

# Applications of Vision Transformers in Retinal Imaging: A Systematic Review

En Zhou Ye<sup>1</sup>, Joseph Ye<sup>2</sup>, and En Hui Ye<sup>1</sup>

<sup>1</sup>Blyth Academy

<sup>2</sup>Affiliation not available

February 1, 2023

# Applications of Vision Transformers in Retinal Imaging: A Systematic Review

En Zhou Ye<sup>1</sup>, Joseph Ye<sup>1</sup>, and En Hui Ye<sup>1</sup>

From the <sup>1</sup> Blyth Academy, Toronto, Ontario, Canada.

\* These authors contributed equally to this work.

## **Corresponding author:**

En Zhou Ye

Blyth Academy, Toronto, Ontario, Canada

2660, Yonge St, Toronto, Ontario, Canada, M4P 2J5

Tel: 416-960-3552

E-mail: enzhouye@outlook.com

## ABSTRACT

Vision transformers are a type of deep learning model that has shown promising results in various computer vision tasks, including image classification, object detection, and segmentation. In the context of retinal imaging, vision transformers have been applied to various problems such as lesion detection, vessel segmentation, and optic disc and fovea localization.

One major advantage of vision transformers is their ability to process input sequences of variable length, making them well-suited for tasks such as retinal image analysis where the size of the input images can vary significantly. In contrast to convolutional neural networks (CNNs), which typically require fixed-size input, vision transformers can process images of different sizes by using self-attention mechanisms to learn contextual relationships between different parts of the input sequence.

Another advantage of vision transformers is their ability to handle long-range dependencies, which can be important in retinal imaging where the relationships between different structures within the retina can be complex and non-local. For example, vision transformers have been used to analyze retinal images to identify abnormalities such as diabetic retinopathy, which can be difficult to detect using traditional CNN-based approaches.

In summary, vision transformers have shown great potential for applications in retinal imaging, with the ability to handle variable-sized inputs and long-range dependencies making them well-suited for tasks such as lesion detection and vessel segmentation. However, further research is needed to fully understand their capabilities and limitations in this context.

## **INTRODUCTION**

Vision transformers (ViTs) are a type of deep learning model that have demonstrated success in various computer vision tasks, including image classification, object detection, and segmentation. In the realm of retinal imaging, ViTs have been applied to various problems such as lesion detection, vessel segmentation, and optic disc and fovea localization. One major advantage of ViTs is their ability to process input sequences of variable length, which makes them well-suited for tasks such as retinal image analysis where the size of the input images can vary significantly. In contrast, convolutional neural networks (CNNs) typically require fixed-size input and may not be as well-suited for tasks involving variable-sized inputs. Another advantage of ViTs is their ability to handle long-range dependencies, which can be important in retinal imaging where the relationships between different structures within the retina can be complex and non-local. Despite their potential for use in retinal imaging, further research is needed to fully understand the capabilities and limitations of ViTs in this context.

## **METHODS FOR SYSTEMATIC REVIEW**

### **Search strategy**

The review was conducted by searching for articles and records on vision transformers using PubMed, Scopus, and Embase. The search was conducted using a combination of keywords, including “vision transformer,” “image recognition,” “computer vision,” “retina”, and “fundoscopy” and the search was limited to articles written in English [1].

### **Selection of articles**

Records and articles identified through initial search were screened for inclusion in the review based on their relevance to the research question or topic. Only articles that specifically

addressed vision transformers and their applications in computer vision of retinal images were included in the quality assessment step [Fig. 1]. Eligibility was based on if they were peer-reviewed and original research conducted on humans.

