

# LEVERAGING MULTIMODAL FUSION FOR ENHANCED DIAGNOSIS OF MULTIPLE RETINAL DISEASES IN ULTRA-WIDE OCTA

Hao Wei<sup>1,\*</sup> Peilun Shi<sup>1,\*</sup> Guitao Bai<sup>2</sup> Minqing Zhang<sup>1</sup> Shuangli Li<sup>2,†</sup> Wu Yuan<sup>1,†</sup>

<sup>1</sup> Department of Biomedical Engineering, The Chinese University of Hong Kong, Hong Kong SAR

<sup>2</sup> Department of Ophthalmology, Zigong First People’s Hospital, Zigong, China

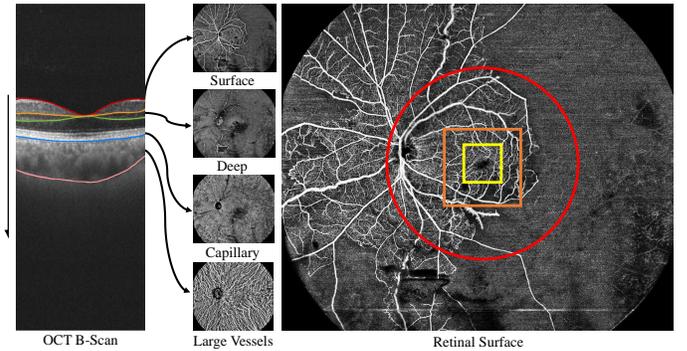
## ABSTRACT

Ultra-wide optical coherence tomography angiography (UW-OCTA) is an emerging imaging technique that offers significant advantages over traditional OCTA by providing an exceptionally wide scanning range of up to  $24 \times 20 \text{ mm}^2$ , covering both the anterior and posterior regions of the retina. However, the currently accessible UW-OCTA datasets suffer from limited comprehensive hierarchical information and corresponding disease annotations. To address this limitation, we have curated the pioneering M3OCTA dataset, which is the first multimodal (i.e., multilayer), multi-disease, and widest field-of-view UW-OCTA dataset. Furthermore, the effective utilization of multi-layer ultra-wide ocular vasculature information from UW-OCTA remains underdeveloped. To tackle this challenge, we propose the first cross-modal fusion framework that leverages multi-modal information for diagnosing multiple diseases. Through extensive experiments conducted on our openly available M3OCTA dataset, we demonstrate the effectiveness and superior performance of our method, both in fixed and varying modalities settings. The construction of the M3OCTA dataset, the first multimodal OCTA dataset encompassing multiple diseases, aims to advance research in the ophthalmic image analysis community.

**Index Terms**— Ultra-wide Optical coherence tomography angiography, Open access dataset

## 1. INTRODUCTION

A thorough assessment of ocular vasculature is crucial for ophthalmologists to evaluate eye health. Angiography is the most effective modality for this purpose. Fundus fluorescein angiography (FFA) provides clear visualization of vessels with a wide field of view (FOV) but is limited by its invasive nature. Conventional Optical coherence tomography angiography (OCTA) overcomes the invasiveness of FFA but has a limited FOV [1, 2]. Ultra-wide OCTA (UW-OCTA) combines the advantages of a wide FOV and non-invasiveness, making it the optimal angiographic technique. It enables comprehensive and non-invasive visualization of ocular blood vessels,



**Fig. 1.** Illustration of proposed M3OCTA Dataset. The selected four-modal sample, scanned in  $24 \times 20 \text{ mm}^2$ , contains four layer projection maps: retinal surface (inner limiting membrane-the inner plexiform layer), retinal deep (the inner plexiform layer-the outer plexiform layer), choroid capillary (bruch’s membrane) and choroid large vessel (choroid layer) (The vertical arrow denotes the projection direction). The red, orange, and yellow regions denote the scan region of the regular fundus,  $3 \times 3 \text{ mm}$ , and  $6 \times 6 \text{ mm}$  OCTA image.

providing valuable information for diagnostic and therapeutic decisions related to ophthalmic diseases.

