

# TRANSOP: TRANSFORMER-BASED MULTIMODAL CLASSIFICATION FOR STROKE TREATMENT OUTCOME PREDICTION

Zeynel A. Samak<sup>1</sup> Philip Clatworthy<sup>2,3</sup> Majid Mirmehdi<sup>1</sup>

<sup>1</sup> Department of Computer Science, University of Bristol, Bristol, UK

<sup>2</sup> Translational Health Sciences, University of Bristol, Bristol, UK

<sup>3</sup> Stroke Neurology, Southmead Hospital, North Bristol NHS Trust, Bristol, UK

## ABSTRACT

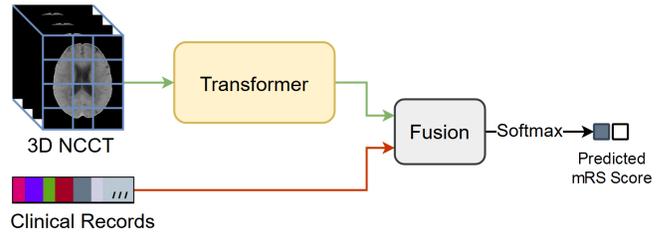
Acute ischaemic stroke, caused by an interruption in blood flow to brain tissue, is a leading cause of disability and mortality worldwide. The selection of patients for the most optimal ischaemic stroke treatment is a crucial step for a successful outcome, as the effect of treatment highly depends on the time to treatment. We propose a transformer-based multimodal network (TranSOP) for a classification approach that employs clinical metadata and imaging information, acquired on hospital admission, to predict the functional outcome of stroke treatment based on the modified Rankin Scale (mRS). This includes a fusion module to efficiently combine 3D non-contrast computed tomography (NCCT) features and clinical information. In comparative experiments using unimodal and multimodal data on the MRCLEAN dataset, we achieve a state-of-the-art AUC score of 0.85.

**Index Terms**— Transformer, Multimodal, Stroke, Ischaemic, NCCT, Outcome.

## 1. INTRODUCTION

Acute ischaemic stroke is the most common type of stroke and a leading cause of disability and mortality worldwide [1]. It is a condition caused by the formation of clots, following interruption of blood flow to the brain. If the blockage is not resolved, the extent of dead tissue increases and the irreversible ischaemic core expands over time. As Saver [2] stated, "Time is brain" for stroke diagnosis and treatment, and it is essential to carry out the appropriate treatment in a timely manner. Although thrombectomy is the most effective treatment for ischaemic stroke cases, there is a risk of brain haemorrhage and death. Therefore, determining if a patient just admitted can benefit from mechanical thrombectomy leading to a good functional outcome, is an important step towards reducing risk and improving the quality of life for stroke patients.

Methods for automatic outcome prediction of stroke treatment have been proposed using logistic regression [3, 4], random forests [5, 6], support vector machines [4, 7], and recently, convolutional neural networks (CNNs) [8, 9, 10]. Some use clinical records [3, 5, 4], imaging information



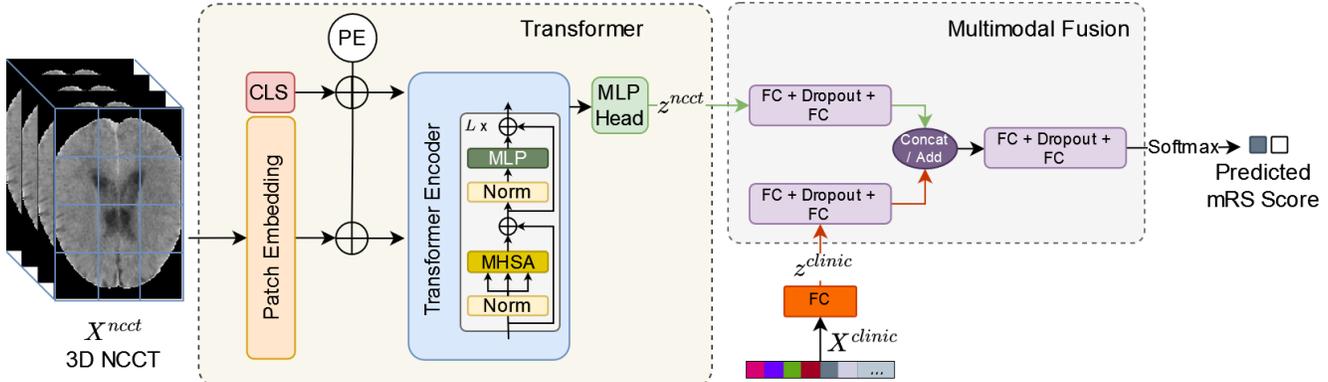
**Fig. 1:** TranSOP predicts functional outcome of ischaemic stroke treatment leveraging only the baseline NCCT scan and clinical records available on hospital admission.

