

1 **FFA Sora: generating fundus fluorescein angiography videos for**
2 **healthcare data sharing**

3

4 **Authors:**

5 Xinyuan Wu, MD^{1#}, Lili Wang, MS^{2#}, Ruoyu Chen, MD^{1#}, Bowen Liu, MS¹, Weiyi Zhang, MS¹,

6 Xi Yang, MD, PhD³, Yifan Feng, MD, PhD³, Mingguang He, MD, PhD^{1,4,5} and Danli Shi, MD,

7 PhD^{1,4}✉

8

9 **Affiliations:**

10 1. School of Optometry, The Hong Kong Polytechnic University, Hong Kong.

11 2. Department of Computing, The Hong Kong Polytechnic University, Hong Kong.

12 3. Department of Ophthalmology, Zhongshan Hospital, Fudan University, Shanghai, China.

13 4. Research Centre for SHARP Vision (RCSV), The Hong Kong Polytechnic University, Hong
14 Kong.

15 5. Centre for Eye and Vision Research (CEVR), 17W Hong Kong Science Park, Hong Kong.

16

17 # Xinyuan Wu, Lili Wang and Ruoyu Chen contributed equally to this work.

18

19 **Correspondence:**

20 Dr. Danli Shi, The Hong Kong Polytechnic University, Hong Kong, China. Email:

21 danli.shi@polyu.edu.hk

22 Prof. Mingguang He, Chair Professor of Experimental Ophthalmology, The Hong Kong

23 Polytechnic University, Hong Kong, China. Email: mingguang.he@polyu.edu.hk

24

25 **Manuscript word count:** 2995 words.

26

27 **Running title:** FFA Sora generates fundus fluorescein angiography videos

28

29 **Key Points:**

30 **Question** Can artificial intelligence (AI) models generate dynamic fundus fluorescein
31 angiography (FFA) videos directly from textual descriptions of retinal conditions?

32

33 **Findings** This study developed FFA Sora, an AI-driven text-to-video model capable of
34 generating high-quality FFA videos from medical textual descriptions. The generated videos
35 depicted retinal abnormalities, such as non-perfusion areas, leakage, and microaneurysms.
36 Objective and subjective assessments confirmed their alignment with corresponding lesion or
37 diseases descriptions. Retrieval Evaluation demonstrated the model's privacy-preserving
38 performance.

39

40 **Meaning** AI-driven text-to-video models offer an approach to visualizing various retinal
41 conditions from text prompts, potentially promoting privacy-preserving data sharing and
42 improving comprehension in medical education.

43

44 **Abstract**

45 **Importance** Medical data sharing faces strict restrictions. Text-to-video generation shows
46 potential for creating realistic medical data while preserving privacy, offering a solution for
47 cross-centre data sharing and medical education.

48

49 **Objective** To develop and evaluate a text-to-video generative model that converts report texts
50 into dynamic FFA videos, enabling visualization of retinal vascular and structural abnormalities.

51

52 **Design, Setting, and Participants** This study retrospectively collected anonymized FFA data
53 from a tertiary hospital in China, including 3625 FFA videos. Using this dataset, a text-to-video
54 model named FFA Sora was developed and evaluated, with 80% for training, 10% for validation
55 and the remaining 10% for testing. The model integrates two key components: the Wavelet-
56 Flow Variational Autoencoder (WF-VAE) and the diffusion transformer (DiT).

57

58 **Main Outcomes and Measures** The model's performance was assessed through objective
59 metrics, including Fréchet Video Distance (FVD), Learned Perceptual Image Patch Similarity
60 (LPIPS), and Visual-question-answering Score (VQAScore). FFA domain-specific evaluation
61 was measured by Bidirectional Encoder Representations from Transformers Score
62 (BERTScore), and image retrieval was evaluated by Recall@K. Human evaluations were
63 conducted on a scale of 1 (best) to 5 (worst).

64

65 **Results** The generated FFA videos demonstrated retinal abnormalities from the input text, as

66 confirmed by objective metrics: FVD = 2273, LPIPS = 0.48 ± 0.043 , and VQAScore = 0.61
67 ± 0.078 . Domain-specific evaluations showed alignment between the generated videos and
68 textual prompts, with BERTScore of 0.35 ± 0.091 . Additionally, the model demonstrated
69 privacy-preserving performance in retrieval evaluations, achieving an average Recall@K of
70 0.073 . Human assessments indicated visual quality, with an average score of 1.570 ± 0.440 .

71

72 **Conclusions and Relevance** This study demonstrated that FFA Sora, an AI-driven text-to-
73 video model, generated FFA videos from textual descriptions, potentially improving
74 visualization for clinical and educational purposes. Its privacy-preserving design may address
75 key challenges in data sharing while trying to ensure compliance with confidentiality standards.

76

77 **Introduction**

78 Fundus fluorescein angiography (FFA) enables the dynamic visualization of retinal blood flow
79 and lesional changes through the intravenous injection of fluorescein sodium. This technique
80 can be important for assessing blood-retina barrier function and diagnosing various retinal
81 vascular diseases, including diabetic retinopathy (DR) and choroidal neovascularization
82 (CNV).¹⁻⁵ Despite its importance, accurately interpreting FFA images demands expertise. For
83 ophthalmologists and medical students, the challenge of correlating textual description on text
84 book with corresponding FFA findings may be compounded by the complexity of disease
85 presentations and the limited availability of representative FFA cases.⁶

86

87 The growing application of artificial intelligence (AI) offers a potential solution for automatic
88 image interpretation. However, developing and validating these tools requires extensive, high-
89 quality labeled datasets, which can be difficult to obtain. Strict regulations governing the use
90 and exchange of personal health information protect patient privacy but can inadvertently
91 hinder the collaborative data sharing necessary to build effective AI models.⁷ Additionally,
92 healthcare providers often work in isolation, making them reluctant to share the clinical data
93 for boosting AI model training. Therefore, innovative strategies are needed to overcome the
94 challenges.

95

96 Open AI's Sora has been reported to generate realistic videos based on text inputs,⁸ paving the
97 way for the synthesis of well-labeled, diversified, and privacy-preserving medical data.

98 Recently, diffusion models (DMs) have gained attention for their capacity to generate high-

99 quality, temporally consistent images and videos through iterative denoising processes.⁹ These
100 models have been applied to generate medical data, including chest X-ray images¹⁰, brain MRI
101 scans,¹¹ multi-modal retinal images,¹² endoscopy videos,¹³ and facilitate neurosurgery.¹⁴ These
102 technological advancements have expanded applications across multiple domains, including
103 medical education, clinical training, surgical simulation, medical knowledge dissemination,
104 doctor-patient communication, and biomedical data augmentation.⁸ However, generating
105 videos from textual description has yet to be explored for creating lesion-preserving dynamic
106 FFA videos in ophthalmology.

107

108 To bridge the gap, we developed FFA Sora, an AI-driven model capable of producing dynamic
109 FFA videos from textual reports. This cross-modality generative model might assist in privacy-
110 preserving data sharing and medical education.

111

112 **Methods**

113 We utilized de-identified existing data for our study, which received approval from the
114 Institutional Review Board of the Hong Kong Polytechnic University. Our study adhered to
115 Strengthening the Reporting of Observational Studies in Epidemiology (STROBE) reporting
116 guidelines.

117

118 **Dataset**

119 This study uses a subset of a retrospective FFA dataset.¹⁵ FFA data were obtained from a tertiary
120 hospital in China. The dataset included both the medical records and FFA examinations of

121 patients assessed between November 2016 and December 2019. All FFA videos were captured
122 using Zeiss FF450 Plus and Heidelberg Spectralis systems (Heidelberg, Germany) at a
123 resolution of 768×768 pixels. The collection encompassed a wide variety of ocular conditions,
124 such as DR, retinal vein occlusion (RVO), and central serous chorioretinopathy (CSC). To try
125 to ensure high-quality data, we implemented a filtering process based on vessel area ratios
126 within the FFA frames, excluding any video whose frame vessel area ratios fell below 0.005.

127

128 To address the issue of inconsistent and insufficient frame counts in the original FFA videos,
129 we standardized all videos to a fixed frame count of 21 frames. For videos with more than 21
130 frames, we selected 21 frames in reverse chronological order. For videos with fewer than 21
131 frames, additional frames were interpolated to reach the required count. This was achieved
132 using a frame interpolation technique based on linear interpolation, implemented via the
133 OpenCV library.

134

135 The process involved extracting all frames from the input video and computing the interpolation
136 positions based on the target frame count. For each interpolated frame, a weighted linear
137 interpolation between two adjacent frames was employed to generate smooth transition frames.

138 Specifically, let an interpolated frame lie between the i -th frame and the $(i+1)$ -th frame of the
139 original video. The pixel values of the interpolated frame are computed as follows:

140

$$141 \quad F_{interpolated} = (1 - t) \times F_i - t \times F_{i+1}$$

142 where:

143 F_i and F_{i+1} denote the pixel values of the i -th and the $(i+1)$ -th frames, respectively,
144 t is the interpolation ratio, which takes values in the range $[0,1]$ and determines the relative
145 contribution of the two adjacent frames to the interpolated frame.

146

147 The interpolation ratio t is calculated based on the relative temporal position of the interpolated
148 frame within the sequence. Once all interpolated frames were generated, they were combined
149 with the original frames to produce a standardized video consisting of exactly 21 frames. All
150 patient data were stratified by patient ID and randomly divided into 80% for training, 10% for
151 validation and 10% for testing.

152

153 **Model Architecture**

154 The framework of FFA Sora was based on the Open-Sora Plan repository ([v1.3.0](#)),¹⁶ which we
155 specifically adapted to address the task of generating FFA images from medical textual
156 descriptions. It consisted of two primary components: the Wavelet-Flow Variational
157 Autoencoder (WF-VAE) and the Diffusion Transformer (DiT). Both were trained on a single
158 NVIDIA A800 80,000 MB (80 GB) PCIe GPU with CUDA 12.1. The programming
159 environment used for this project is Python 3.8.20 with PyTorch 2.1.0 (CUDA 12.1).