More recently, a new type of UW-OCTA technique (BMizar 400KHz Full-Range SS-OCT, TowardPi. Inc) affords a widest imaging range up to  $24 \times 20 \text{ mm}^2$  [3], as shown in Fig.1. It has been demonstrated for various clinical applications such as disease diagnosis and screening [4, 5]. Projections from the three-dimensional (3D) volume in different layers show the multiple-layer visualization of both the retina and choroid, as shown in Fig.1. Considering that ocular pathologies can manifest as abnormalities across different layers of the retina, it is crucial to take into account lesions on each affected layer in order to facilitate a multi-modality-based diagnosis of diseases. Therefore, it is of great importance to develop and validate methods that effectively utilize the multi-modality information. Recently, the convolutional neural network (CNN) based two-branch paradigm [6] or the hybrid manner of CNN with Transformer [7] achieved comparable performance in multi-modal disease diagnosis.

\*Equal Contribution; <https://github.com/hwei-hw/M3OCTA>

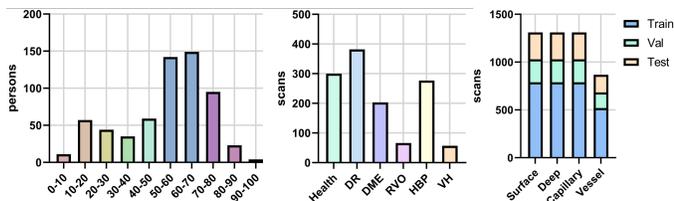
†Corresponding Author: [wuyuan@cuhk.edu.hk](mailto:wuyuan@cuhk.edu.hk), [985750247@qq.com](mailto:985750247@qq.com)

However, the parameters of these methods usually increase multiplicatively with the number of involved modalities, and the trained model cannot dynamically suit the varying modalities’ inference, limiting practical usage. Moreover, the existing public OCTA datasets are limited by either the narrow scan range or the single modality.

To tackle these challenges, we compile the first retinal UW-OCTA dataset of widest FOV and containing five diseases and four modalities, termed as M3OCTA, and further propose a novel cross-modal fusion framework (CMF-Net) to leverage multi-modal information in this dataset for multiple diseases diagnosis in multi-label setting. Specifically, the unlabeled images and the whole train set are utilized to pre-train the ViT-based encoder by the multi-modal based masked image modeling, which learns relationships between any two modalities by the global self-attention. Then, we propose the attention-based cross-modal fusion (CMF) block to reinforce and extract the multi-modal semantics for the disease diagnosis. Moreover, our design enables a varying number of modal inputs but without performance drop during the inference stage, increasing suitability and compatibility in clinical use. In summary, our main contributions are as follows:

- We introduce the first multi-modal UW-OCTA dataset with multiple disease annotations, i.e., M3OCTA, aiming to promote advances in ophthalmic image analysis.
- We propose the cross-modal fusion network (CMF-Net) to involve multi-layer OCTA images for accurate and robust diagnosis of multiple retinal diseases. Its support for dynamic input modalities during inference greatly improves the suitability for the practical use.
- Extensive experiments have verified the effectiveness and superiority of the proposed method. Investigation of the impact on varying input modalities also proves the value of the multiple modalities of the proposed dataset.

## 2. DATASET



**Fig. 2.** Age, diseases and modalities statistics of M3OCTA dataset.

Our proposed M3OCTA is the first multi-modal based ultra-wide retinal OCTA dataset, involving 1637 scans from 1046 eyes of 620 individuals imaged in Zigong First People’s Hospital through  $24 \times 20$  scan mode. Specifically, 1067 scans contains choroid large vessel image; images of 1310 scans from 496 people are labeled as six classes in multi-label setting, including healthy, diabetic retinopathy (DR), diabetic

**Table 1.** Summary of the public OCTA datasets and ours, where  $R$  and  $C$  denotes the retina and choroid

Dataset	Modalities	Subjects	Diseases	Regions	Resolution	FOV(mm)
Giarrarano et al. [8]	1	11	1	R	$91 \times 91$	$3 \times 3$
ROSE [9]	1	151	-	R	$304 \times 304$ $512 \times 512$	$3 \times 3$
OCTAGON [10]	1	213	2	R	$320 \times 320$	$6 \times 6$ $3 \times 3$
FAZID [11]	1	304	3	R	$420 \times 420$	$6 \times 6$
OCTA-500 [12]	2	500	>12	R	$640 \times 400$	$6 \times 6$ $3 \times 3$
DRAC [13]	1	<611	1	R	$1024 \times 1024$	$12 \times 12$
<b>M3OCTA</b>	<b>4</b>	<b>1067</b>	<b>5</b>	<b>R &amp; C</b>	<b><math>1536 \times 1280</math></b>	<b><math>24 \times 20</math></b>

