, et al. Multiscale transunet++: dense hybrid u-net with transformer for medical image segmentation. *Signal Image Video Process.*, pages 1–8, 2022. `https://doi.org/10.1007/s11760-021-02115-w`.

[378] Geert Litjens, Robert Toth, Wendy van de Ven, Caroline Hoeks, Sjoerd Kerkstra, Bram van Ginneken, Graham Vincent, Gwenael Guillard, Neil Birbeck, Jindang Zhang, et al. Evaluation of prostate segmentation algorithms for mri: the promise12 challenge. *Med. Image Anal.*, 18(2):359–373, 2014. `https://doi.org/10.1016/j.media.2013.12.002`.

[379] Numan Saeed, Ikboljon Sobirov, Roba Al Majzoub, and Mohammad Yaqub. Tmss: An end-to-end transformer-based multimodal network for segmentation and survival prediction. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 319–329. Springer, 2022. `https://doi.org/10.1007/978-3-031-16449-1_31`.

[380] Hecktor 2021. https://www.aicrowd.com/challenges/miccai-2021-hecktor. Accessed 19 February 2023.

[381] Reza Azad, Moein Heidari, Moein Shariatnia, Ehsan Khodapanah Aghdam, Sanaz Karimijafarbigloo, Ehsan Adeli, and Dorit Merhof. Transdeeplab: Convolution-free transformer-based deeplab v3+ for medical image segmentation. In *Proceedings of the Predictive Intelligence in Medicine*, pages 91–102. Springer, 2022. https://doi.org/10.1007/978-3-031-16919-9_9.

[382] Bo Wang, Qian Li, and Zheng You. Self-supervised learning based transformer and convolution hybrid network for one-shot organ segmentation. *Neurocomputing*, 527:1–12, 2023. https://doi.org/10.1016/j.neucom.2022.12.028.

[383] Shuying Xu and Hongyan Quan. Ect-nas: Searching efficient cnn-transformers architecture for medical image segmentation. In *Proceedings of the IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1601–1604. IEEE, 2021. https://doi.org/10.1109/BIBM52615.2021.9669734.

[384] Eli Gibson, Francesco Giganti, Yipeng Hu, Ester Bonmati, Steve Bandula, Kurinchi Gurusamy, Brian Davidson, Stephen P. Pereira, Matthew J. Clarkson, and Dean C. Barratt. Multi-organ abdominal ct reference standard segmentations. https://doi.org/10.5281/zenodo.1169361, 2018. Accessed 5 August 2022.

[385] Jue Jiang, Neelam Tyagi, Kathryn Tringale, Christopher Crane, and Harini Veeraraghavan. Self-supervised 3d anatomy segmentation using self-distilled masked image transformer (smit). In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 556–566. Springer, 2022. https://doi.org/10.1007/978-3-031-16440-8_53.

[386] Shen Jiang, Jinjiang Li, and Zhen Hua. Transformer with progressive sampling for medical cellular image segmentation. *Math. Biosci. Eng.*, 19(12):12104–12126, 2022. https://doi.org/10.3934/mbe.2022563.

[387] Yuanyuan Li, Ziyu Wang, Li Yin, Zhiqin Zhu, Guanqiu Qi, and Yu Liu. X-net: a dual encoding–decoding method in medical image segmentation. *The Visual Computer*, pages 1–11, 2021.

[388] Peter Naylor, Marick Laé, Fabien Reyal, and Thomas Walter. Segmentation of nuclei in histopathology images by deep regression of the distance map. *IEEE transactions on medical imaging*, 38(2):448–459, 2018.

[389] Mashood Mohammad Mohsan, Muhammad Usman Akram, Ghulam Rasool, Norah Saleh Alghamdi, Muhammad Abdullah Aamer Baqai, and Muhammad Abbas. Vision transformer and language model based radiology report generation. *IEEE Access*, 11:1814–1824, 2022. https://doi.org/10.1109/ACCESS.2022.3232719.

[390] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *J. Am. Med. Inf. Assoc.*, 23(2):304–310, 2016. https://doi.org/10.1093/jamia/ocv080.

[391] Hojun Lee, Hyunjun Cho, Jieun Park, Jinyeong Chae, and Jihie Kim. Cross

encoder-decoder transformer with global-local visual extractor for medical image captioning. *Sensors*, 22(4):1429, 2022. `https://doi.org/10.3390/s22041429`.