160

161 The WF-VAE architecture¹⁷, derived from the Stable-Diffusion Image VAE, was fine-tuned and
162 optimized for the specialized task of text-to-FFA image generation in the medical domain. A
163 key architectural enhancement involved the transformation of Conv2D layers into
164 CausalConv3D layers, enabling the model to process temporal video data, which was essential

165 for the sequential nature of FFA imaging. Furthermore, multi-level wavelet transforms were
166 employed in the encoder to decompose video signals into frequency sub-bands, allowing low-
167 frequency components, which is critical for FFA imaging details, to bypass the backbone
168 network via the Main Energy Flow Pathway, thereby reducing computational redundancy. The
169 model was trained with a latent dimension of 8 for 200,000 steps, with GPU memory
170 consumption reaching approximately 19,954 MB.

171

172 The second major component was the DiT architecture,¹⁸ which employed a Transformer-based
173 architecture with cross-modal attention mechanisms to optimize text-to-FFA image generation.
174 The denoising process was performed in a low-dimensional latent space compressed by the
175 VAE, where each Cross-DiT Block, consisting of self-attention, cross-attention, and feed-
176 forward networks, was enhanced with gating mechanisms to refine feature extraction and fusion.
177 This design improved the model's ability to capture both high-level semantics and fine-grained
178 details, ensuring that the generated FFA images accurately reflected clinically relevant features
179 such as global vascular structures and local pathological regions. The fine-tuning process was
180 conducted on the same dataset for 50,000 steps over approximately 72 hours, with GPU
181 memory consumption reaching approximately 73,728 MB.

182

183 **Model Inference**

184 During the inference phase, the model processed input textual descriptions and generated
185 corresponding FFA images with a resolution of $21 \times 512 \times 512$. The input text typically contained
186 detailed medical descriptions of retinal conditions, such as microaneurysms, non-perfusion area,

187 leakage, and other abnormalities, which were crucial for FFA-based diagnosis. To ensure that
188 the generated images align closely with the described retinal conditions, our model utilized the
189 pre-trained text encoder to capture the semantic nuances of the input effectively. These latent
190 features, conditioned on the text, were passed through the CausalConv3D layers in the decoder,
191 which reconstructed the FFA images by generating temporal and spatial patterns consistent with
192 real FFA images. The inference process employed a tiled convolution approach, which
193 optimized memory usage and ensured high-quality image output. By applying this technique,
194 the model generated coherent and clinically relevant FFA images that accurately depict the
195 specified retinal conditions, while maintaining computational efficiency during inference. The
196 overview of FFA Sora is presented in **Figure 1**.

197

198 **Objective Evaluation**

199 Several objective metrics were employed to assess the quality of the generated FFA videos,
200 including Fréchet Video Distance (FVD),¹⁹ Learned Perceptual Image Patch Similarity
201 (LPIPS),²⁰ and Visual-question-answering Score (VQAScore).²¹ FVD discerns the likeness in
202 feature distribution between authentic and synthetic videos, encapsulating the holistic
203 excellence and unity. FVD typically ranges from 0 to higher values, with lower values
204 indicating better similarity. LPIPS, utilizing deep learning-inspired features, measures the
205 perceptual resemblance by comparing image segments, thus potentially revealing nuances
206 beyond the capabilities of conventional pixel-based assessments. LPIPS scores range from 0 to
207 1, with lower scores considered higher perceptual similarity. VQAScore evaluates the
208 alignment of text-to-video models. This score ranges from 0 to 1, where higher values represent

209 better quality. Collectively, these evaluation criteria potentially provide a comprehensive
210 analysis of the generated FFA video quality.

211

212 For FFA domain-specific evaluation, we used FFA-GPT to translate the generated FFA video
213 into text report and compared it with the original text prompt using Bidirectional Encoder
214 Representations from Transformers Score (BERTScore)²² as similarity measure. BERTScore
215 utilizes the pre-trained BERT language model to measure conceptual overlap between
216 generated and reference content. Unlike purely lexical metrics, it ideally captures nuanced
217 alignments within broader contexts, providing a deeper assessment of textual congruence. This
218 score ranges from 0 to 1, with higher values indicating better alignment between the generated
219 and original content.

220

221 For privacy preservation evaluation, we assessed FFA Sora's ability to generate content-
222 preserving yet deidentifiable data. We employed an image-to-image retrieval method, using
223 Recall@K²³ to measure the proportion of matched visuals retrieved among the top K results. It
224 ranges from 0 to 1, with lower values indicating that the original lesion information cannot be
225 captured in the generated videos. Since the accuracy of the generated videos can be validated
226 through other parameters mentioned above, Recall@K serves as an indicator of FFA Sora's
227 ability in preserving privacy.

228

229 **Human Assessment**

230 This evaluation followed the methodology described previously.^{24,25} Three ophthalmologists

231 (X.W., X.Y. and Y.F.) with 3, 5, and 10 years of experience respectively, representing trainee,
232 attending, and expert levels, reviewed 50 randomly selected FFA videos generated from the test
233 set, comparing them against the corresponding ground-truth FFA videos. Each video was rated
234 on a scale from 1 to 5, where 1 represented excellent quality and 5 indicated very poor quality.
235 The specific criteria were defined in eFigure 3 in the Supplement.

236

237 **Results**

238 The final dataset used in the study consists of 3625 FFA videos, among which 2851 videos were
239 randomly selected for training, 387 videos for validation, and 387 videos for testing.
240 Demographic characteristics of the dataset were presented in eTable1 in the Supplement. FFA
241 reports in the dataset included description of various lesions, such as microaneurysms, leakage,
242 neovascularization, capillary non-perfusion, and macular edema (eTable 2 in the Supplement).

243

244 Figure 2 demonstrates representative examples of generated FFA videos, including
245 corresponding prompts, ground-truth videos, and the generated video frames. FFA Sora could
246 generate detailed FFA videos from text prompts, illustrating retinal lesions and vascular
247 abnormalities, such as leakage, neovascularization, and microaneurysms, etc. Our model could
248 produce FFA videos for other retinal and choroidal diseases as well, such as uveitis, retinitis
249 pigmentosa (RP), etc. Representative frames of these videos are shown in eFigure 1 and eFigure
250 2 in the Supplement.

251

252 **Objective evaluation showed excellent quality of generated FFA videos**

253 The generated videos were comprised of 21 frames that precisely document the entire process
254 of FFA examination, with a particular emphasis on the venous and late phases. The model we
255 employed exhibited satisfactory performance in producing FFA videos when applied to the test
256 dataset, as evidenced by FVD = 2273, LPIPS = 0.48 ± 0.043 , and VQAScore = 0.61 ± 0.078
257 (Figure 4A). These metrics served as benchmarks in assessing the fidelity and quality of the
258 generated FFA videos.

259

260 In addition to the evaluation with above standard generative metrics, we also proposed domain-
261 specific evaluation strategies. Based on previous development of FFA-GPT,¹⁵ a two-stage
262 system which can generate FFA report automatically, our model was evaluated to determine
263 whether the content of the generated videos aligns with the actual characteristics of retinal
264 diseases as denoted by specified textual prompts, with BERTScore = 0.35 ± 0.091 (Figure 4A).

265 To visualize these findings, an example of reports by FFA-GPT based on FFA Sora's generated
266 FFA videos was presented in Figure 3. Additionally, Figure 4C demonstrated an example where
267 the objective evaluation results were relatively poor, with ranking in the lowest 20% across all
268 parameters.

269

270 **Image retrieval demonstrated FFA Sora's excellent performance on privacy preserving**

271 To further explore FFA Sora's robust feature representation, we investigated its performance in
272 image retrieval tasks, specifically examining the relationship between the generated videos and
273 their ground-truth counterparts. The results of this investigation were summarized in Figure 4B,
274 where we reported Recall@K scores that include 0.02, 0.04, 0.16 for K = 5, 10, and 50,

275 respectively. Furthermore, the mean recall was calculated to be 0.073, reflecting disparities
276 among generated and ground-truth videos. These relatively low Recall@K values indicated that
277 FFA Sora could prevent the leakage of confidential image information while reserving the
278 essential pathological features during video generation.

279

280 **Human assessment confirmed high visual quality of generated FFA videos**

281 A comprehensive analysis was conducted on fifty randomly selected videos produced via our
282 model, which were subjected to meticulous review by three ophthalmologists. The generated
283 videos were systematically compared with the corresponding textual reports and the ground-
284 truth FFA videos to rigorously assess their accuracy and fidelity. The subjective visual quality
285 score averaged at 1.570 ± 0.440 , with a Kappa value of 0.761. The specific results are presented
286 in **Figure 5A** and **5B**, while **Figure 5C** illustrated an example where the generated video received
287 consistently low subjective ratings. To clarify the scoring scale, examples corresponding to
288 ratings from 1 to 5 were presented in **eFigure 3** in the Supplement. Additionally, **eFigure 4** in
289 the Supplement demonstrated cases where discrepancies arose between objective and
290 subjective evaluations.

291

292 **Discussion**

293 In this study, we developed FFA Sora, a text-to-video diffusion transformer designed to generate
294 FFA videos from textual descriptions. This research demonstrated the potential of text-to-video
295 models for visualizing FFA-related terms and associated manifestations of several retinal
296 disease, offering a solution for privacy-preserving data sharing in machine learning model

297 training and medical AI applications.