macular edema (DME), Retinal Vein Occlusion (RVO), Hypertension (HBP) and Vitreous Hemorrhage (VH), and then split into train, validation and test set as 6:2:2. The remaining unlabeled data are only used in the pretraining step. Details of our M3OCTA and other public ones are listed in Table.1. Compared with others, M3OCTA dataset demonstrates superiorities in several aspects including the number of modalities, number of patients, image resolution, and FOV.

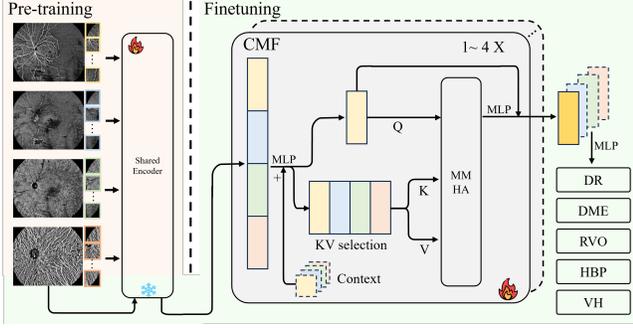
## 3. METHOD

We design a novel framework that leverages the pre-trained transformer encoder and fuses multi-modal information for diagnostics, as shown in Fig. 3. The proposed framework is divided into two stages: transformer encoder pretraining and multi-modal fusion.

### 3.1. Multi-modal Encoder Pretraining

This stage follows the paradigm of masked image modeling [14]: masks the patch of input images are used to train an encoder to recover the patch for better multi-modal feature representations. Specifically, we first build the pretraining set to include all the unlabeled samples and the whole train set. Then, the multi-modal UW-OCTA images are split into image patches, and 25% of these patches are randomly selected as the inputs while the remaining ones would be utilized as the targets to train the encoder [15]. During this selection process, we adopt the Dirichlet distribution to assign the number of input patches in each modality. In addition, for samples with missing modalities, we introduce dynamic controls after Dirichlet sampling to set the selected proportion of missing modalities as zeros and also increase the sampled patches in other modal images to achieve the fixed number of input patches.

Following MAE [15], four decoders are designed to reconstruct four modalities, where each decoder includes a linear projection layer to reduce the dimensions, positional embedding, and transformer blocks. For those input samples with missing modality, the missing one will not contribute to the loss of reconstruction. Through this reconstruction process, the encoder can model global and long-range interactions between any two modalities, which would be particularly beneficial in the multi-modal based diagnosis.



**Fig. 3.** Our proposed framework for the cross-modal fusion and multiple retinal diseases diagnosis, including two steps: 1) encoder pretraining and 2) decoder finetuning. We first pre-train the vision transformer through the masked image modeling learning paradigm to learn multi-modal feature presentations. Then, freezing the encoder and adopting the proposed cross-modal fusion (CMF) block for each modality (share weights) to fuse the multi-modal features and extract the high-level semantics for the following diagnosis. During inference, our method can dynamically process inputs of varying modal numbers and provide stable and comparable performances.

### 3.2. Cross-modal Fusion Decoder

To reinforce the multi-modal representations and extract high-level semantics, we design the cross-modal fusion (CMF) block, as shown in Fig.3, and then parallel (but share weights) this block for each modality to build the diagnosis decoder. Specifically, a linear projection layer is incorporated to tailor the encoder outputs to the decoder’s dimension. Following this projection, the decoder inputs are enhanced by the addition of both sin-cosine positional embeddings and modality-specific embeddings (context) that have been previously learned. Subsequently, the process continues with a masked multi-head attention (MMHA) layer, a Multilayer Perceptron (MLP), and the final classifier.