[392] Benjamin Hou, Georgios Kaissis, Ronald Summers, and Bernhard Kainz. Ratchet: Medical transformer for chest x-ray diagnosis and reporting. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2107.02104`.

[393] AEWP Johnson, Tom Pollard, Roger Mark, Seth Berkowitz, and Steven Horng. Mimic-cxr database. *PhysioNet*, 13026:C2JT1Q, 2019. `https://doi.org/10.13026/C2JT1Q`.

[394] Ming Kong, Zhengxing Huang, Kun Kuang, Qiang Zhu, and Fei Wu. Transq: Transformer-based semantic query for medical report generation. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 610–620. Springer, 2022. `https://doi.org/10.1007/978-3-031-16452-1_58`.

[395] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. *arXiv preprint*, 2019. `https://doi.org/10.48550/arXiv.1901.07042`.

[396] Zhanyu Wang, Hongwei Han, Lei Wang, Xiu Li, and Luping Zhou. Automated radiographic report generation purely on transformer: A multicriteria supervised approach. *IEEE Trans. Med. Imaging*, 41(10):2803–2813, 2022. `https://doi.org/10.1109/TMI.2022.3171661`.

[397] Mingjie Li, Wenjia Cai, Karin Verspoor, Shirui Pan, Xiaodan Liang, and Xiaojun Chang. Cross-modal clinical graph transformer for ophthalmic report generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20656–20665, 2022.

[398] Mingjie Li, Wenjia Cai, Rui Liu, Yuetian Weng, Xiaoyun Zhao, Cong Wang, Xin Chen, Zhong Liu, Caineng Pan, Mengke Li, et al. Ffa-ir: Towards an explainable and reliable medical report generation benchmark. In *Proceedings of the Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, pages 1–14, 2021. `https://doi.org/10.13026/ccbh-z832`.

[399] Yiming Cao, Lizhen Cui, Fuqiang Yu, Lei Zhang, Zhen Li, Ning Liu, and Yonghui Xu. Kdtnet: medical image report generation via knowledge-driven transformer. In *Proceedings of the Database Systems for Advanced Applications*, pages 117–132. Springer, 2022. `https://doi.org/10.1007/978-3-031-00129-1_8`.

[400] Chen Lin, Shuai Zheng, Zhizhe Liu, Youru Li, Zhenfeng Zhu, and Yao Zhao. Sgt: Scene graph-guided transformer for surgical report generation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VII*, pages 507–518. Springer, 2022. `https://doi.org/10.1007/978-3-031-16449-1_48`.

[401] Max Allan, Satoshi Kondo, Sebastian Bodenstedt, Stefan Leger, Rahim Kadkhodamohammadi, Imanol Luengo, Felix Fuentes, Evangello Flouty, Ahmed Mohammed, Marius Pedersen, et al. 2018 robotic scene segmentation challenge. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2001.11190`.

[402] Hoang TN Nguyen, Dong Nie, Taivanbat Badamdorj, Yujie Liu, Lingzi Hong, Jason Truong, and Li Cheng. Eddie-transformer: Enriched disease embedding transformer for x-ray report generation. In *Proceedings of the International Symposium on Biomedical Imaging (ISBI)*, pages 1–5. IEEE, 2022. `https://doi.org/10.1109/ISBI52829.2022.9761459`.

[403] Joseph Paul Cohen, Paul Morrison, and Lan Dao. Covid-19 image data collection. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2003.11597`.

[404] Daniel S Marcus, Tracy H Wang, Jamie Parker, John G Csernansky, John C Morris, and Randy L Buckner. Open access series of imaging studies (oasis): cross-sectional mri data in young, middle aged, nondemented, and demented older adults. *J. Cogn. Neurosci.*, 19(9):1498–1507, 2007. `https://doi.org/10.1162/jocn.2007.19.9.1498`.

[405] Mingrui Ma, Yuanbo Xu, Lei Song, and Guixia Liu. Symmetric transformer-based network for unsupervised image registration. *Knowledge-Based Syst.*, 257:109959, 2022. `https://doi.org/10.1016/j.knosys.2022.109959`.