[8, 7, 6], or a combination of both [9, 11, 12]. The CNN-based models have been applied to various imaging modalities, *e.g.* magnetic resonance imaging (MRI), NCCT and CT angiography (CTA). While such deep learning models perform well in medical image analysis, 3D CNN models that exploit 3D brain volumes require numerous parameters and computational resources. Furthermore, they cannot learn long-range relationships due to their limited receptive field. In contrast, more recently, transformers have achieved outstanding results in various applications thanks to their big data and model size scalability and better longer-range attention-based modelling capability [13, 14]. However, pure transformer-based methods have not been widely applied in medical image classification due to their limited performance on small datasets [15].

In this paper, we introduce TranSOP, a transformer-based multimodal architecture to predict functional outcomes of ischaemic stroke patients 90 days after treatment (see Fig 1). We combine clinical metadata (*e.g.* gender, age, hypertension, glucose level) and 3D NCCT obtained at the point of hospital admission for 500 ischaemic stroke patients. We also explore different strategies for this multimodal fusion and conduct extensive experiments on various architectures, including ViT, ViT with CNN, pre-trained ViT (from DeiT [16]) and Swin transformer (SwinT) [17] in our TranSOP model.

## 2. RELATED WORKS

There are only a few studies that have employed CNN-based



**Fig. 2:** Overview of our proposed transformer-based multimodal architecture, TranSOP. PE: positional encoding, CLS: a token/vector that represents the input volume for classification, MHSA: Multi-head self-attention, MLP: multi-layer perceptron, FC: fully connected layer.

multimodal networks to predict the functional outcome of stroke treatments, *e.g.* for thrombolysis [9] and for thrombectomy [11, 10]. Bacchi et al. [9] applied a CNN model to 3D NCCT images and clinical records of patients who underwent thrombolysis treatment. Samak et al. [11] also proposed a multimodal CNN architecture with channel-wise and spatial attentional blocks to predict dichotomised mRS scores from baseline 3D NCCT scans and clinical records of MR CLEAN [18] dataset. Further, in [10], Samak et al. additionally incorporate 1-week follow-up scans during their model training to encode stroke changes over time for better mRS score prediction.

Transformers have shown significant success in natural language processing, *e.g.* machine translation [19], and computer vision, *e.g.* medical imaging tasks [20, 21]. They facilitate a mechanism of self-attention that can model the long-range dependency of sequences and focus on important features. Dosovitskiy et al. [13] proposed the first pure vision transformer (ViT), applied directly to sequences of image patches for image classification. ViTs have obtained comparable and even better results in some tasks than CNNs, *e.g.* for object detection [22, 23].

Since its introduction ViT has been deployed in medical image segmentation using different imaging modalities. UNETR [24] adapts the commonly deployed and successful U-Net architecture [25], by replacing its convolutional encoder with a transformer encoder and modifying its convolutional decoder based on the output of the transformer encoder for image segmentation. Similarly, other studies [26, 21, 27] also replace the convolutional encoder with a transformer encoder, while some integrate the transformer encoder into the bottleneck of a U-Net-like model [28, 29, 30, 31] or use hybrid blocks that combine the convolutional and transformer layers [32, 33]. Such works have been applied to NCCT [28], MRI [29, 32, 21, 33] and microscope [30] images. In another recent work, Amador et al. [34] propose a hybrid model that performs segmentation of the final lesion outcome of ischaemic stroke from baseline spatio-temporal CT perfusion

(CTP) images using a transformer encoder embedded in the U-Net bottleneck.