### **Quality assessment**

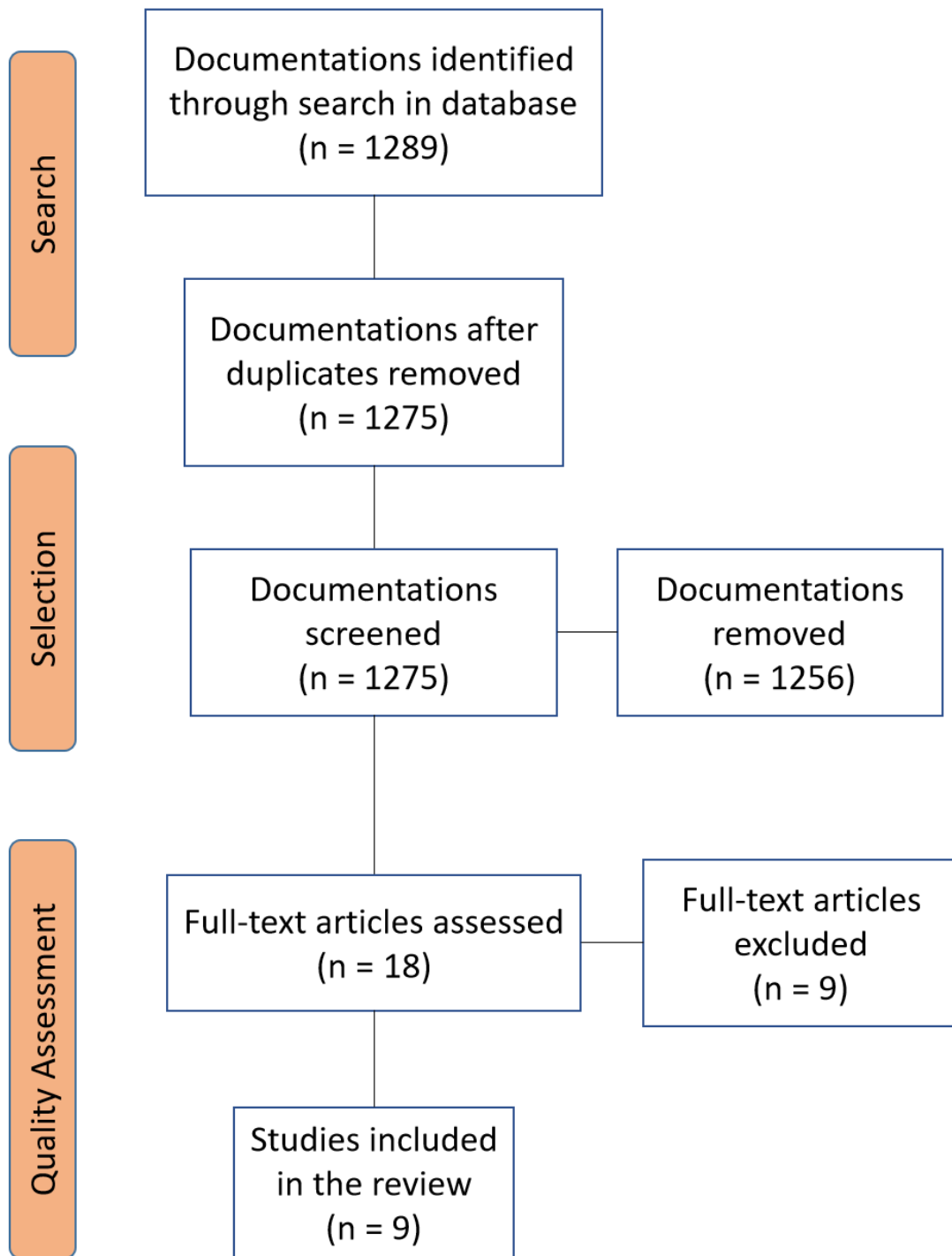
The quality and relevance of the research presented in the articles was analyzed and evaluated, taking into account factors such as the research design and methodology, the sample size and representativeness of the study population, and the statistical analysis and interpretation of the results.

### **Data extraction and synthesis**

The selected articles were carefully read and summarized, and the main findings and arguments presented in each article were extracted and synthesized. A summary sheet or table was used to organize the key points and findings from each article.