298

299 Generative models represent potential advancements in medical data generation of their ability
300 to produce images of high quality, stability, and diversity.²⁶⁻²⁸ Previous research has explored
301 applications of generative adversarial networks (GANs) in retinal image quality enhancement,
302 domain adaptation, vessel segmentation, and cross-modality generation of fundus
303 autofluorescence, FFA, and indocyanine green angiography.²⁹⁻³⁵ More recently, diffusion
304 models beat GANs in image generation⁹ and have been applied to ocular imaging tasks, such
305 as ocular surface structure segmentation for meibomian gland dysfunction grading³⁶ and data
306 augmentation using retinography images to support deep learning model training.³⁷ DiT, a
307 transformer-based architecture for diffusion models, enhances scalability and generative
308 performance by replacing the conventional U-Net backbone with a transformer. Its architectural
309 innovation improves feature representation, synthesis quality, and cross-modality
310 adaptability.³⁸⁻⁴⁰ Building on this foundation, FFA Sora is an advanced application of
311 generative models with a potentially better architectural extensibility. Compared with our
312 previous image-to-video study^{24,25}, FFA Sora exhibits higher LPIPS and FVD values because
313 text-based video generation introduces variations in lesion presentation, whereas image-based
314 video generation retains the overall structural integrity of the original input. Comprehensive
315 evaluations confirmed that the produced videos preserved diagnostic plausibility when visually
316 diversified.

317

318 Data sharing in AI model training is often hindered by rigorous privacy regulations and

319 concerns regarding patient confidentiality, particularly with sensitive medical imaging and
320 clinical videos.^{41,42} These challenges limit collaboration and the creation of large, diverse
321 datasets needed for training and validating AI models. Additionally, there is a risk that AI
322 models may inadvertently generate videos by recalling data from their training set, potentially
323 compromising privacy.⁴³ The notably low Recall@K scores in our image retrieval evaluation
324 demonstrated the visual differences between synthetic videos and the ground-truth ones, thus
325 preserving privacy disease information. The visual distinction in generated outputs ensures that
326 the resulting datasets are not mere replications of existing clinical images, allowing for
327 simulating the diversity in real-world cases of the same disease. By mitigating these concerns,
328 FFA Sora facilitates cross-center data sharing and reproducibility, paving the way for potentially
329 improving collaborative research, clinical training, and the development of AI-driven
330 diagnostic tools.

331

332 AI-generated content in the medical domain has broad implications. Research on AI-assisted
333 medical image generation has demonstrated its value in medical education⁶ and the training of
334 convolutional neural networks (CNNs) for medical purposes such as diagnosis and disease
335 classification.^{10,44} For instance, Tabuchi H, et al.⁶ introduced a text-to-image model to produce
336 ultra-widefield (UWF) retinal images, significantly improving diagnostic accuracy among
337 medical students. This outcome highlights the potential of integrating AI technologies into
338 medical education. Compared with their model, FFA Sora offers a more dynamic demonstration
339 of retinal anatomy, vascular perfusion, and pathology, thereby facilitating ophthalmology
340 training. FFA Sora can simulate rare and complex pathologies, offering an interactive learning

341 tool to improve diagnostic skills, decision-making, and treatment planning, particularly in
342 resource-limited settings. FFA Sora can also contribute to case libraries for board exams and
343 clinical assessments, providing trainees with diverse retinal conditions that may not be
344 encountered during rotations. Its accessibility potentially ensures its broad applicability across
345 educational institutions and healthcare facilities, potentially equipping future ophthalmologists
346 with AI-driven tools for diagnosing and managing retinal diseases. Additionally, FFA Sora has
347 potential in clinical workflows, particularly for patient education and disease progression
348 simulation. For example, clinicians counseling patients with severe NPDR can use AI-
349 generated FFA videos to illustrate disease progression. This dynamic visualization can help
350 patients understand the impact of uncontrolled diabetes and the benefits of timely interventions,
351 potentially improving adherence to treatments such as blood sugar control, laser therapy, and
352 anti-VEGF injections.

353

354 **Limitations**

355 This research has several limitations. First, although FFA Sora demonstrates high accuracy and
356 reliability, its real-world clinical utility requires further validation. Broader external validation
357 using diverse datasets would enhance its practical relevance. Second, our evaluation approach
358 may not fully capture the generative quality. Our evaluators were not masked to the synthetic
359 nature of the generated videos. Future work should incorporate masked evaluations to better
360 quantify realism and ensure synthetic videos closely resemble real-world data. Finally, our
361 current approach focuses solely on visualizing FFA reports, without integrating other imaging
362 modalities. Multimodal imaging could improve AI-based models in ophthalmology, offering

363 more complete diagnostic and treatment support.

364

365 **Conclusion**

366 In conclusion, this study presents FFA Sora, a model that transforms FFA report text into
367 dynamic video representations. These synthetic videos appear to retain critical clinical details
368 while preserving patient privacy, offering a potential solution to data-sharing challenges in
369 healthcare. By enabling potentially secure and efficient visualization, FFA Sora may support
370 collaborative research, enhance medical education, and accelerate the development of AI-
371 driven diagnostic and medical systems.

372

373 **Acknowledgement**

374 We thank the InnoHK HKSAR Government for providing valuable supports.
375 The research work described in this paper was conducted in the JC STEM Lab of Innovative
376 Light Therapy for Eye Diseases funded by The Hong Kong Jockey Club Charities Trust.

377

378 **Funding**

379 M.H. disclose support for the research and publication of this work from the Global STEM
380 Professorship Scheme (P0046113), and Henry G. Leong Endowed Professorship in Elderly Vision
381 Health.

382

383 **Role of the Funder/Sponser**

384 The funders had no role in the design and conduct of the study, collection, management,
385 analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and

386 decision to submit the manuscript for publication.

387

388 **Author contributions**

389 D.S. conceived the study. L.W. built the text-to-video model. D.S., R.C., B.L., W.Z. conducted
390 the literature search, analyzed the data. L.W. performed objective evaluation. X.W., X.Y., Y.F.
391 completed visual evaluation. X.W., L.W. wrote the manuscript. L.W., X.W. organized figures
392 and tables in this study. M.H. provided the data and facilities. All authors critically revised the
393 manuscript.

394

395 **Conflict of Interest Disclosures**

396 The authors declare no conflict interest.

397

398 **Access to data and data analysis**

399 Code is available at <https://github.com/PKU-YuanGroup/Open-Sora-Plan/tree/v1.3.0>.

400 The generated dataset will be available on request, but the original dataset involving human
401 participants will remain confidential.

402

Reference

- 404 1. Burns SA, Elsner AE, Gast TJ. Imaging the Retinal Vasculature. *Annual Review of Vision Science*.
405 2021-09-15 2021;7:129-153. doi:10.1146/annurev-vision-093019-113719
- 406 2. Sulzbacher F, Pollreisz A, Kaider A, Kicking S, Sacu S, Schmidt-Erfurth U. Identification and
407 clinical role of choroidal neovascularization characteristics based on optical coherence
408 tomography angiography. *Acta Ophthalmol*. Jun 2017;95(4):414-420. doi:10.1111/aos.13364
- 409 3. Do DV. Detection of new-onset choroidal neovascularization. *Curr Opin Ophthalmol*. May
410 2013;24(3):244-7. doi:10.1097/ICU.0b013e32835fd7dd
- 411 4. Coscas GJ, Lupidi M, Coscas F, Cagini C, Souied EH. OPTICAL COHERENCE TOMOGRAPHY
412 ANGIOGRAPHY VERSUS TRADITIONAL MULTIMODAL IMAGING IN ASSESSING THE ACTIVITY OF
413 EXUDATIVE AGE-RELATED MACULAR DEGENERATION: A New Diagnostic Challenge. *Retina*. Nov
414 2015;35(11):2219-28. doi:10.1097/iae.0000000000000766
- 415 5. Sabanayagam C, Banu R, Chee ML, et al. Incidence and progression of diabetic retinopathy:
416 a systematic review. *Lancet Diabetes Endocrinol*. Feb 2019;7(2):140-149. doi:10.1016/s2213-
417 8587(18)30128-1
- 418 6. Tabuchi H, Engelmann J, Maeda F, et al. Using artificial intelligence to improve human
419 performance: efficient retinal disease detection training with synthetic images. *The British Journal*
420 *of Ophthalmology*. 2024-09-20 2024;108(10):1430-1435. doi:10.1136/bjo-2023-324923
- 421 7. Price WN, 2nd, Cohen IG. Privacy in the age of medical big data. *Nat Med*. Jan 2019;25(1):37-
422 43. doi:10.1038/s41591-018-0272-7
- 423 8. O'Callaghan J. How OpenAI's text-to-video tool Sora could change science – and society.
424 *Nature*. 2024-03-12 2024;627(8004):475-476. doi:10.1038/d41586-024-00661-0
- 425 9. Dhariwal P, Nichol A. Diffusion models beat gans on image synthesis. *Advances in neural*
426 *information processing systems*. 2021;34:8780-8794.
- 427 10. Bluethgen C, Chambon P, Delbrouck J-B, et al. A vision-language foundation model for the
428 generation of realistic chest X-ray images. *Nature Biomedical Engineering*. 2024-08-26
429 2024;doi:10.1038/s41551-024-01246-y
- 430 11. Dorjsembe Z, Pao H-K, Odonchimed S, Xiao F. Conditional Diffusion Models for Semantic 3D
431 Brain MRI Synthesis. *IEEE Journal of Biomedical and Health Informatics*. 7/2024 2024;28(7):4084-
432 4093. doi:10.1109/JBHI.2024.3385504
- 433 12. Chen R, Zhang W, Liu B, et al. EyeDiff: text-to-image diffusion model improves rare eye
434 disease diagnosis. *arXiv preprint arXiv:241110004*. 2024;
- 435 13. Li C, Liu H, Liu Y, et al. Endora: Video Generation Models as Endoscopy Simulators. 2024-03-
436 17 2024;
- 437 14. Sener F, Saraf R, Yao A. Transferring Knowledge From Text to Video: Zero-Shot Anticipation
438 for Procedural Actions. *IEEE transactions on pattern analysis and machine intelligence*. 2023-06
439 2023;45(6):7836-7852. doi:10.1109/TPAMI.2022.3218596
- 440 15. Chen X, Zhang W, Xu P, et al. FFA-GPT: an automated pipeline for fundus fluorescein
441 angiography interpretation and question-answer. *npj Digital Medicine*. 2024-05-03 2024;7(1):1-
442 9. doi:10.1038/s41746-024-01101-z
- 443 16. Lin B, Ge Y, Cheng X, et al. Open-Sora Plan: Open-Source Large Video Generation Model.
444 2024:
- 445 17. Li Z, Lin B, Ye Y, et al. WF-VAE: Enhancing Video VAE by Wavelet-Driven Energy Flow for