[406] Junyu Chen, Eric C Frey, Yufan He, William P Segars, Ye Li, and Yong Du. Transmorph: Transformer for unsupervised medical image registration. *Med. Image Anal.*, 82:102615, 2022. `https://doi.org/10.1016/j.media.2022.102615`.

[407] WP Segars, Jason Bond, Jack Frush, Sylvia Hon, Chris Eckersley, Cameron H Williams, Jianqiao Feng, Daniel J Tward, JT Ratnanather, MI Miller, et al. Population of anatomically variable 4d xcat adult phantoms for imaging research and optimization. *Med. Phys.*, 40(4):043701, 2013. `https://doi.org/10.1118/1.4794178`.

[408] Da Hu. Fusing cnns and transformers for deformable medical image registration. In *Proceedings of the International Conference on Computer Science, Electronic Information Engineering and Intelligent Control Technology*, pages 19–23. IEEE, 2022. `https://doi.org/10.1109/CEI57409.2022.9950077`.

[409] Lpba40 dataset. `https://resource.loni.usc.edu/resources/atlases-downloads/`. Accessed 22 February 2023.

[410] Yongpei Zhu and Shi Lu. Swin-voxelmorph: A symmetric unsupervised learning model for deformable medical image registration using swin transformer. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 78–87. Springer, 2022. `https://doi.org/10.1007/978-3-031-16446-0_8`.

[411] Susanne G Mueller, Michael W Weiner, Leon J Thal, Ronald C Petersen, Clifford R Jack, William Jagust, John Q Trojanowski, Arthur W Toga, and Laurel Beckett. Ways toward an early diagnosis in alzheimer's disease: the alzheimer's disease neuroimaging initiative (adni). *Alzheimers. Dement.*,

1(1):55–66, 2005. `https://doi.org/10.1016/j.jalz.2005.06.003`.

[412] Jiacheng Shi, Yuting He, Youyong Kong, Jean-Louis Coatrieux, Huazhong Shu, Guanyu Yang, and Shuo Li. Xmorpher: Full transformer for deformable medical image registration via cross attention. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pages 217–226. Springer, 2022. `https://doi.org/10.1007/978-3-031-16446-0_21`.

[413] Xiahai Zhuang and Juan Shen. Multi-scale patch and multi-modality atlases for whole heart segmentation of mri. *Med. Image Anal.*, 31:77–87, 2016. `https://doi.org/10.1016/j.media.2016.02.006`.

[414] Ramtin Gharleghi, Gihan Samarasinghe, Arcot Sowmya, and Susann Beier. Automated segmentation of coronary arteries. In *Proceedings of the International Conference on Medical Image Computing and Computer Assisted Intervention*, pages 1–13, 2022. `https://doi.org/10.5281/zenodo.3819799`.

[415] Tony CW Mok and Albert Chung. Affine medical image registration with coarse-to-fine vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20835–20844, 2022.

[416] David W Shattuck, Mubeena Mirza, Vitria Adisetiyo, Cornelius Hojatkashani, Georges Salamon, Katherine L Narr, Russell A Poldrack, Robert M Bilder, and Arthur W Toga. Construction of a 3d probabilistic atlas of human cortical structures. *Neuroimage*, 39(3):1064–1080, 2008. `https://doi.org/10.1016/j.neuroimage.2007.09.031`.

[417] Amparo S Betancourt Tarifa, Claudio Marrocco, Mario Molinara, Francesco Tortorella, and Alessandro Bria. Transformer-based mass detection in digital mammograms. *J. Ambient Intell. Humaniz. Comput.*, pages 1–15, 2023. `https://doi.org/10.1007/s12652-023-04517-9`.

[418] Mark D Halling-Brown, Lucy M Warren, Dominic Ward, Emma Lewis, Alistair Mackenzie, Matthew G Wallis, Louise S Wilkinson, Rosalind M Given-Wilson, Rita McAvinchey, and Kenneth C Young. Optimam mammography image database: a large-scale resource of mammography images and clinical data. *Radiology: Artificial Intelligence*, 3(1):e200103, 2020. `https://doi.org/10.1148/ryai.2020200103`.

[419] Bing Leng, Chunqing Wang, Min Leng, Mingfeng Ge, and Wenfei Dong. Deep learning detection network for peripheral blood leukocytes based on improved detection transformer. *Biomed. Signal Process. Control*, 82:104518, 2023. `https://doi.org/10.1016/j.bspc.2022.104518`.