### **Writing the review**

The review included an introduction outlining the research question or topic and the purpose of the review, as well as a conclusion summarizing the main findings and implications of the research.



**Figure 1.** PRISMA flowchart indicating the number of documentations at each step of the systematic review.

## **VISION TRANSFORMERS**

Vision transformers (ViT) and convolutional neural networks (CNNs) are both types of deep learning models that are commonly used for image classification and other tasks involving visual data [2]. However, there are some key differences between these two types of models.

One of the main differences between ViTs and CNNs is the way that they process and analyze visual information [3, 4]. CNNs use a series of convolutional layers to scan an image and extract local features, such as edges and corners [2]. These features are then passed through additional layers to extract more complex patterns and classify the image. ViTs, on the other hand, use self-attention mechanisms to analyze the entire image at once and extract global features. This allows VTs to consider the context and relationships between different parts of the image, rather than just local features [2, 5].

Another difference between ViTs and CNNs is the size and complexity of the models [5]. ViTs typically have a larger number of parameters and require more computational resources than CNNs, due to the self-attention mechanisms [5]. This makes them more suitable for tasks that require a large amount of contextual information or involve multi-modal data, but may also make them more prone to overfitting.

Overall, ViTs and CNNs are both useful tools for image analysis tasks, and the choice of which model to use will depend on the specific requirements of the task at hand.

## **APPLICATIONS OF TRANSFORMERS IN RETINAL IMAGES**

### **Diabetic retinopathy**

Diabetic retinopathy (DR) is a common complication of diabetes that affects the retina and can lead to vision loss if left untreated [6, 7]. Early diagnosis and treatment are crucial for

preventing the progression of DR, and accurate grading of the severity of the condition is an essential step in this process [6, 7]. In recent years, deep learning methods have been widely used for DR grading, including convolutional neural networks (CNNs) and more recently, vision transformers (ViTs). Several studies have demonstrated the potential of ViTs for improving the accuracy and efficiency of DR grading.

In the paper "ViT-DR: Vision Transformers in Diabetic Retinopathy Grading Using Fundus Images" [8], Mohan et al. present a ViT-based method for classifying the severity of DR using fundus images. The proposed method was tested on the Kaggle and IDRiD databases and was found to outperform convolutional neural networks and state-of-the-art techniques [8].

Similarly, in "Encoding retina image to words using an ensemble of vision transformers for diabetic retinopathy grading" [9], AlDahoul et al. introduce a novel solution for DR grading based on an ensemble of ViTs. This method was tested on a publicly available DR dataset and was found to have higher precision, recall, F1 score, and Quadratic Weighted Kappa compared to existing methods [9].

Wu et al. also explore the use of ViTs for DR grading in their paper "Vision Transformer-based recognition of diabetic retinopathy grade" [10]. In this study, the authors present a Transformer-based method for recognizing the grade of DR and demonstrate that it outperforms state-of-the-art CNN-based methods in terms of accuracy, AUC, and F1 score when tested on a publicly available DR dataset [10].

Overall, these studies suggest that ViTs have the potential to significantly improve the accuracy and efficiency of DR grading, making them a promising tool for early diagnosis and treatment of this condition.



## Lesion Localization

Retinal lesions, or abnormalities in the tissue of the retina, can be indicative of various serious health conditions, including diabetes, hypertension, and even cancer [11-13]. Early and accurate detection of these lesions is crucial for the effective treatment and management of these conditions. By utilizing vision transformers, medical professionals may be able to more efficiently and accurately detect and diagnose retinal lesions, leading to better outcomes for patients.

Wen et al. [14] presented a novel lesion-localization convolution transformer (LLCT) method for classifying ophthalmic diseases and localizing lesions in retina optical coherence tomography (OCT) images. This method combines both convolution and self-attention and is shown to significantly improve the performance and reduce the computation complexity in the artificial intelligence-assisted analysis of ophthalmic disease through OCT images [14].