446 Latent Video Diffusion Model. 2024:
447 18. Peebles WS, Xie S. Scalable Diffusion Models with Transformers. *2023 IEEE/CVF International*
448 *Conference on Computer Vision (ICCV)*. 2022:4172-4182.
449 19. Unterthiner T, van Steenkiste S, Kurach K, Marinier R, Michalski M, Gelly S. FVD: A new metric
450 for video generation. 2019;
451 20. Zhang R, Isola P, Efros AA, Shechtman E, Wang O. The Unreasonable Effectiveness of Deep
452 Features as a Perceptual Metric. presented at: 2018 IEEE/CVF Conference on Computer Vision and
453 Pattern Recognition (CVPR); 2018; <https://doi.ieeecomputersociety.org/10.1109/CVPR.2018.00068>
454 21. Lin Z, Pathak D, Li B, et al. Evaluating Text-to-Visual Generation with Image-to-Text
455 Generation. 2024:
456 22. Zhang T, Kishore V, Wu F, Weinberger KQ, Artzi Y. BERTScore: Evaluating Text Generation
457 with BERT. *ArXiv*. 2019;abs/1904.09675
458 23. Song HO, Xiang Y, Jegelka S, Savarese S. Deep Metric Learning via Lifted Structured Feature
459 Embedding. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
460 2015:4004-4012.
461 24. Zhang W, Yang J, Chen R, et al. Fundus to Fluorescein Angiography Video Generation as a
462 Retinal Generative Foundation Model. *ArXiv*. 2024;abs/2410.13242
463 25. Zhang W, Huang S, Yang J, et al. Fundus2Video: Cross-Modal Angiography Video Generation
464 from Static Fundus Photography with Clinical Knowledge Guidance. Springer Nature Switzerland;
465 2024:689-699.
466 26. Tian Y, Fan L, Isola P, Chang H, Krishnan D. StableRep: Synthetic Images from Text-to-Image
467 Models Make Strong Visual Representation Learners. *ArXiv*. 2023;abs/2306.00984
468 27. Rombach R, Blattmann A, Lorenz D, Esser P, Ommer B. High-Resolution Image Synthesis with
469 Latent Diffusion Models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
470 *(CVPR)*. 2021:10674-10685.
471 28. Croitoru FA, Hondru V, Ionescu RT, Shah M. Diffusion Models in Vision: A Survey. *IEEE Trans*
472 *Pattern Anal Mach Intell*. Sep 2023;45(9):10850-10869. doi:10.1109/tpami.2023.3261988
473 29. He S, Joseph S, Bulloch G, et al. Bridging the Camera Domain Gap With Image-to-Image
474 Translation Improves Glaucoma Diagnosis. *Transl Vis Sci Technol*. 2023;12(12):20-20.
475 doi:10.1167/tvst.12.12.20
476 30. Song F, Zhang W, Zheng Y, Shi D, He M. A deep learning model for generating fundus
477 autofluorescence images from color fundus photography. *Adv Ophthalmol Pract Res*. Nov-Dec
478 2023;3(4):192-198. doi:10.1016/j.aopr.2023.11.001
479 31. Shi D, He S, Yang J, Zheng Y, He M. One-shot Retinal Artery and Vein Segmentation via
480 Cross-modality Pretraining. *Ophthalmol Sci*. Mar-Apr 2024;4(2):100363.
481 doi:10.1016/j.xops.2023.100363
482 32. Chen R, Xu K, Zheng K, et al. Generating Multi-frame Ultrawide-field Fluorescein Angiography
483 from Ultrawide-field Color Imaging Improves Diabetic Retinopathy Stratification. *arXiv preprint*
484 *arXiv:240810636*. 2024;
485 33. Shi D, Zhang W, He S, et al. Translation of Color Fundus Photography into Fluorescein
486 Angiography Using Deep Learning for Enhanced Diabetic Retinopathy Screening. *Ophthalmol Sci*.
487 Dec 2023;3(4):100401. doi:10.1016/j.xops.2023.100401
488 34. Zhang W, Huang S, Yang J, et al. Fundus2Video: Cross-Modal Angiography Video Generation
489 from Static Fundus Photography with Clinical Knowledge Guidance. Springer Nature Switzerland;

490 2024:689-699.

491 35. Chen R, Zhang W, Song F, et al. Translating color fundus photography to indocyanine green
492 angiography using deep-learning for age-related macular degeneration screening. *npj Digital*
493 *Medicine*. 2024-02-12 2024;7(1):34. doi:10.1038/s41746-024-01018-7

494 36. Guo X, Wen H, Hao H, et al. Randomness-restricted Diffusion Model for Ocular Surface
495 Structure Segmentation. *IEEE Trans Med Imaging*. Nov 11 2024;Ppdoi:10.1109/tmi.2024.3494762

496 37. Aktas B, Ates DD, Duzyel O, Gumus A. Diffusion-based data augmentation methodology for
497 improved performance in ocular disease diagnosis using retinography images. *International*
498 *Journal of Machine Learning and Cybernetics*. 2024/12/11 2024;doi:10.1007/s13042-024-02485-
499 w

500 38. Peebles W, Xie S. Scalable diffusion models with transformers. 2023:4195-4205.

501 39. Dosovitskiy A, Beyer L, Kolesnikov A, et al. An Image is Worth 16x16 Words: Transformers for
502 Image Recognition at Scale. *ArXiv*. 2020;abs/2010.11929

503 40. Bao F, Nie S, Xue K, et al. All are Worth Words: A ViT Backbone for Diffusion Models. 2023
504 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022:22669-22679.

505 41. Yang Y, Chen X, Lin H. Privacy preserving technology in ophthalmology. *Current Opinion in*
506 *Ophthalmology*. 11/2024 2024;35(6):431-437. doi:10.1097/ICU.0000000000001087

507 42. Aa S, M G, C A-R, et al. Privacy-preserving federated machine learning on FAIR health data:
508 A real-world application. *Computational and structural biotechnology journal*. 02/17/2024
509 2024;24doi:10.1016/j.csbj.2024.02.014

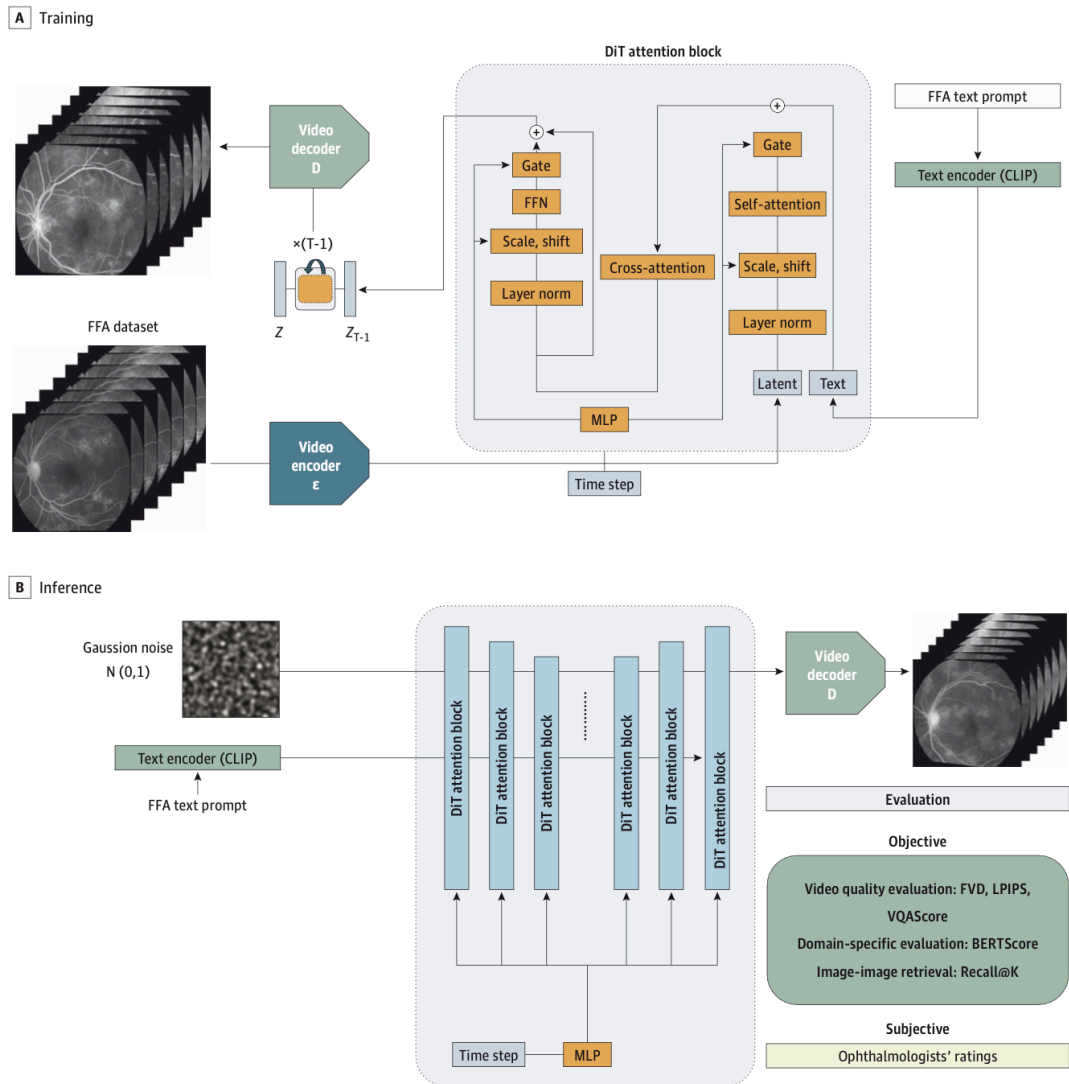
510 43. Collado-Mesa F, Alvarez E, Arheart K. The Role of Artificial Intelligence in Diagnostic
511 Radiology: A Survey at a Single Radiology Residency Training Program. *J Am Coll Radiol*. Dec
512 2018;15(12):1753-1757. doi:10.1016/j.jacr.2017.12.021