Based on the mentioned process, for the  $M$  modal encoded features  $z_0$  from the encoder:

$$z_0 = [CLS; x_0^0; x_1^0; \dots, x_P^0; x_1^1; \dots; x_P^M], P = 196 \quad (1)$$

$$z_{proj} = proj(z_0) + E_{pos} + [E^0; \dots E^M] \quad (2)$$

$$Q^i = [x_0^i; x_1^i; \dots, x_P^i], x_j^i \in z_{proj}, i = 0, \dots, M \quad (3)$$

$$KV = \sum_{i=0}^M Q^i \quad (4)$$

$$z^i = MMHA(Q^i, KV, mask_i), i = 0, \dots, M \quad (5)$$

$$z_{fuse}^i = z^i + MLP(z^i) + Q^i, i = 0, \dots, M \quad (6)$$

where  $E_{pos}$  and  $E^i$  denote the sin-cosine positional embeddings and modality  $i$  embeddings. To learn the relationships between any modalities, We utilize the tokens of each modality as the *key* and the summation of all that of each modality

as the *key-value* pair and then apply the multi-head attention to interact with each other, where other implementations about *key-value* pair are discussed in the following ablation study.  $mask_i \in \{0, 1\}$  represents whether the modality  $i$  in each input sample is missing (0) or not (1). The corresponding position of the attention computation in the MMHA layer will be set to negative infinity before the *softmax* operation. Finally, all fused features  $z_{fuse}$  are concatenated together as the inputs of the classifier.

## 4. EXPERIMENT

In order to investigate the capabilities of our proposed method, we performed extensive experiments utilizing the datasets with multimodal fusion algorithms. Furthermore, we evaluated foundation models trained on natural and retinal-related images, as feature extractors in an encoder role.

### 4.1. Implement Details and Evaluation Metrics

The proposed method was implemented in Pytorch using a ViT-B [16] with  $16 \times 16$  pixel patches as the backbone encoder, and input images with the size of  $224 \times 224$ . During pretraining, random cropping and random horizontal flipping are used as the data augmentations, with the Cosine Anneal scheduler and AdamW [17] optimizer ( $lr = 0.001$ , warm up= 40 with total = 1600 epochs, batch size = 160) and MSE loss. For the finetuning stages, the decoder was trained using an initial learning rate of  $2.5e - 4$ , binary cross entropy loss, with the AdamW optimizer for 100 epochs and batch size = 128. The whole approach was implemented and trained using Pytorch on NVIDIA 4090 GPUs.

Accuracy (ACC), Area Under the Receiver Operating Characteristic (AUROC), Average-precision (AP), and F1-score (F1) are adopted in the multi-label setting to evaluate all involved experiments by using the TorchMetrics library [18].

### 4.2. Comparison

To evaluate and compare the performance and effectiveness of the proposed framework, we select the most commonly used classification network: ResNet-18 [19], ConvNeXT[20] and multi-modal image classification network in the medical domain: TFormer [7], MMC[6]. All methods are trained and tested on all four modalities except the TFormer and MMC, which natively and only support two-modal images (retinal surface and deep layer images are used in our experiments). Table 2 lists all results of compared methods, where the lower boundary of theoretical performance by random predictions is shown in the first row. Compared with other approaches, our network achieves state-of-the-art (SOTA) performance across all evaluation metrics in four-modal and two-modal settings.

To validate the performance gain brought by the multiple modalities, we pre-train another encoder on only one modal-

**Table 2.** The classification results comparison between our method and others. Where \* denotes method was trained and tested on two modalities (retinal surface and deep images).

Method	ACC	AUROC	AP	F1
Random	50.39	48.90	19.43	25.56
Resnet[19]	81.85	82.17	49.91	37.65
Convnext-Tiny[20]	83.20	83.68	54.65	50.88
MMC*[6]	84.27	84.26	50.69	46.21
TFormer*[7]	84.41	69.11	37.21	50.48
Ours*	83.49	84.88	56.44	51.52
Ours	<b>84.98</b>	<b>84.77</b>	<b>59.99</b>	<b>59.60</b>

ity, i.e., the retinal surface image, and then finetune the decoder on these two encoders by 1 (retinal surface), 2 (retinal surface + deep), 3 (retinal surface + deep + choroid Capillary), or all four modalities. From the results in Table.3, the 4-modal pretrained encoder outperforms the 1-modal encoder in all 4 scenarios with an average improvement of 10.9, which demonstrates the effectiveness of additional modalities. For the 4-modal encoder-based experiments, the performance gain is observed as the number of input modalities increases. More comprehensive details brought by more modalities boost the performance of our model. However, the opposite trend is found in another group, which may result from the weak multi-modal features of the encoder. In summary, we determined that our method attains an optimal solution as the number of input modalities increases and multi-modal images can boost diagnosis performance.