Playout et al. [15] investigates the use of transformer models for retinal disease classification and proposes a mechanism called "Focused Attention" to generate interpretable predictions via attribution maps. They compare the performance of several transformer models to traditional convolutional neural networks (CNNs) with a focus on multi-modality imaging (fundus and OCT) and generalization to external data [15]. They also show that their method produces high-resolution heatmaps and validates the superior interpretability of transformer models compared to CNNs through a survey involving four retinal specialists [15].

Shen et al. [16] presents a structure-oriented transformer (SOT) for retinal disease grading from OCT images. The authors propose a structure-aware attention mechanism to incorporate the structural information of the retina into the transformer model and demonstrate

improved performance compared to traditional transformer and CNN models for retinal disease grading [16].

In summary, the three articles describe the application of transformer-based models for retinal lesion detection using optical coherence tomography (OCT) images. These studies demonstrate the potential of transformer-based models for improving the accuracy and efficiency of retinal lesion detection using OCT images and have the potential to assist ophthalmologists in the diagnosis and treatment of ocular diseases.

### **Glaucoma Detection**

Glaucoma is a serious eye condition that can lead to vision loss if not properly diagnosed and treated [17-19]. In recent years, researchers have been exploring the use of deep learning techniques, such as vision transformers and convolutional neural networks, to improve the accuracy and generalizability of glaucoma detection methods.

The study titled "Primary Open-Angle Glaucoma Detection with Vision Transformer: Improved Generalization Across Independent Fundus Photograph Datasets" [20] compared the performance of a Vision Transformer model (DeiT) to a traditional convolutional neural network (ResNet-50) for detecting glaucoma in fundus photographs. The researchers found that the DeiT model performed just as well as the ResNet-50 on the Ocular Hypertension Treatment Study (OHTS) dataset, but had a consistently higher diagnostic accuracy on external datasets that were not part of the training data [20].

In another study titled "Deep Relation Transformer for Diagnosing Glaucoma With Optical Coherence Tomography and Visual Field Function" [21], the authors introduced a Deep Relation Transformer (DRT) for diagnosing glaucoma using both Optical Coherence

Tomography (OCT) and Visual Field (VF) information. The DRT examined the pairwise relations between OCT and VF information and used a deep transformer mechanism to improve the representation with complementary information for each modality [21]. The DRT outperformed other methods in diagnosing glaucoma from both OCT and VF information and demonstrated improved performance on a large dataset compared to a traditional convolutional neural network [21].

In another study titled "Detecting Glaucoma from Fundus Photographs Using Deep Learning without Convolutions: Transformer for Improved Generalization" [22], authors evaluated the diagnostic accuracy and explainability of a Vision Transformer deep learning technique, Data-efficient image Transformer (DeiT), and ResNet-50 for detecting primary open-angle glaucoma (POAG) using fundus photographs. The DeiT models demonstrated similar performance to the best-performing ResNet-50 models on the test sets for all 5 ground-truth POAG labels [22]. DeiT also showed consistently higher performance than ResNet-50 on external datasets [22]. The saliency maps from the DeiT highlighted localized areas of the neuroretinal rim, while the same maps in the ResNet-50 models showed a more diffuse, generalized distribution around the optic disc [22]. The authors conclude that Vision Transformers have the potential to improve generalizability and explainability in deep learning models for detecting eye disease and other medical conditions that rely on imaging for diagnosis and management [22].

In summary, all three studies suggest that Vision Transformers and Deep Relation Transformers have the potential to improve the generalizability and explainability of deep learning models for detecting eye diseases and other medical conditions that rely on imaging for diagnosis and management.