513 44. Coyner AS, Chen JS, Chang K, et al. Synthetic Medical Images for Robust, Privacy-Preserving
514 Training of Artificial Intelligence: Application to Retinopathy of Prematurity Diagnosis.
515 *Ophthalmology Science*. 2022-06 2022;2(2):100126. doi:10.1016/j.xops.2022.100126

516

517

518 **Figure legends**
 519



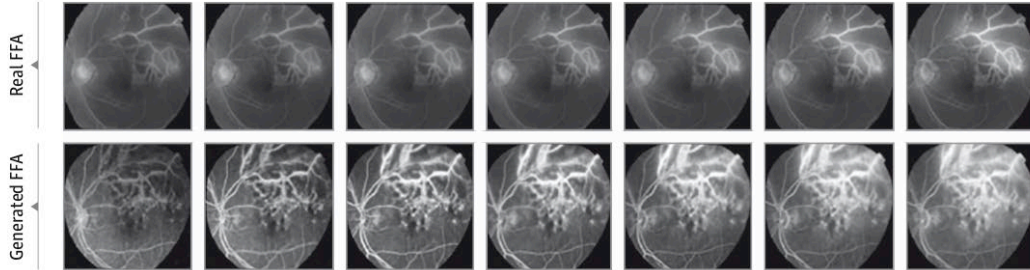
520

521 **Figure 1. Overview of the text-to-video diffusion model.** This figure illustrates a text-to-video
 522 framework for generating fundus fluorescein angiography (FFA) sequences using a diffusion
 523 transformer (DiT). BERTScore indicates bidirectional encoder representations from transformers
 524 score; FFN, feed-forward network; FVD, Fréchet video distance; LPIPS, learned perceptual image
 525 patch similarity; MLP, multilayer perceptron; VQAScore, visual question answering score. In A
 526 (during training), the input FFA video sequences are first encoded into latent representations by the
 527 video encoder ϵ . These latent representations undergo progressive refinement through multiple DiT
 528 attention blocks, which integrate self-attention for capturing spatial-temporal dependencies, cross-

529 attention for conditioning on textual embeddings, and diffusion-based denoising guided by time-
530 step embeddings. A pretrained text encoder (contrastive language-image pretraining [CLIP])
531 translates medical text prompts into semantic embeddings, providing textual conditioning that
532 guides the denoising process. The refined latent representations are subsequently decoded by the
533 video decoder D to reconstruct realistic FFA video sequences aligned with the input textual
534 descriptions, potentially enabling clinically meaningful ophthalmic imaging. In B (during inference),
535 the model synthesizes FFA sequences from Gaussian noise conditioned solely on clinical text
536 prompts. The framework is evaluated through objective metrics and subjective evaluations by
537 ophthalmologists to assess visual quality and clinical relevance. The details of the FFA text prompt
538 are as follows: Severe nonproliferative diabetic retinopathy at right eye FFA: (1) The arterial filling
539 and venous return time were basically normal. (2) Multiple microaneurysmal hyperfluorescence can
540 be seen in all quadrants, accompanied by mild leakage. No perfusion defects were formed around
541 the periphery, and obvious vascular leakage was seen in the nonperfusion zone. Diffuse capillary
542 leakage in the posterior pole, obvious leakage around the macula, and no obvious macular edema.
543 (3) No obvious abnormal fluorescence was found in the optic disc.
544

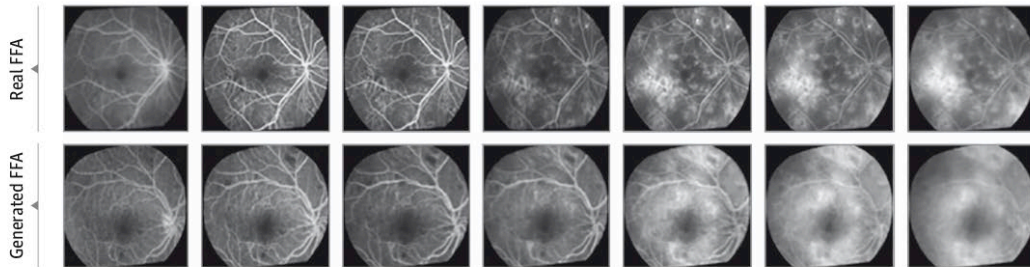
A Branch retinal vein occlusion

Text prompt
Ischemic branch retinal vein occlusion, nonperfusion area at left eye FFA:
1. Retinal artery diameter is narrowed, with the formation of temporal branch arteriovenous pressure imprint, slow flow of temporal branch venous reflux, and visible venous anastomosis on the temporal side of the macula; small blood vessels in the area governed by the temporal branch vein are dilated and leaky, with scattered microaneurysms, large areas of nonperfusion forming, affecting the temporal upper quadrant of the macula.
2. No obvious abnormalities were observed in the inferior and nasal midperiphery.
3. No obvious abnormalities were observed in the fluorescence of the optic disc.



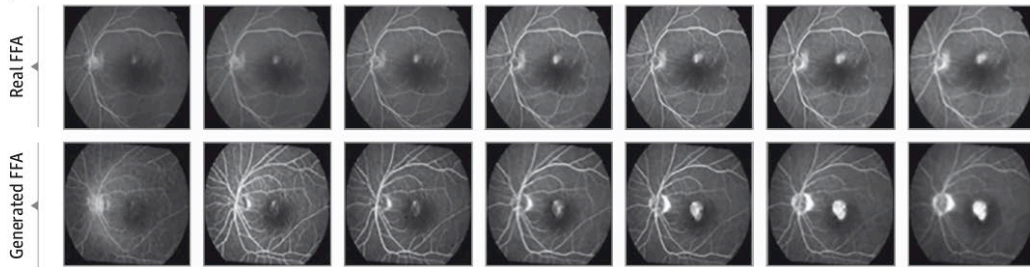
B Severe nonproliferative diabetic retinopathy

Text prompt
Severe nonproliferative diabetic retinopathy, arteriosclerosis, microaneurysm at right eye FFA:
1. The arterial filling and venous return time were basically normal.
2. Multiple microaneurysmal hyperfluorescence can be seen in all quadrants, accompanied by mild leakage; no perfusion defects were formed around the periphery; and obvious vascular leakage was seen in the nonperfusion zone. Diffuse capillary leakage in the posterior pole, obvious leakage below the temporal macula, and no obvious macular edema.
3. No obvious abnormal fluorescence was found in the optic disc.



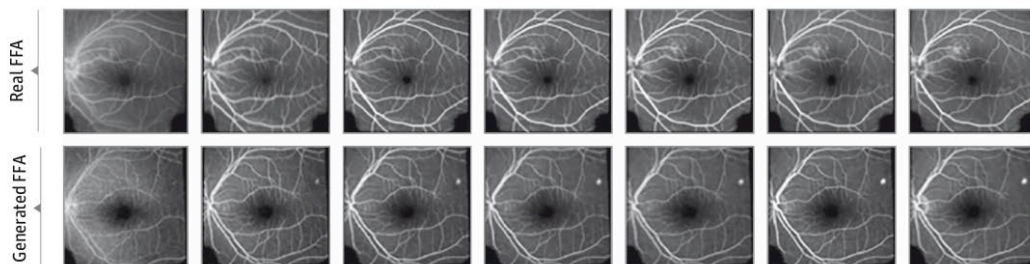
C Age-related macular degeneration

Text prompt
Depigmentation, pigmentary change, hypopigmentation at left eye FFA:
1. Focal chorioidal atrophy and pigment epithelial damage with focal leakage of fluorescent dye.
2. No obvious abnormal fluorescence was observed in the midperipheral and peripheral retina.
3. No obvious abnormal fluorescence was observed in the optic disc.