Although most transformer-based models in medical image analysis are in the *segmentation* domain, there are some studies that have employed them on medical image *classification*, *e.g.* for COVID-19 [35, 36], retinal disease [37], cell analysis [38, 39], brain tumour [20, 40], Alzheimer’s disease [41, 15] classification and age estimation [40, 42]. These methods are based on a pure transformer [35, 43, 44, 45, 46] or a hybrid model that uses ResNet [20, 41, 37], DenseNet [36] or a CNN module [38, 47, 39, 42, 15] followed by a transformer encoder on 2D imaging modalities like X-Rays [46, 48], microscope images [38, 39] and 3D MRI volumes [20, 41, 15, 40]. To the best of our knowledge, there are no studies using the transformer in 3D NCCT classification and prediction of functional stroke outcomes from unimodal or multimodal data.

### 3. PROPOSED METHOD

An overview of the proposed architecture, TranSOP, is shown in Fig. 2, which includes a transformer encoder and a multimodal fusion module to predict mRS scores. Transformers can process 1D input sequences, as originally used in the NLP domain where each word is embedded in a 1D vector as a token. Similarly, we split a 3D NCCT volume,  $X_i^{ncct} \in \mathbb{R}^{1 \times D \times W \times H}$ , into 1D vectors via patch embedding where  $D$ ,  $W$ , and  $H$  are depth, width and height, and a volume is divided into non-overlapping patches of size  $P^3$ , which generate a sequence of 1D patch vectors of length  $L = \lfloor \frac{D}{P} \rfloor \times \lfloor \frac{W}{P} \rfloor \times \lfloor \frac{H}{P} \rfloor$ .

We use a convolutional layer to project each patch into a  $K$  dimensional embedding space [16, 21]. We add a learnable parameter  $[CLS] \in \mathbb{R}^{1 \times K}$ , to the patch embedding sequence to represent the entire volume for classification. In addition, a learnable positional encoding,  $(PE \in \mathbb{R}^{(L+1) \times K})$  is added to the sequences, so that the spatial information of

the patches can be preserved (see Fig. 2). Next, a series of transformer blocks, each including a normalisation layer followed by multi-head self-attention (MHSA), a normalisation layer, and a multi-layer perceptron (MLP) head are utilised in the transformer encoder. Then, an MLP head is applied to the classification token to extract NCCT volume features  $z^{ncct}$  for the fusion process. Clinical metadata features  $z^{clinic}$  are computed by a fully connected layer (FC) (orange box in Fig. 2).

In the multimodal fusion module, a stack of two FCs with a dropout layer in-between prepare the input scan’s  $z^{ncct}$  and  $z^{clinic}$  for fusion (see right box in Fig. 2). We use two methods for the fusion of these image volume and clinical features, (i) concatenation where both features are joined to make a larger 1D vector and (ii) addition where both features are added element-wise with each feature vector multiplied by a learnable weight. Finally, another stack of FC, Dropout, and FC layers is applied to the fused features before being passed to a *Softmax* layer for final predictions. These predictions are dichotomised mRS scores, where  $mRS \leq 2$  indicates a good outcome and  $mRS > 2$  expresses a bad outcome. Note that, dropout layers are deactivated during inference.

#### 4. EXPERIMENTS & RESULTS

**Dataset** – We used the MR CLEAN Trial dataset<sup>1</sup>, collected from a multi-centre study, which is one of the most comprehensive datasets of patients who underwent ischaemic stroke treatment. Five hundred patients (233 assigned to mechanical thrombectomy and 267 to usual care) were treated in 16 medical centres in the Netherlands. We refer the reader to the MR CLEAN study protocol [49, 18] for more detailed information on the dataset.

Through pre-processing, some of the apparent variations due to various acquisition protocols at different clinical centres were reduced to allow our model to deal with more similar standard input. First, all scans were re-sampled to the same voxel size of  $3 \times 1 \times 1 \text{ mm}^3$ , followed by clipping the intensity range of 0-80HU. The skull structure was then removed in the NCCT scans and the volumes were cropped to  $32 \times 192 \times 128$  from the centre.