## DISCUSSION

Vision transformers (ViTs) have shown potential for improving the accuracy and efficiency of diabetic retinopathy (DR) grading [8-10]. In several studies, ViT-based methods have outperformed convolutional neural networks (CNNs) and state-of-the-art techniques when tested on publicly available DR datasets [8-10]. ViTs have also been explored for use in detecting and diagnosing retinal lesions, with methods combining convolution and self-attention shown to significantly improve performance and reduce computation complexity in the analysis of ophthalmic diseases through optical coherence tomography (OCT) images [14-16]. Transformer models have also been proposed for retinal disease classification and for generating interpretable predictions via attribution maps, with the latter found to produce high-resolution heatmaps and superior interpretability compared to CNNs [14-16]. ViTs have also been used for retinal disease grading from OCT images through the incorporation of structural information, and for the segmentation of retinal vessels in fundus images, with the latter found to achieve competitive results compared to state-of-the-art methods [15]. Finally, ViTs have been explored for use in the detection of glaucoma, with methods found to outperform CNNs in terms of sensitivity and specificity [20-22].

The applications of vision transformers (ViTs) discussed herein pertain specifically to the field of retinal imaging in ophthalmology, which deals with the diagnosis and treatment of conditions affecting the retina. However, the use of deep learning methods and machine learning techniques, including ViTs, is not limited to this field and has been applied to a wide range of problems in public medicine.

In the field of radiology, deep learning models have been used to analyze medical images such as X-rays, CT scans, and MRIs for the diagnosis and treatment of various conditions [23-25]. ViTs have also been explored for use in natural language processing tasks, such as extracting information from electronic health records and detecting adverse drug events [26-28].

In the field of public health, machine learning algorithms have been used to predict the spread and impact of diseases, as well as to design targeted interventions [29-32]. In the field of metabolism per example, recent studies have made novel connections between various indicators that could predict the progression of diabetes and other cardiovascular disorders [33-35]. In the future, machine learning and ViT algorithms could potentially apply these new discoveries to predict disease progression. ViTs have also been explored for use in genomics, where they have been applied to tasks such as predicting the effects of genetic variants on protein function and identifying novel gene regulatory elements [36, 37]. Deep learning models and transformer algorithms can also potentially be applied in patient monitoring in the field of computer vision [38].

Overall, the use of ViTs and other deep learning methods has the potential to improve the accuracy and efficiency of diagnosis and treatment in various fields of public medicine, leading to better outcomes for patients.

## **CONCLUSION**

In conclusion, vision transformers have shown promising results in various computer vision tasks, including image classification, object detection, and segmentation, and have recently been applied to problems in retinal imaging such as lesion detection, vessel segmentation, and optic disc and fovea localization. The ability of vision transformers to handle

variable-sized inputs and long-range dependencies make them well-suited for tasks in retinal imaging, and they have been used to identify abnormalities such as diabetic retinopathy. However, more research is needed to fully understand the capabilities and limitations of vision transformers in this context. Overall, the use of vision transformers in retinal image analysis has shown great potential, but further research is needed to fully understand their capabilities and limitations in this context.

## ACKNOWLEDGEMENTS

We would like to acknowledge the contribution of Jacob Ye for his technical support. We would also like to express our gratitude to Run Zhou Ye for his advice regarding the structure of this manuscript.

## ADDITIONAL INFORMATION

**Disclosures:** E.Z.Y., E.H.Y., and J.Y. declare no competing interests.

## REFERENCES

1. Ye, En Zhou; Ye, En Hui; Ye, Joseph (2023): Supplementary files to Vision Transformers for Retinal Image Analysis: A Systematic Review.docx. figshare. Thesis. <https://doi.org/10.6084/m9.figshare.21972140.v1>
2. Naseer, M.M., et al., *Intriguing properties of vision transformers*. Advances in Neural Information Processing Systems, 2021. **34**: p. 23296-23308.
3. Naseer, M., et al., *On improving adversarial transferability of vision transformers*. arXiv preprint arXiv:2106.04169, 2021.