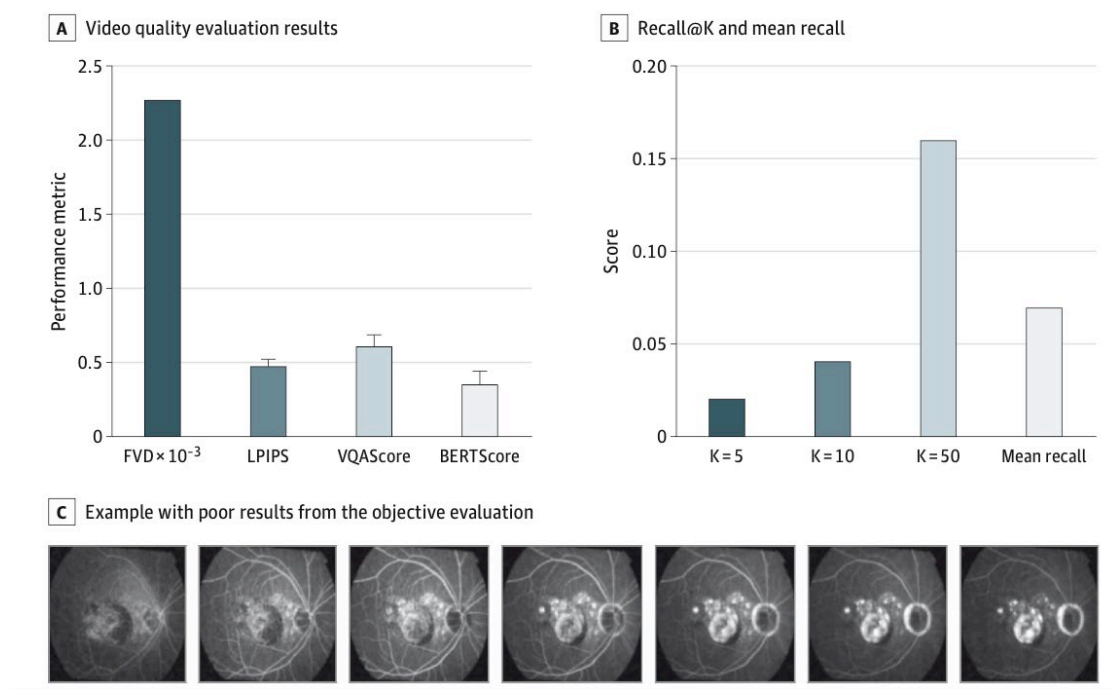


D Central serous chorioretinopathy

Text prompt
Central serous chorioretinopathy, pigmentary change at left eye FFA:
Small early hypofluorescent dots are seen at 1.5 disc diameter temporal to the optic disc, gradually fluorescing with no leakage, and with fluorescein visibility around the lesion. Several small pinpoint areas of hyperfluorescence are seen in the temporal macula. Optic disc appears normal.

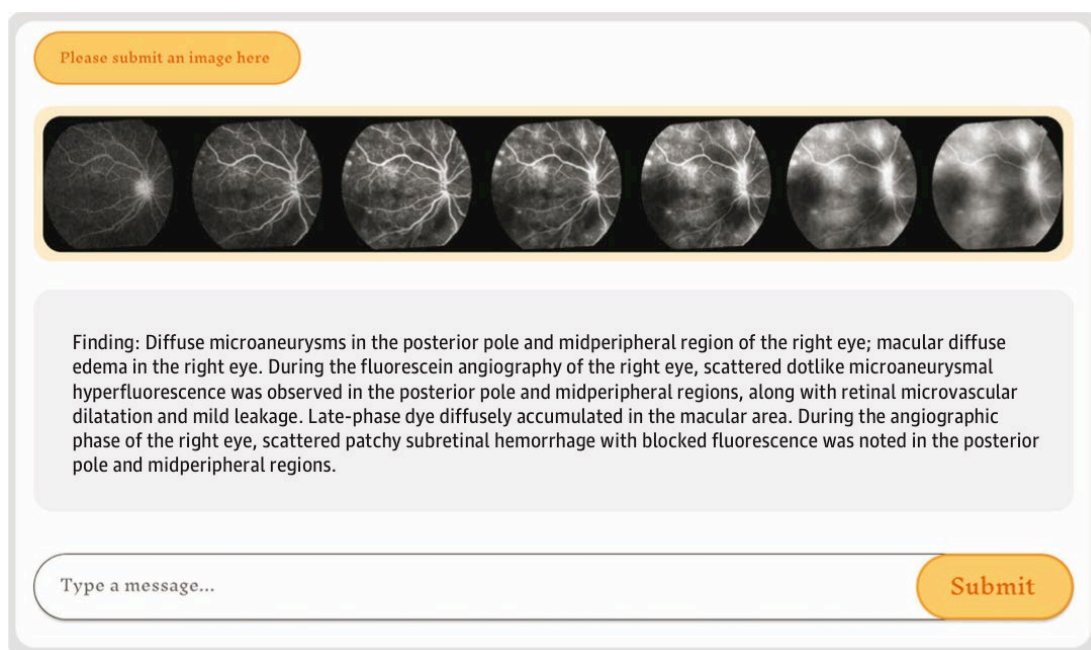


546 **Figure 2. Representative Examples of Real vs Artificial Intelligence–Generated Fundus**
 547 **Fluorescein Angiography (FFA) Videos.** BRVO = branch retinal vein occlusion, NPDR = non-
 548 proliferative diabetic retinopathy, AMD = age-related macular degeneration, CSC = central serous
 549 chorioretinopathy.
 550

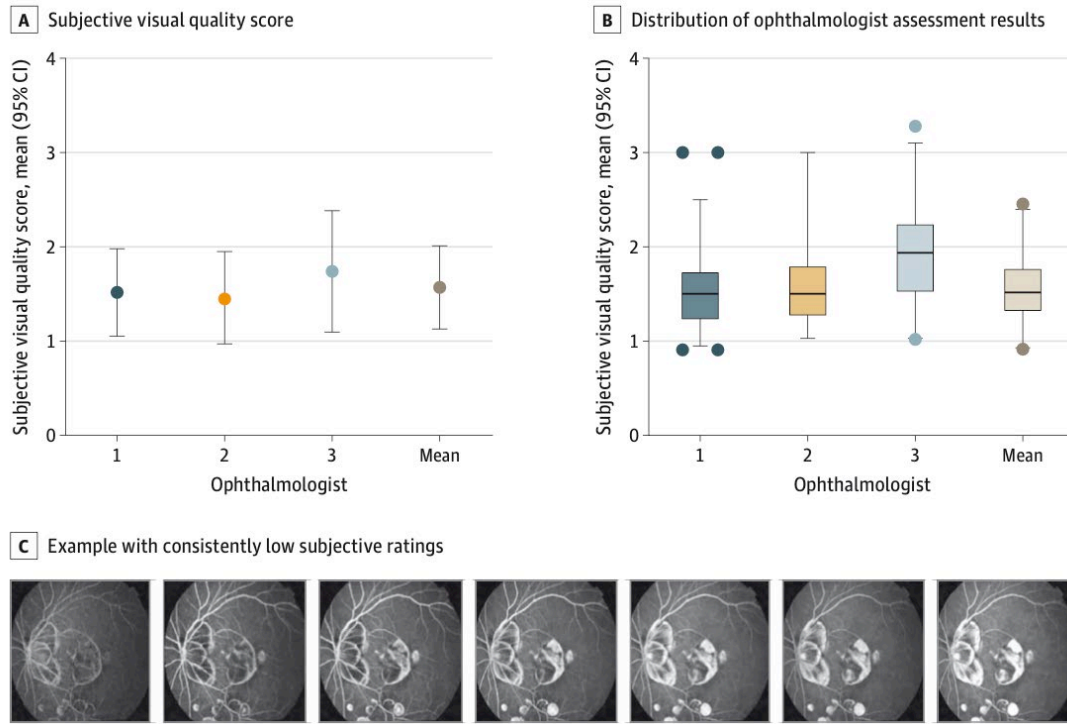


551 **Figure 3. Results From the Objective Evaluations of the Fundus Fluorescein Angiography**
 552 **(FFA) Video Quality.** BERTScore indicates bidirectional encoder representations from
 553 transformers score; FVD, Fréchet video distance; LPIPS, learned perceptual image patch similarity;
 554 VQAScore, visual question answering score. The details of the FFA text prompt are as follows:
 555 Depigmentation, pigmentary change, atrophy at right eye FFA: (1) Macular limited retinal pigment
 556 epithelial damage, fluorescein leakage and tissue hyperfluorescence staining. Macular limited
 557 retinal choroidal atrophy and pigmented spot deposition. Focal loss of retinal pigment epithelium
 558 pigmentation and pigmented deposition in the posterior pole. (2) Multifocal retinal pigment

559 epithelial depigmentation and pigmented deposition in the mid-peripheral and peripheral regions.
560 (3) No obvious abnormal fluorescence was observed in the optic disc. In B, the lower the Recall@K
561 value, the smaller the overlapping part between the artificial intelligence-generated FFA videos and
562 the real clinical videos, demonstrating the better performance in preserving patient privacy. In C,
563 the video quality rankings were poor (scores fell within the lowest 20%) across all parameters: the
564 scores were 0.53 for the LPIPS; 0.52 for the VQAScore; and 0.27 for the BERTScore.
565



566
567 **Figure 4. An example of domain-specific evaluation using a Fundus Fluorescein Angiography**
568 **Generative Pretrained Transformer (FFA-GPT).** This figure illustrates an example of the
569 domain-specific evaluation process. An FFA-GPT was used to implement a domain-specific
570 evaluation for determining whether the generated FFA videos accurately depict retinal disease
571 characteristics described in the given prompts



572

573 **Figure 5. Results From the Ophthalmologists' Assessments.** Three ophthalmologists (X.W.,

574 X.Y., and Y.F.) with 3, 5, and 10 years of experience, respectively (representing trainee,

575 attending, and expert levels), reviewed 50 randomly selected FFA videos generated from the

576 test set, comparing them against the corresponding real clinical videos. In A, the 3

577 ophthalmologists assessed video quality on a scale from 1 (excellent) to 5 (very poor). The

578 mean score was 1.52 (SD, 0.46) for the first ophthalmologist; 1.46 (SD, 0.49) for the second

579 ophthalmologist; and 1.74 (SD, 0.65) for the third ophthalmologist. The mean score was 1.57

580 (SD, 0.44) for all 3 ophthalmologists. In B, the first quartile was the 25th percentile of data; the

581 second quartile, the 50th percentile; and the third quartile, the 75th percentile. The whiskers

582 indicate the 2.5th and 97.5th percentiles, representing the range encompassing 95% of the data,

583 and the outliers (shown as filled circles) denote the values outside this range. In C, the artificial

584 intelligence-generated FFA video received consistently low subjective ratings due to its failure

585 to accurately generate the lesions described in the prompt. The details of the FFA text prompt

586 are as follows: Depigmentation, pigmentary change, atrophy at right eye FFA: (1) Macular
587 limited retinal pigment epithelial damage, fluorescein leakage, and tissue hyperfluorescence
588 staining. Macular limited retinal choroidal atrophy and pigmented spot deposition. Focal loss
589 of retinal pigment epithelium pigmentation and pigmented deposition in the posterior pole. (2)
590 Multifocal retinal pigment epithelial depigmentation and pigmented deposition in the
591 midperipheral and peripheral regions. (3) No obvious abnormal fluorescence was observed in
592 the optic disc.
593

594

Supplementary Documents

595 eFigure 1: Representative FFA videos of other retinal vascular diseases generated by FFA Sora.