Data augmentations, such as horizontal/vertical flips and Gaussian noise, were applied to increase the variation and number of input samples to help improve the robustness of the network. Finally, the voxels of the NCCT scans were normalised to zero mean and one standard deviation.

**Implementation Details** – We split the dataset into three subsets, training (70%, 350 patients), validation (15%, 75 patients) and testing (15%, 75 patients). The proposed model was trained for 500 epochs using an Adam optimiser with a weight decay of 0.0001, a learning rate of 0.0003, and a batch size of 24. A cosine learning rate scheduler was used. The

experiments were implemented in PyTorch and MONAI [50] on a single NVIDIA P100 16GB GPU.

**Details of Experiments** – We evaluated the performance of our proposed approach against two existing methods and various transformer architectures that also operate on 3D NCCT volumes and predict the functional outcome of stroke treatment. The methods of Bacchi et al. [9] and Samak et al. [11] which both use imaging and clinical information, were re-trained on the registered MR CLEAN dataset and our data split from scratch. Although, the FeMA [10] model performs a similar task, it additionally uses 1-week follow-up scans that contain information on stroke changes after treatment during model training. Hence, in the interest of direct comparability, we do not include that work in the present evaluation.

We also evaluated our TranSOP approach using different transformer architectures for its encoder part. These are referred to as  $\text{TranSOP}_{ViT}$ ,  $\text{TranSOP}_{DeiT}$ ,  $\text{TranSOP}_{ConViT}$  and  $\text{TranSOP}_{SwinT}$ .  $\text{TranSOP}_{ViT}$  uses the ViT network and is trained from scratch,  $\text{TranSOP}_{DeiT}$  utilises the ImageNet pre-trained DeiT model to demonstrate the effect of transfer learning, and  $\text{TranSOP}_{ConViT}$  uses the first three layers of convolutional blocks before the input is fed into the ViT model to explore the performance of a hybrid model. These three models have the same ViT network which consist of 12 layers of transformer blocks, 12 heads, a hidden MLP feature size of 768 and 3072. In  $\text{TranSOP}_{SwinT}$ , four stages each consisting of two Swin transformer blocks and  $N$  MHSA heads, where  $N = \{3, 6, 12, 24\}$  for each stage respectively, were used. The ClinicDNN model only consumed clinical information to show the expected benefit from imaging information. Note, the multimodal fusion step is the same for all the models.

We evaluated the classification performance of the models with three commonly used metrics, Accuracy, *F1-score* and Area Under ROC Curve (AUC). Table 1 reports the evaluations of the transformer-based and convolution-based networks, along with confidence intervals, for two fusion methods. Broadly, the CNN-based state of the art works [9, 11] outperformed the transformer methods when only imaging information was used, for example, [9] and [11] performed best and second best in accuracy at 0.75 and 0.72 respectively. On the other hand, transformer-based methods exceeded Bacchi et al. [9] and Samak et al. [11] when clinical records were included for multimodal analysis, with the best result obtained by  $\text{TranSOP}_{SwinT}$  at 0.85 AUC. These variations in performance by the transformer could be attributed to both the transformer’s appetite for larger datasets (see [14]), and its already established superiority in handling 1D natural language data. As  $\text{TranSOP}_{SwinT}$  achieved the best AUC score, and it is more efficient thanks to its hierarchical architecture and shifted windowing, it can be a more preferable approach.

The results on the use of fusion methods (concat and addition) are inconclusive and further investigation on more efficient fusion methods is necessary.

<sup>1</sup><https://www.mrclean-trial.org/home.html>

**Table 1:** Results of the models with and without clinical records. The best and second best results are shown in bold and underlined respectively. The second and third rows are convolutional-based models. CI is confidence interval.