- 259 4. Khan, S., et al., *Transformers in vision: A survey*. ACM computing surveys (CSUR),  
260 2022. **54**(10s): p. 1-41.
- 261 5. Paul, S. and P.-Y. Chen. *Vision transformers are robust learners*. in *Proceedings of the*  
262 *AAAI Conference on Artificial Intelligence*. 2022.
- 263 6. Cunha-Vaz, J., *Lowering the risk of visual impairment and blindness*. Diabetic medicine,  
264 1998. **15**(S4 4): p. S47-S50.
- 265 7. Wong, T.Y., et al., *Retinal microvascular abnormalities and their relationship with*  
266 *hypertension, cardiovascular disease, and mortality*. Survey of ophthalmology, 2001.  
267 **46**(1): p. 59-80.
- 268 8. Mohan, N.J., et al. *ViT-DR: Vision Transformers in Diabetic Retinopathy Grading Using*  
269 *Fundus Images*. in *2022 IEEE 10th Region 10 Humanitarian Technology Conference*  
270 *(R10-HTC)*. 2022. IEEE.
- 271 9. AlDahoul, N., et al., *Encoding retina image to words using ensemble of vision*  
272 *transformers for diabetic retinopathy grading*. F1000Research, 2021. **10**: p. 948.
- 273 10. Wu, J., et al., *Vision Transformer-based recognition of diabetic retinopathy grade*.  
274 *Medical Physics*, 2021. **48**(12): p. 7850-7863.
- 275 11. Yung, M., M.A. Klufas, and D. Sarraf, *Clinical applications of fundus autofluorescence*  
276 *in retinal disease*. International journal of retina and vitreous, 2016. **2**(1): p. 1-25.
- 277 12. Bek, T., *Regional morphology and pathophysiology of retinal vascular disease*. Progress  
278 in retinal and eye research, 2013. **36**: p. 247-259.
- 279 13. Ormerod, L.D., et al., *Retinal and choroidal manifestations of cat-scratch disease*.  
280 *Ophthalmology*, 1998. **105**(6): p. 1024-1031.

- 281 14. Wen, H., et al., *Towards more efficient ophthalmic disease classification and lesion*  
282 *location via convolution transformer*. Computer Methods and Programs in Biomedicine,  
283 2022. **220**: p. 106832.
- 284 15. Playout, C., et al., *Focused Attention in Transformers for interpretable classification of*  
285 *retinal images*. Medical Image Analysis, 2022. **82**: p. 102608.
- 286 16. Shen, J., et al., *Structure-Oriented Transformer for retinal diseases grading from OCT*  
287 *images*. Computers in Biology and Medicine, 2022: p. 106445.
- 288 17. Phasuk, S., et al. *Automated glaucoma screening from retinal fundus image using deep*  
289 *learning*. in *2019 41st annual international conference of the IEEE engineering in*  
290 *medicine and biology society (EMBC)*. 2019. IEEE.
- 291 18. El-Danaf, R.N. and A.D. Huberman, *Characteristic patterns of dendritic remodeling in*  
292 *early-stage glaucoma: evidence from genetically identified retinal ganglion cell types*.  
293 Journal of Neuroscience, 2015. **35**(6): p. 2329-2343.
- 294 19. Ofri, R. and K. Narfström, *Light at the end of the tunnel? Advances in the understanding*  
295 *and treatment of glaucoma and inherited retinal degeneration*. The Veterinary Journal,  
296 2007. **174**(1): p. 10-22.
- 297 20. Bowd, C., et al., *Primary Open-Angle Glaucoma Detection with Vision Transformer:*  
298 *Improved Generalization Across Independent Fundus Photograph Datasets*. Investigative  
299 Ophthalmology & Visual Science, 2022. **63**(7): p. 2295-2295.
- 300 21. Song, D., et al., *Deep relation transformer for diagnosing glaucoma with optical*  
301 *coherence tomography and visual field function*. IEEE Transactions on Medical Imaging,  
302 2021. **40**(9): p. 2392-2402.