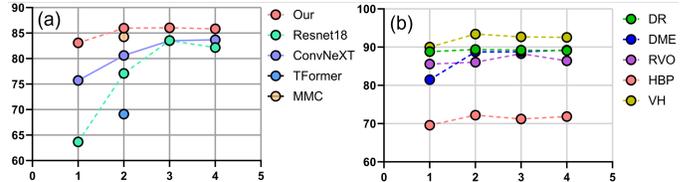
**Table 3.** Influence of modality variation on performance (F1 Score) during pretraining and finetuning stages

Modality Number	1	2	3	4
Pretrained on 1 Modality	46.48	44.85	44.57	41.72
Pretrained on 4 Modality	<b>52.73</b>	<b>51.52</b>	<b>57.37</b>	<b>59.60</b>

We subjected our approach (which incorporates a 4-modal fine-tuned decoder with a 4-modal pre-trained encoder) to an evaluation across diverse modalities for validating stability when facing varying modalities shown in Fig.4(a). With the number of input modalities increases, the performance gain is also observed in almost all comparable methods, where our model shows more stability than others, demonstrating the effectiveness of learned multi-modal features in our design. Further, this stability also can be observed in all five disease types in Fig.4(b).

### 4.3. Ablation Study

We also conduct ablation experiments to demonstrate the influence of pretraining datasets. Specifically, we alternately replaced the encoders from pre-trained or foundation models which are listed in Table 4. The experimental results demon-



**Fig. 4.** Stability of model performance (F1-Score) versus variation number of test modalities (arranged in order from retinal surface to choroid vessels) for different models (a) and five diseases (b)

strate that the greater the similarity between pretrained and target images (smaller domain gap[21]), the better the performance [22]. In addition, different implementations in Eq.(4) also have an impact on the final results. In Table.4, V3 achieve the best performance.

**Table 4.** Ablation studies on pretrained images and implementations of KV in Eq.(4). V1 refers to randomly sampling and addition, V2 refers to concatenation, and V3 refers to addition.

Weights	ImageNet[14]	RETFound [23]	VisionFM [24]	Ours
Images	Natural	Fundus	FFA	UW-OCTA
F1-score	21.66	31.43	37.89	<b>59.60</b>
Implementation	V1	V2	V3	-
F1-score	55.42	57.98	<b>59.60</b>	-

## 5. CONCLUSION

In this study, we first introduce a new multi-modal ultra-wide retinal OCTA dataset (M3OCTA) with the multi-diseases annotations. Then the cross-modal fusion network (CMF-Net) is proposed to leverage multi-modal OCTA images for efficiently diagnosing a range of ophthalmological diseases. Moreover, CMF-Net enables a varying number of input modalities without performance drop during the inference phase, which greatly increases the applicability and compatibility in different clinical scenarios. Both quantified and visualized results indicate that our model exhibits robust performance. Our fusion method successfully provides a new view for leveraging the multi-modalities to innovative ophthalmic imaging technology. Beyond the fusion method, we delved deeper into the potential application of retina-related foundational models on UW-OCTA. Our findings also suggest that incorporating more information derived from multi-modalities for foundational models is also essential for robust performance.

## 6. ACKNOWLEDGEMENT

This work was supported in part by the Research Grants Council (RGC) of Hong Kong SAR (ECS24211020, GRF14203821, GRF14216222), the Innovation and Technology Fund (ITF) of Hong Kong SAR (ITS/240/21), the Science, Technology and Innovation Commission (STIC) of Shenzhen Municipality (SGDX20220530111005039)

## 7. COMPLIANCE WITH ETHICAL STANDARDS

This study was performed in line with the principles of the Declaration of Helsinki and approved by the local institutional review board. Informed written consent was obtained from all institutional patients.