Method	w/o Clinical Records			Fusion	with Clinical Records		
	ACC (95% CI)	F1-score (95% CI)	AUC (95% CI)		ACC (95% CI)	F1-score (95% CI)	AUC (95% CI)
ClinicDNN*	-	-	-	-	0.75 (0.65-0.85)	0.44 (0.19-0.64)	0.73 (0.57-0.86)
Samak et al. [11]	<u>0.72</u> (0.62-0.82)	0.33 (0.09-0.53)	0.63 (0.44-0.81)	concat	0.77 (0.66-0.87)	0.47 (0.18-0.67)	0.78 (0.63-0.91)
				add	<u>0.79</u> (0.69-0.89)	0.44 (0.17-0.67)	0.71 (0.51-0.88)
Bacchi et al. [9]	<b>0.75</b> (0.65-0.85)	0.40 (0.16-0.60)	<u>0.66</u> (0.48-0.80)	concat	0.73 (0.62-0.83)	0.51 (0.29-0.68)	0.78 (0.62-0.90)
				add	0.73 (0.62-0.83)	0.51 (0.29-0.68)	0.78 (0.62-0.90)
TranSOP <sub>ConViT</sub>	0.58 (0.46-0.69)	0.40 (0.21-0.56)	<b>0.67</b> (0.46-0.85)	concat	0.77 (0.68-0.87)	<u>0.58</u> (0.36-0.74)	0.83 (0.72-0.93)
				add	0.77 (0.68-0.87)	<u>0.58</u> (0.36-0.74)	0.82 (0.71-0.92)
TranSOP <sub>DeiT</sub>	0.58 (0.46-0.69)	0.40 (0.21-0.56)	0.63 (0.44-0.80)	concat	0.77 (0.68-0.86)	0.53 (0.30-0.71)	0.82 (0.68-0.93)
				add	<u>0.79</u> (0.69-0.89)	0.52 (0.27-0.71)	<u>0.84</u> (0.71-0.94)
TranSOP <sub>ViT</sub>	0.58 (0.46-0.69)	0.40 (0.21-0.56)	0.60 (0.40-0.78)	concat	<b>0.80</b> (0.70-0.89)	0.53 (0.28-0.74)	<u>0.84</u> (0.72-0.94)
				add	<b>0.80</b> (0.70-0.89)	<b>0.59</b> (0.35-0.76)	0.83 (0.71-0.93)
TranSOP <sub>SwinT</sub>	0.58 (0.46-0.69)	0.40 (0.21-0.56)	0.64 (0.44-0.82)	concat	0.76 (0.66-0.86)	0.54 (0.32-0.71)	0.83 (0.71-0.93)
				add	<u>0.79</u> (0.69-0.89)	0.55 (0.31-0.73)	<b>0.85</b> (0.75-0.94)

\* A method that uses only clinical metadata information.

## 5. CONCLUSION

In this work, we investigated the performance of various networks in predicting the functional outcome of ischaemic stroke treatment based on 3D NCCT scans and clinical information, such as age, sex, and demographic data from the patient’s medical history records. Transformer models outperformed convolutional architectures in multimodal settings. This suggests that transformer models, although not performing as well on only imaging data, can learn better complementary imaging information when combined with clinical metadata. In future work, we plan to investigate and explore a data-efficient transformer model for small image datasets. In addition, we would like to extend the proposed architecture to use follow-up scans, such as used in the FeMA [10] method during model training.

## 6. ACKNOWLEDGMENTS

The authors would like to thank the Principal Investigators of the MR CLEAN trial: Profs Aad van der Lugt, Diederik W.J. Dippel, Charles B.L.M. Majoie, Yvo B. W.E.M. Roos, Wim H. van Zwam and Robert J. van Oostenbrugge for providing the data. Z.A. Samak is funded by the Ministry of Education (1416/YLSY), the Republic of Türkiye.