596 eFigure 2: Representative FFA videos of other RPE-choroidal diseases generated by FFA Sora.

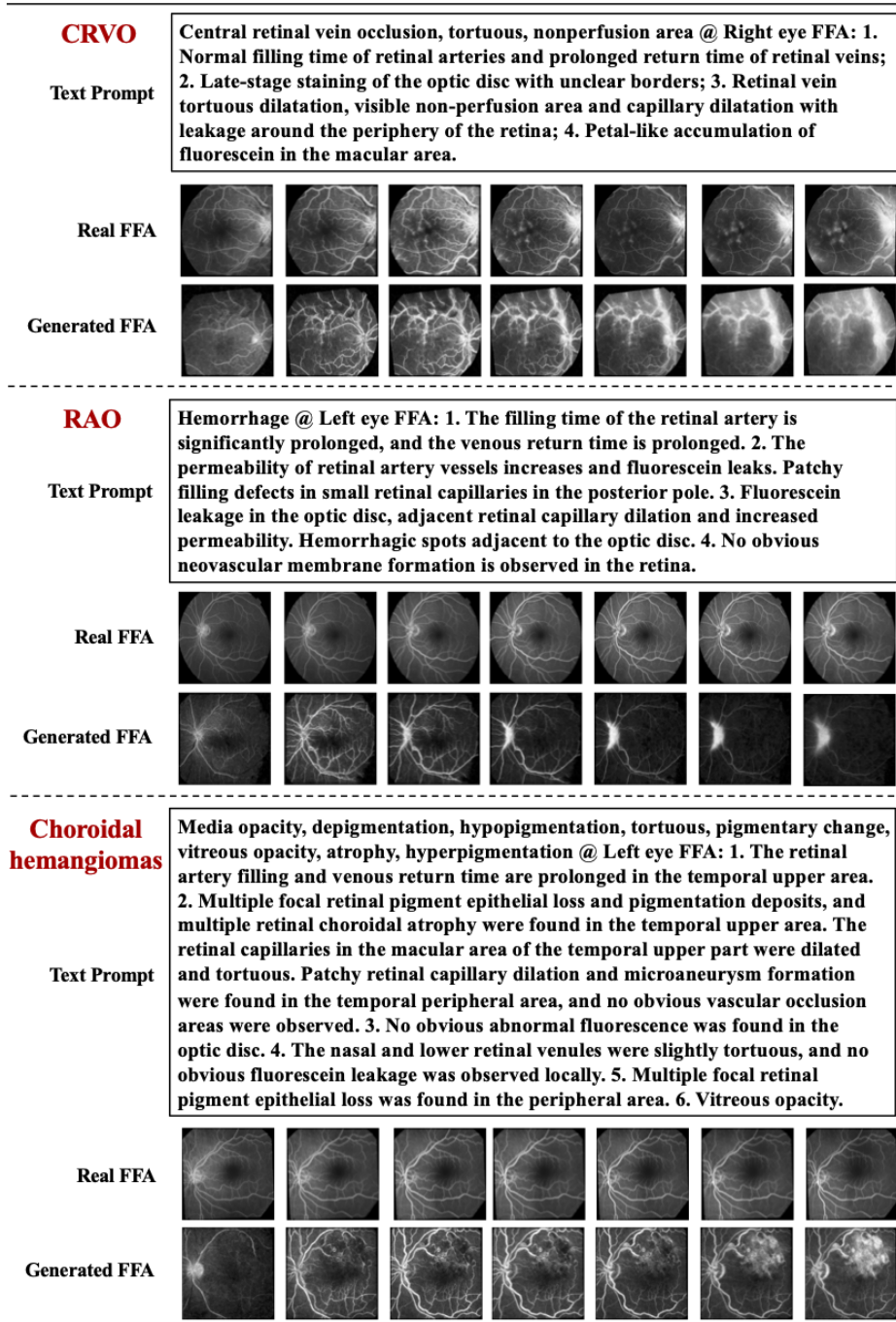
597 eFigure 3: Examples of **subjective evaluations** across levels (from 1 to 5).

598 eFigure 4: Examples of inconsistencies between objective and subjective evaluations

599 eTable 1: Demographic characteristics of the dataset.

600 eTable 2: The main eye conditions extracted from the fundus fluorescein angiography reports
601 (total N = 3625).

602

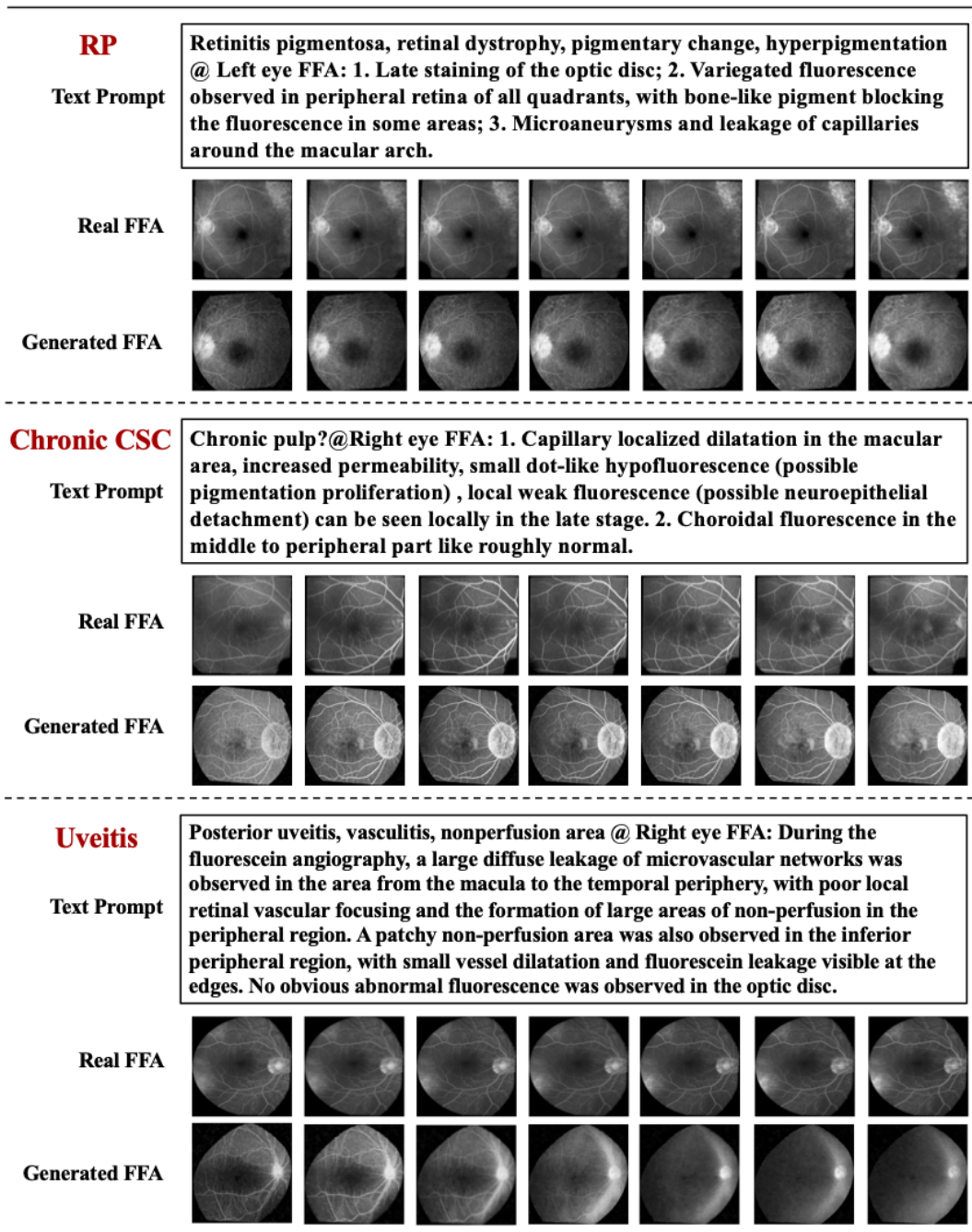


603

604 **eFigure 1: Representative FFA videos of other retinal vascular diseases generated by FFA Sora.**

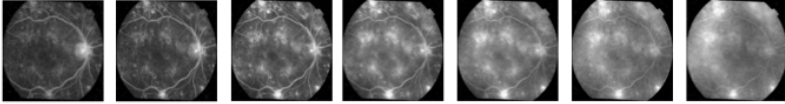

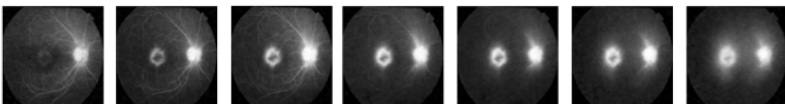
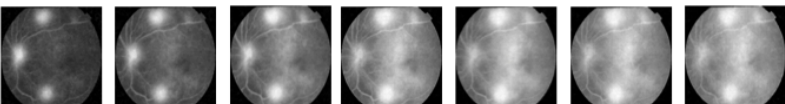
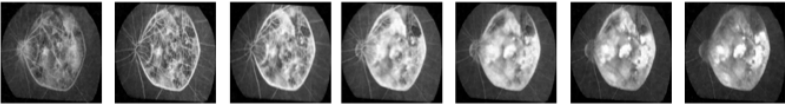
605 CRVO = central retinal vein occlusion, RAO = retinal artery occlusion.