[420] Zahra Mousavi Kouzehkanan, Sepehr Saghari, Sajad Tavakoli, Peyman Rostami, Mohammadjavad Abaszadeh, Farzaneh Mirzadeh, Esmaeil Shahabi Satlsar, Maryam Gheidishahran, Fatemeh Gorgi, Saeed Mohammadi, et al. A large dataset of white blood cells containing cell locations and types, along with segmented nuclei and cytoplasm. *Sci Rep*, 12(1):1123, 2022. `https://doi.org/10.1038/s41598-021-04426-x`.

[421] Ahmad Obeid, Taslim Mahbub, Sajid Javed, Jorge Dias, and Naoufel Werghi. Nucdetr: End-to-end transformer for nucleus detection in histopathology im-

ages. In *Proceedings of the Computational Mathematics Modeling in Cancer Analysis*, pages 47–57. Springer, 2022. `https://doi.org/10.1007/978-3-031-17266-3_5`.

[422] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and Nasir Rajpoot. Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Med. Image Anal.*, 58:101563, 2019. `https://doi.org/10.1016/j.media.2019.101563`.

[423] Jevgenij Gamper, Navid Alemi Koohbanani, Ksenija Benes, Simon Graham, Mostafa Jahanifar, Syed Ali Khurram, Ayesha Azam, Katherine Hewitt, and Nasir Rajpoot. Pannuke dataset extension, insights and baselines. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2003.10778`.

[424] Yifan Zhang, Haoyu Dong, Nicholas Konz, Hanxue Gu, and Maciej A Mazurowski. Lightweight transformer backbone for medical object detection. In *Proceedings of the Cancer Prevention Through Early Detection*, pages 47–56. Springer, 2022. `https://doi.org/10.1007/978-3-031-17979-2_5`.

[425] M Buda, A Saha, R Walsh, S Ghate, N Li, A Święcicki, JY Lo, J Yang, and MA Mazurowski. Data from the breast cancer screening–digital breast tomosynthesis (bcs-dbt). *Data from The Cancer Imaging Archive*, 2020.

[426] Yuntao Shou, Tao Meng, Wei Ai, Canhao Xie, Haiyan Liu, and Yina Wang. Object detection in medical images based on hierarchical transformer and mask mechanism. *Comput. Intell. Neurosci.*, 2022, 2022. `https://doi.org/10.1155/2022/5863782`.

[427] Ke Yan, Xiaosong Wang, Le Lu, and Ronald M Summers. Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning. *J. Med. Imaging*, 5(3):036501–036501, 2018. `https://doi.org/10.1117/1.JMI.5.3.036501`.

[428] Quankai Liu, Guangyuan Zhang, Kefeng Li, Fengyu Zhou, and Dexin Yu. Sfodtrans: semi-supervised fine-grained object detection framework with transformer module. *Med. Biol. Eng. Comput.*, 60(12):3555–3566, 2022. `https://doi.org/10.1007/s11517-022-02682-1`.

[429] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *Int. J. Comput. Vis.*, 111:98–136, 2015. `https://doi.org/10.1007/s11263-014-0733-5`.

[430] Huidong Xie, Stephanie Thorn, Yi-Hwa Liu, Supum Lee, Zhao Liu, Ge Wang, Albert J Sinusas, and Chi Liu. Deep learning based few-angle cardiac spect reconstruction using transformer. *IEEE Trans. Radiat. Plasma Med. Sci.*, 2022. `https://doi.org/10.1109/TRPMS.2022.3187595`.

[431] Yilmaz Korkmaz, Mahmut Yurt, Salman Ul Hassan Dar, Muzaffer Özbey, and Tolga Cukur. Deep mri reconstruction with generative vision transformers. In *Proceedings of the Machine Learning for Medical Image Reconstruction*, pages 54–64. Springer, 2021. `https://doi.org/10.1007/978-3-030-88552-6_6`.

[432] fastmri dataset. `https://fastmri.org/`. Accessed 13 February 2023.

[433] Dayang Wang, Zhan Wu, and Hengyong Yu. Ted-net: Convolution-free t2t vision transformer-based encoder-decoder dilation network for low-dose ct denoising. In *Proceedings of the International Workshop on Machine Learning in Medical Imaging*, pages 416–425. Springer, 2021. `https://doi.org/10.1007/978-3-030-87589-3_43`.