## 7. REFERENCES

- [1] AS Neethi, et al., “Stroke classification from computed tomography scans using 3d convolutional neural network,” *BSPC*, vol. 76, pp. 103720, 2022.
- [2] Jeffrey L Saver, “Time is brain—quantified,” *Stroke*, vol. 37, no. 1, pp. 263–266, 2006.
- [3] Esmee Venema, et al., “Selection of patients for intra-arterial treatment for acute ischaemic stroke: development and validation of a clinical decision tool in two randomised trials,” *BMJ*, vol. 357, 2017.
- [4] Lucas A Ramos, et al., “Predicting poor outcome before endovascular treatment in patients with acute ischemic stroke,” *Front. Neurol.*, vol. 11, pp. 580957, 2020.
- [5] Hendrikus J. A. van Os, et al., “Predicting outcome of endovascular treatment for acute ischemic stroke: Potential value of machine learning algorithms,” *Front. Neurol.*, vol. 9, pp. 784, 2018.
- [6] Jawed Nawabi, et al., “Imaging-based outcome prediction of acute intracerebral hemorrhage,” *TSR*, vol. 12, no. 6, pp. 958–967, Feb. 2021.
- [7] Jeremy Hofmeister, et al., “Clot-based radiomics predict a mechanical thrombectomy strategy for successful recanalization in acute ischemic stroke,” *Stroke*, vol. 51, no. 8, pp. 2488–2494, 2020.
- [8] A Hilbert, et al., “Data-efficient deep learning of radiological image data for outcome prediction after endovascular treatment of patients with acute ischemic stroke,” *CBM*, p. 103516, 2019.
- [9] Stephen Bacchi, et al., “Deep learning in the prediction of ischaemic stroke thrombolysis functional outcomes: A pilot study,” *Academic Radiology*, 4 2019.
- [10] Zeynel A. Samak, et al., “FeMA: Feature matching auto-encoder for predicting ischaemic stroke evolution and treatment outcome,” *CMIG*, vol. 99, pp. 102089, 2022.
- [11] Zeynel A. Samak, et al., “Prediction of thrombectomy functional outcomes using multimodal data,” in *MIUA*, Cham, 2020, pp. 267–279.
- [12] N Kappelhof, et al., “Evolutionary algorithms and decision trees for predicting poor outcome after endovascular treatment for acute ischemic stroke,” *CBM*, vol. 133, pp. 104414, 2021.