606



607
608
609
610

eFigure 2: Representative FFA videos of other RPE-choroidal diseases generated by FFA Sora.
RP = retinal pigmentosa, CSC = central serous chorioretinopathy.

Rank 1	Hemorrhage, macular edema, depigmentation, laser scar, hypopigmentation, pigmentary change, edema, microaneurysm @ Right eye FFA: 1. Scattered retinal capillary dilatation and microaneurysm formation with increased segmental vascular permeability and leakage of fluorescein. 2. Scattered retinal capillary occlusion with no obvious neovascular membrane formation in the retina. 3. Macular retinal capillary dilatation, increased permeability, macular edema, and macular focal pigment epithelial damage. 4. Scattered hemorrhagic spots in the retina. 5. Staining of the optic disc. 6. Scattered multifocal pigment epithelial damage in the retina with tissue staining and pigment proliferation, scattered laser spots.
Text Prompt	
Generated FFA	
Rank 2	Fibroproliferation, traction, tortuous @ Right eye FFA: 1. Retinal artery perfusion and venous return time are normal. 2. The retinal vessels in the macula are tortuous, and fibrous proliferation membranes have formed in the macula. 3. There is no obvious fluorescein leakage in the mid-peripheral and peripheral retina. 4. The optic disc shows no significant abnormal fluorescence.
Text Prompt	
Generated FFA	
Rank 3	Media opacity, macular edema @ Right eye FFA: 1. Increased permeability of retinal capillaries and leakage of fluorescein. No obvious retinal vascular occlusion or neovascular membrane formation observed. 2. Increased permeability of macular retinal capillaries, fluorescein leakage, macular edema. 3. Enlarged optic disc and cup, relative filling defect at the disc edge, optic disc fluorescein leakage, strong fluorescein staining of the optic disc.
Text Prompt	
Generated FFA	
Rank 4	Microaneurysm, hemorrhage @ Left eye FFA: 1. The retina has scattered capillary dilatation, microvascular proliferation, increased vascular permeability, and leakage of fluorescein. 2. There is retinal capillary occlusion area without neovascularization membrane formation in the retina. Staining of the optic disc. 3. Macular retinal capillary dilatation and increased permeability. 4. The retina has scattered hemorrhagic spots.
Text Prompt	
Generated FFA	
Rank 5	Choroidal neovascularization, depigmentation, fibroproliferation, hypopigmentation, pigmentary change, atrophy, traction, hyperpigmentation @ Left eye FFA: 1. Limited macular retinal pigment epithelial pigment loss and pigment deposition, focal choroidal atrophy of the retina. Formation of choroidal neovascular fibrous proliferative membrane in macular lesions, tissue staining. 2. Disc staining, localized choroidal atrophy of the retina next to the optic disc. 3. Small-focal macular retinal pigment epithelial pigment loss in the middle and peripheral areas. Focal choroidal atrophy of the retina in the nasal and inferior peripheral areas.
Text Prompt	
Generated FFA	

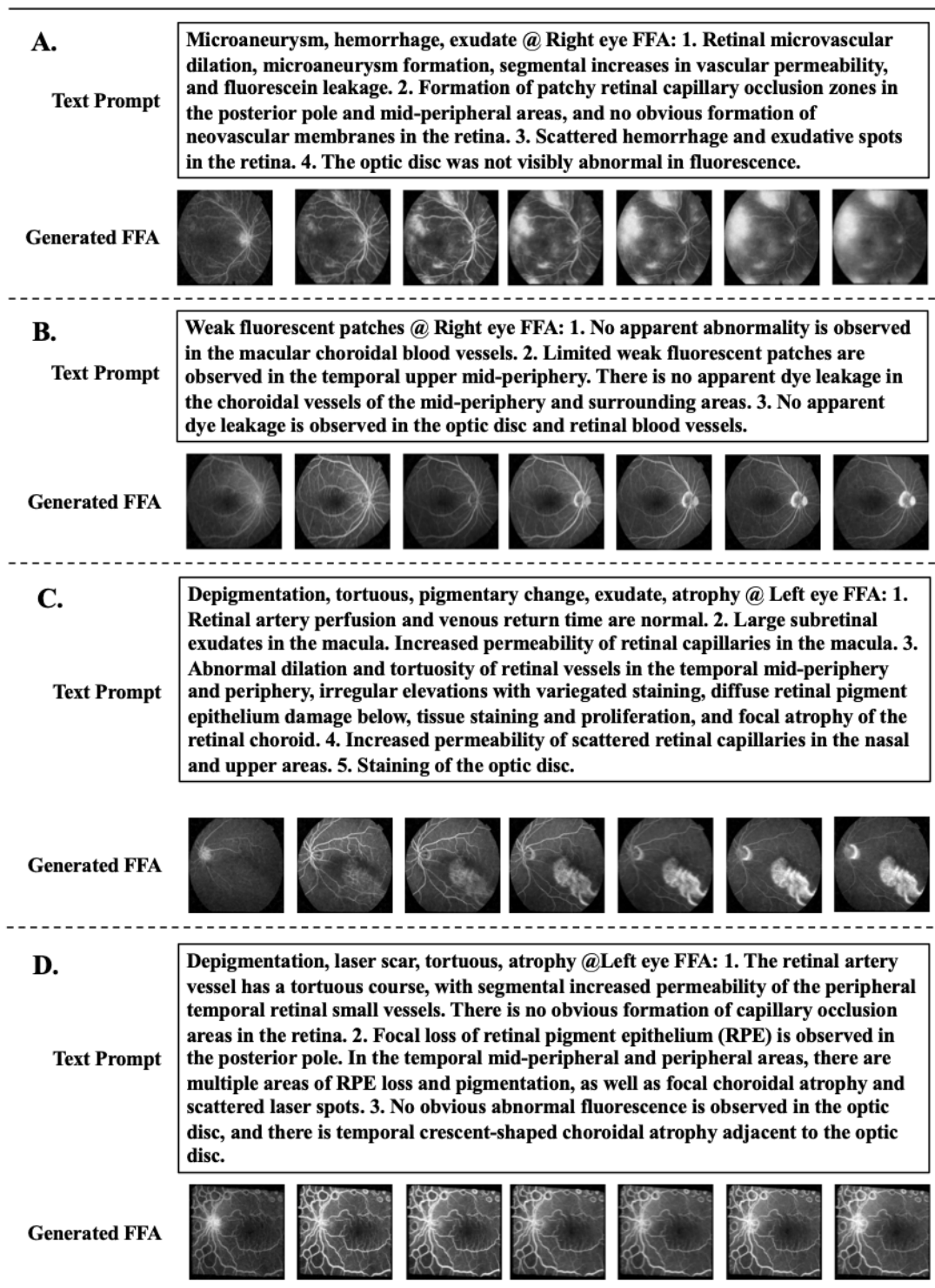
611

612 eFigure 2: Examples of subjective evaluations across levels (1 to 5). 1 = The retinal structure

613 and lesion characteristics of the generated videos exactly match the text prompts; 2 = The retinal

614 structure of the generated videos corresponds to the text prompts, and the lesion characteristics

615 are basically consistent with the text prompts. In this example, fibroproliferation, a key feature
616 in the text description, is missing from the generated images; 3 = The retinal structure of the
617 generated videos corresponds to the text prompt, and the lesion characteristics are slightly
618 consistent with the text prompts. In this example, the fluorescence leakage is overly pronounced,
619 and the optic disc structure appears indistinct, reducing the clarity of anatomical details; 4 =
620 The retinal structure of the generated videos corresponds to the text prompt, but the lesion
621 features cannot be generated. In this example, only microaneurysms are generated, and the
622 fluorescence signal is excessively strong, failing to reflect the hemorrhagic changes described
623 in the prompt; 5 = Unable to generate all text-oriented features.
624



625

626

627 eFigure 3: Examples of inconsistencies between objective and subjective evaluations. (A-B)

628 Examples where objective evaluation failed but subjective evaluation results were good. (C-D)

629

630

631

632

633 **eTable 1: Demographic characteristics of the dataset.**

	Total	Train	Validation	Test	P value
No.	3625	2851	387	387	
Age, median (IQR)	50 (36, 61)	50 (36, 61)	48 (34, 62)	50 (37, 62)	0.544
Sex, n (%)					0.685
Female	1636 (45.1)	1309 (45.9)	171 (44.2)	183 (47.3)	
Male	1931 (53.3)	1542 (54.1)	216 (55.8)	204 (52.7)	

634

635 **eTable 2: The main eye conditions extracted from the fundus fluorescein angiography reports**
 636 **(total N = 3625).**

<i>Conditions</i>	<i>N (%)</i>
Retinitis	2757(30.3%)
Leakage	2010(22.1%)
Neovascularization	720(7.9%)
Atrophy	673(7.4%)
Microaneurysms	670(7.4%)
Edema	654(7.2%)
Hemorrhage	571(6.3%)
Retinopathy	511(5.6%)
Exudates	166(1.8%)
Uveitis	103(1.1%)
Drusen	86(0.9%)
Chorioretinitis	73(0.8%)
Ischemia	38(0.4%)
Scleritis	24(0.3%)
Vein Sheathing	20(0.2%)
Fibrosis	7(0.1%)
Foveal Cyst	3(0.0%)

637

638

639

640

641