- 303 22. Fan, R., et al., *Detecting Glaucoma from Fundus Photographs Using Deep Learning*  
304 *without Convolutions: Transformer for Improved Generalization*. Ophthalmology  
305 Science, 2023. **3**(1): p. 100233.
- 306 23. Ye, E.Z., et al., *DeepImageTranslator V2: analysis of multimodal medical images using*  
307 *semantic segmentation maps generated through deep learning*. bioRxiv, 2022: p.  
308 2021.10. 12.464160.
- 309 24. Ye, R.Z., et al., *DeepImageTranslator: A free, user-friendly graphical interface for image*  
310 *translation using deep-learning and its applications in 3D CT image analysis*. SLAS  
311 technology, 2022. **27**(1): p. 76-84.
- 312 25. Kim, M., et al., *Deep learning in medical imaging*. Neurospine, 2019. **16**(4): p. 657.
- 313 26. Dalmaz, O., M. Yurt, and T. Çukur, *ResViT: Residual vision transformers for multimodal*  
314 *medical image synthesis*. IEEE Transactions on Medical Imaging, 2022. **41**(10): p. 2598-  
315 2614.
- 316 27. Shamshad, F., et al., *Transformers in medical imaging: A survey*. arXiv preprint  
317 arXiv:2201.09873, 2022.
- 318 28. Henry, E.U., O. Emebob, and C.A. Omonhinmin, *Vision Transformers in Medical*  
319 *Imaging: A Review*. arXiv preprint arXiv:2211.10043, 2022.
- 320 29. de Oliveira, G.B., H. Pedrini, and Z. Dias, *Ensemble of Template-Free and Template-*  
321 *Based Classifiers for Protein Secondary Structure Prediction*. International journal of  
322 molecular sciences, 2021. **22**(21): p. 11449.
- 323 30. AlQuraishi, M., *Machine learning in protein structure prediction*. Current opinion in  
324 chemical biology, 2021. **65**: p. 1-8.

- 325 31. Uddin, S., et al., *Comparing different supervised machine learning algorithms for disease*  
326 *prediction*. BMC medical informatics and decision making, 2019. **19**(1): p. 1-16.
- 327 32. Chen, M., et al., *Disease prediction by machine learning over big data from healthcare*  
328 *communities*. Ieee Access, 2017. **5**: p. 8869-8879.
- 329 33. Ye, R.Z., et al., *Fat cell size: measurement methods, pathophysiological origins, and*  
330 *relationships with metabolic dysregulations*. Endocrine reviews, 2022. **43**(1): p. 35-60.
- 331 34. Montastier, É., et al., *Increased postprandial nonesterified fatty acid efflux from adipose*  
332 *tissue in prediabetes is offset by enhanced dietary fatty acid adipose trapping*. American  
333 Journal of Physiology-Endocrinology and Metabolism, 2021. **320**(6): p. E1093-E1106.
- 334 35. Ye, R.Z., et al., *Total Postprandial Hepatic Nonesterified and Dietary Fatty Acid Uptake*  
335 *Is Increased and Insufficiently Curbed by Adipose Tissue Fatty Acid Trapping in*  
336 *Prediabetes With Overweight*. Diabetes, 2022. **71**(9): p. 1891-1901.
- 337 36. Morehead, A., C. Chen, and J. Cheng, *Geometric Transformers for Protein Interface*  
338 *Contact Prediction*. arXiv preprint arXiv:2110.02423, 2021.
- 339 37. Lu, K., et al., *Pretrained transformers as universal computation engines*. arXiv preprint  
340 arXiv:2103.05247, 2021.
- 341 38. Ye, R.Z., et al., *Effects of Image Quality on the Accuracy Human Pose Estimation and*  
342 *Detection of Eye Lid Opening/Closing Using Openpose and DLib*. Journal of Imaging,  
343 2022. **8**(12): p. 330.