- [13] Alexey Dosovitskiy, et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [14] Salman Khan, et al., “Transformers in vision: A survey,” *ACM Comput. Surv.*, 2021.
- [15] Jinseong Jang and Dosik Hwang, “M3T: Three-dimensional medical image classifier using multi-plane and multi-slice transformer,” in *CVPR*, 2022, pp. 20718–20729.
- [16] Hugo Touvron, et al., “Training data-efficient image transformers & distillation through attention,” in *ICML*. PMLR, 2021, pp. 10347–10357.
- [17] Ze Liu, et al., “Swin transformer: Hierarchical vision transformer using shifted windows,” in *CVPR*, 2021, pp. 10012–10022.
- [18] Puck SS Fransen, et al., “MR CLEAN, a multicenter randomized clinical trial of endovascular treatment for acute ischemic stroke in the Netherlands: study protocol for a randomized controlled trial,” *Trials*, vol. 15, no. 1, pp. 343, 2014.
- [19] Ashish Vaswani, et al., “Attention is all you need,” *ANIPS*, vol. 30, 2017.
- [20] Yin Dai, et al., “Transmed: Transformers advance multi-modal medical image classification,” *Diagnostics*, vol. 11, no. 8, pp. 1384, 2021.
- [21] Ali Hatamizadeh, et al., “Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images,” in *MICCAI*, 2022, pp. 272–284.
- [22] Nicolas Carion, et al., “End-to-end object detection with transformers,” in *ECCV*, 2020, pp. 213–229.
- [23] Wen Wang, et al., “FP-DETR: Detection transformer advanced by fully pre-training,” in *ICLR*, 2022.
- [24] Ali Hatamizadeh, et al., “Unetr: Transformers for 3d medical image segmentation,” in *WACV*, 2022, pp. 574–584.
- [25] Özgün Çiçek, et al., “3D U-Net: learning dense volumetric segmentation from sparse annotation,” in *MICCAI*, 2016, pp. 424–432.
- [26] Ali Hatamizadeh, et al., “UNetFormer: A unified vision transformer model and pre-training framework for 3d medical image segmentation,” *arXiv preprint arXiv:2204.00631*, 2022.
- [27] Xuejian Li, et al., “TranSiam: Fusing multimodal visual features using transformer for medical image segmentation,” *arXiv preprint arXiv:2204.12185*, 2022.
- [28] Chun Luo, et al., “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 *EMBC*, 2021, pp. 3565–3568.
- [29] Wenxuan Wang, et al., “Transbts: Multimodal brain tumor segmentation using transformer,” in *MICCAI*, 2021, pp. 109–119.
- [30] Yuanfeng Ji, et al., “Multi-compound transformer for accurate biomedical image segmentation,” in *MICCAI*, 2021, pp. 326–336.
- [31] Jing Wang, et al., “METrans: Multi-encoder transformer for ischemic stroke segmentation,” *Electron. Lett.*, vol. 58, no. 9, pp. 340–342, 2022.
- [32] Yunhe Gao, et al., “UTNet: a hybrid transformer architecture for medical image segmentation,” in *MICCAI*, 2021, pp. 61–71.
- [33] Di Liu, et al., “Transfusion: multi-view divergent fusion for medical image segmentation with transformers,” in *MICCAI*, 2022, pp. 485–495.
- [34] Kimberly Amador, et al., “Hybrid spatio-temporal transformer network for predicting ischemic stroke lesion outcomes from 4d ct perfusion imaging,” in *MICCAI*, 2022, pp. 644–654.
- [35] Arnab Kumar Mondal, et al., “xViTCOS: explainable vision transformer based COVID-19 screening using radiography,” *IEEE JTEHM*, vol. 10, pp. 1–10, 2021.
- [36] Sangjoon Park, et al., “Multi-task vision transformer using low-level chest x-ray feature corpus for covid-19 diagnosis and severity quantification,” *MIA*, vol. 75, pp. 102299, 2022.
- [37] Jianfang Wu, et al., “Vision transformer-based recognition of diabetic retinopathy grade,” *Medical Physics*, vol. 48, no. 12, pp. 7850–7863, 2021.
- [38] Xiyue Wang, et al., “Transpath: Transformer-based self-supervised learning for histopathological image classification,” in *MICCAI*, 2021, pp. 186–195.
- [39] Zeyu Gao, et al., “Instance-based vision transformer for subtyping of papillary renal cell carcinoma in histopathological image,” in *MICCAI*, 2021, pp. 299–308.
- [40] Eunji Jun, et al., “Medical transformer: Universal brain encoder for 3d mri analysis,” *arXiv preprint arXiv:2104.13633*, 2021.
- [41] Chao Li, et al., “Trans-ResNet: Integrating transformers and cnns for alzheimer’s disease classification,” in *ISBI*, 2022, pp. 1–5.
- [42] Sheng He, et al., “Global-local transformer for brain age estimation,” *IEEE TMI*, vol. 41, no. 1, pp. 213–224, 2021.
- [43] Shuang Yu, et al., “Mil-vt: Multiple instance learning enhanced vision transformer for fundus image classification,” in *MICCAI*, 2021, pp. 45–54.
- [44] Christos Matsoukas, et al., “Is it time to replace cnns with transformers for medical images?,” *arXiv preprint arXiv:2108.09038*, 2021.
- [45] Behnaz Gheflati and Hassan Rivaz, “Vision transformers for classification of breast ultrasound images,” in *EMBC*, 2022, pp. 480–483.
- [46] Moinak Bhattacharya, et al., “RadioTransformer: A cascaded global-focal transformer for visual attention-guided disease classification,” *arXiv preprint arXiv:2202.11781*, 2022.
- [47] Anuroop Sriram, et al., “Covid-19 prognosis via self-supervised representation learning and multi-image prediction,” *arXiv preprint arXiv:2101.04909*, 2021.
- [48] Grzegorz Jacenków, et al., “Indication as prior knowledge for multimodal disease classification in chest radiographs with transformers,” in *ISBI*, 2022, pp. 1–5.
- [49] Olvert A. Berkhemer, et al., “A randomized trial of intraarterial treatment for acute ischemic stroke,” *NEJM*, vol. 372, no. 1, pp. 11–20, 2015.
- [50] The MONAI Consortium, “Project monai,” 2020.