[434] Cynthia H McCollough, Adam C Bartley, Rickey E Carter, Baiyu Chen, Tammy A Drees, Phillip Edwards, David R Holmes III, Alice E Huang, Farhana Khan, Shuai Leng, et al. Low-dose ct for the detection and classification of metastatic liver lesions: results of the 2016 low dose ct grand challenge. *Med. Phys.*, 44(10):e339–e352, 2017. `https://doi.org/10.1002/mp.12345`.

[435] Liutao Yang, Zhongnian Li, Rongjun Ge, Junyong Zhao, Haipeng Si, and Daoqiang Zhang. Low-dose ct denoising via sinogram inner-structure transformer. *IEEE Trans. Med. Imaging*, 2022. `https://doi.org/10.1109/TMI.2022.3219856`.

[436] Taylor R Moen, Baiyu Chen, David R Holmes III, Xinhui Duan, Zhicong Yu, Lifeng Yu, Shuai Leng, Joel G Fletcher, and Cynthia H McCollough. Low-dose ct image and projection dataset. *Med. Phys.*, 48(2):902–911, 2021. `https://doi.org/10.1002/mp.14594`.

[437] Achleshwar Luthra, Harsh Sulakhe, Tanish Mittal, Abhishek Iyer, and Santosh Yadav. Eformer: Edge enhancement based transformer for medical image denoising. *arXiv preprint arXiv:2109.08044*, 2021.

[438] Mario Viti, Hugues Talbot, and Nicolas Gogin. Transformer graph network for coronary plaque localization in ccta. In *Proceedings of the International Symposium on Biomedical Imaging*, pages 1–5. IEEE, 2022. `https://doi.org/10.1109/ISBI52829.2022.9761646`.

[439] Majd Zreik, Robbert W Van Hamersvelt, Jelmer M Wolterink, Tim Leiner, Max A Viergever, and Ivana Išgum. A recurrent cnn for automatic detection and classification of coronary artery plaque and stenosis in coronary ct angiography. *IEEE Trans. Med. Imaging*, 38(7):1588–1598, 2018. `https://doi.org/10.1109/TMI.2018.2883807`.

[440] Onat Dalmaz, Mahmut Yurt, and Tolga Çukur. Resvit: Residual vision transformers for multi-modal medical image synthesis. *arXiv preprint*, 2021. `https://doi.org/10.48550/arXiv.2106.16031`.

[441] Tufve Nyholm, Stina Svensson, Sebastian Andersson, Joakim Jonsson, Maja Sohlin, Christian Gustafsson, Elisabeth Kjellén, Karin Söderström, Per Albertsson, Lennart Blomqvist, et al. Mr and ct data with multiobserver delineations of organs in the pelvic area—part of the gold atlas project. *Med. Phys.*, 45(3):1295–1300, 2018. `https://doi.org/10.1002/mp.12748`.

[442] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.

[443] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[444] Hu Chen, Yi Zhang, Mannudeep K Kalra, Feng Lin, Yang Chen, Peixi Liao, Jiliu Zhou, and Ge Wang. Low-dose ct with a residual encoder-decoder convolutional neural network. *IEEE Trans Med Imaging*, 36(12):2524–2535, 2017. `https://doi.org/10.1109/TMI.2017.2715284`.

[445] Qingsong Yang, Pingkun Yan, Yanbo Zhang, Hengyong Yu, Yongyi Shi, Xuanqin Mou, Mannudeep K Kalra, Yi Zhang, Ling Sun, and Ge Wang. Low-dose ct image denoising using a generative adversarial network with wasserstein distance and perceptual loss. *IEEE Trans Med Imaging*, 37(6):1348–1357, 2018. `https://doi.org/10.1109/TMI.2018.2827462`.

[446] Hongming Shan, Atul Padole, Fatemeh Homayounieh, Uwe Kruger, Ruhani Doda Khera, Chayanin Nitiwarangkul, Mannudeep K Kalra, and Ge Wang. Competitive performance of a modularized deep neural network compared to commercial algorithms for low-dose ct image reconstruction. *Nat. Mach. Intell*, 1(6):269–276, 2019. `https://doi.org/10.1038/s42256-019-0057-9`.