## 8. REFERENCES

- [1] R. J. Antcliff *et al.*, “Comparison between optical coherence tomography and fundus fluorescein angiography for the detection of cystoid macular edema in patients with uveitis,” *Ophthalmology*, vol. 107, no. 3, pp. 593–599, 2000.
- [2] Z. Chen *et al.*, “Dual-consistency semi-supervision combined with self-supervision for vessel segmentation in retinal octa images,” *Biomedical Optics Express*, vol. 13, no. 5, pp. 2824–2834, 2022.
- [3] F. Zheng *et al.*, “Advances in swept-source optical coherence tomography and optical coherence tomography angiography,” *Advances in Ophthalmology Practice and Research*, 2022.
- [4] M. Nawaz *et al.*, “Unravelling the complexity of optical coherence tomography image segmentation using machine and deep learning techniques: A review,” *Computerized Medical Imaging and Graphics*, p. 102269, 2023.
- [5] D. M. Sampson *et al.*, “Towards standardizing retinal optical coherence tomography angiography: a review,” *Light: science & applications*, vol. 11, no. 1, p. 63, 2022.
- [6] W. Wang *et al.*, “Learning two-stream cnn for multi-modal age-related macular degeneration categorization,” *IEEE Journal of Biomedical and Health Informatics*, vol. 26, no. 8, pp. 4111–4122, 2022.
- [7] Y. Zhang, F. Xie, and J. Chen, “Tformer: A throughout fusion transformer for multi-modal skin lesion diagnosis,” *Computers in Biology and Medicine*, vol. 157, p. 106712, 2023.
- [8] Y. Giarratano *et al.*, “Automated segmentation of optical coherence tomography angiography images: benchmark data and clinically relevant metrics,” *Translational vision science & technology*, vol. 9, no. 13, pp. 5–5, 2020.
- [9] Y. Ma *et al.*, “Rose: a retinal oct-angiography vessel segmentation dataset and new model,” *IEEE TMI*, vol. 40, no. 3, pp. 928–939, 2020.
- [10] M. Díaz, J. Novo, P. Cutrín, F. Gómez-Ulla, M. G. Penedo, and M. Ortega, “Automatic segmentation of the foveal avascular zone in ophthalmological oct-a images,” *PloS one*, vol. 14, no. 2, p. e0212364, 2019.
- [11] A. Agarwal *et al.*, “The foveal avascular zone image database (fazid),” in *Applications of Digital Image Processing XLIII*, vol. 11510. SPIE, 2020, pp. 507–512.
- [12] M. Li *et al.*, “Image projection network: 3d to 2d image segmentation in octa images,” *IEEE TMI*, vol. 39, no. 11, pp. 3343–3354, 2020.
- [13] J. Hou *et al.*, “Deep-octa: Ensemble deep learning approaches for diabetic retinopathy analysis on octa images,” in *MICCAI Challenge on Mitosis Domain Generalization*. Springer, 2022, pp. 74–87.
- [14] R. Bachmann *et al.*, “Multimae: Multi-modal multi-task masked autoencoders,” in *European Conference on Computer Vision*. Springer, 2022, pp. 348–367.
- [15] K. He *et al.*, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 000–16 009.
- [16] Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [17] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [18] Nicki Skafte Detlefsen *et al.*, “TorchMetrics - Measuring Reproducibility in PyTorch,” Feb. 2022. [Online]. Available: <https://github.com/Lightning-AI/torchmetrics>
- [19] K. He *et al.*, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [20] Z. Liu *et al.*, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11 976–11 986.
- [21] H. Wei *et al.*, “Caudr: A causality-inspired domain generalization framework for fundus-based diabetic retinopathy grading,” *arXiv preprint arXiv:2309.15493*, 2023.
- [22] P. Shi *et al.*, “Generalist vision foundation models for medical imaging: A case study of segment anything model on zero-shot medical segmentation,” *Diagnostics*, vol. 13, no. 11, p. 1947, 2023.
- [23] Y. Zhou *et al.*, “A foundation model for generalizable disease detection from retinal images,” *Nature*, pp. 1–8, 2023.
- [24] J. Qiu *et al.*, “Visionfm: a multi-modal multi-task vision foundation model for generalist ophthalmic artificial intelligence,” *arXiv preprint arXiv:2310.04992*, 2023.