[447] Chunwei Tian, Yong Xu, Zuoyong Li, Wangmeng Zuo, Lunke Fei, and Hong Liu. Attention-guided cnn for image denoising. *Neural Netw.*, 124:117–129, 2020. `https://doi.org/10.1016/j.neunet.2019.12.024`.

[448] Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14318–14328, 2021.

[449] Pedro Silva, Eduardo Luz, Guilherme Silva, Gladston Moreira, Rodrigo Silva, Diego Lucio, and David Menotti. Covid-19 detection in ct images with deep learning: A voting-based scheme and cross-datasets analysis. *Informatics in medicine unlocked*, 20:100427, 2020.

[450] Mohammed A Al-Masni, Mugahed A Al-Antari, Hye Min Park, Na Hyeon Park, and Tae-Seong Kim. A deep learning model integrating frcn and residual convolutional networks for skin lesion segmentation and classification. In *2019 IEEE Eurasia Conference on Biomedical Engineering, Healthcare and Sustainability (ECBIOS)*, pages 95–98. IEEE, 2019.

[451] Alexander Kirillov, Kaiming He, Ross Girshick, and Piotr Dollár. A unified architecture for instance and semantic segmentation, 2017.

[452] Xiangrui Yin, Qianlong Zhao, Jin Liu, Wei Yang, Jian Yang, Guotao Quan, Yang Chen, Huazhong Shu, Limin Luo, and Jean-Louis Coatrieux. Domain progressive 3d residual convolution network to improve low-dose ct imaging. *IEEE transactions on medical imaging*, 38(12):2903–2913, 2019.

[453] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.

[454] Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. Exploring and distilling posterior and prior knowledge for radiology report generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13753–13762, 2021.

[455] Tony CW Mok and Albert Chung. Fast symmetric diffeomorphic image registration with convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4644–4653, 2020.

[456] Mohammad D Alahmadi. Multiscale attention u-net for skin lesion segmentation. *IEEE Access*, 10:59145–59154, 2022.

[457] Qiming Zhang, Yufei Xu, Jing Zhang, and Dacheng Tao. Vitaev2: Vision transformer advanced by exploring inductive bias for image recognition and beyond. *Proceedings of the International Journal of Computer Vision*, pages 1–22, 2023. `https://doi.org/10.1007/s11263-022-01739-w`.

[458] Tianlong Chen, Yu Cheng, Zhe Gan, Lu Yuan, Lei Zhang, and Zhangyang Wang. Chasing sparsity in vision transformers: An end-to-end exploration. *Adv Neural Inf Process Syst*, 34:19974–19988, 2021.

[459] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 558–567, 2021.

[460] Haofei Zhang, Jiarui Duan, Mengqi Xue, Jie Song, Li Sun, and Mingli Song. Bootstrapping vits: Towards liberating vision transformers from pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8944–8953, 2022.

[461] Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3531–3539, 2021.

[462] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint*, 2020. `https://doi.org/10.48550/arXiv.2006.04768`.

[463] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Adv Neural Inf Process Syst*, 33:12104–12114, 2020.

[464] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. `https://doi.org/10.1109/CVPR.2009.5206848`.

[465] Hang Su, Dong Zhao, Hela Elmannai, Ali Asghar Heidari, Sami Bourouis, Zongda Wu, Zhennao Cai, Wenyong Gui, and Mayun Chen. Multilevel threshold image segmentation for covid-19 chest radiography: a framework using horizontal and vertical multiverse optimization. *Computers in Biology and*

*Medicine*, 146:105618, 2022.

[466] Ailiang Qi, Dong Zhao, Fanhua Yu, Ali Asghar Heidari, Zongda Wu, Zhennao Cai, Fayadh Alenezi, Romany F Mansour, Huiling Chen, and Mayun Chen. Directional mutation and crossover boosted ant colony optimization with application to covid-19 x-ray image segmentation. *Computers in biology and medicine*, 148:105810, 2022.

[467] Keli Hu, Liping Zhao, Sheng Feng, Shengdong Zhang, Qianwei Zhou, Xiaozhi Gao, and Yanhui Guo. Colorectal polyp region extraction using saliency detection network with neutrosophic enhancement. *Computers in biology and medicine*, 147:105760, 2022.