

# A Concept-based Interpretable Model for the Diagnosis of Choroid Neoplasias using Multimodal Data

Yifan Wu<sup>1†</sup>, Yang Liu<sup>2†</sup>, Yue Yang<sup>1</sup>, Michael S. Yao<sup>1</sup>, Wenli Yang<sup>3</sup>, Xuehui Shi<sup>3</sup>,  
Lihong Yang<sup>3</sup>, Dongjun Li<sup>3</sup>, Yueming Liu<sup>3</sup>, James C. Gee<sup>1</sup>, Xuan Yang<sup>3\*</sup>,  
Wenbin Wei<sup>3\*</sup>, Shi Gu<sup>2\*</sup>

<sup>1</sup>University of Pennsylvania, Philadelphia, PA, US.

<sup>2</sup>University of Electronic Science and Technology, Chengdu, China.

<sup>3</sup>Beijing Tongren Eye Center, Beijing Key Laboratory of Intraocular Tumor Diagnosis and  
Treatment, Beijing Ophthalmology and Visual Sciences Key Lab, Beijing Tongren Hospital,  
Capital Medical University, Beijing, China.

\*Corresponding author(s). E-mail(s): [yangxuan153@126.com](mailto:yangxuan153@126.com); [weiwenbintr@163.com](mailto:weiwenbintr@163.com);  
[gus@uestc.edu.cn](mailto:gus@uestc.edu.cn);

†These authors contributed equally to this work.

## Abstract

Diagnosing rare diseases presents a common challenge in clinical practice, necessitating the expertise of specialists for accurate performance. While the advent of machine learning offers potentially promising solutions, developing such technologies for rare disease diagnosis is hindered by a scarcity of data on rare conditions and the unmet demand for models that are both interpretable and trustworthy for clinical usage. Interpretable AI, with its capacity for producing human-readable outputs, can facilitate clinician validation, and contribute to the education of junior clinicians through exposure to a broad spectrum of cases. In this, we investigate how interpretable AI methods can be leveraged for rare disease diagnosis. We focus on choroid neoplasias, the most prevalent form of eye cancer in adults, albeit with a low prevalence of 5.1 per million. We build the largest multimodal dataset to date of choroid neoplasm imaging data from over 750 patients collected between 2004 to 2022. Furthermore, we introduce a multimodal concept-based interpretable model (MMCBM) that distinguishes between three types of choroidal tumors, integrating insights from domain experts via radiological reports. Our model not only achieves an  $F_1$  score of 0.91, rivaling that of black-box models, but also boosts the diagnostic accuracy of physicians by 42%. This study highlights the significant potential of interpretable machine learning in improving the diagnosis of rare diseases, laying a groundwork for future breakthroughs in medical AI that could tackle a wider array of complex health scenarios.

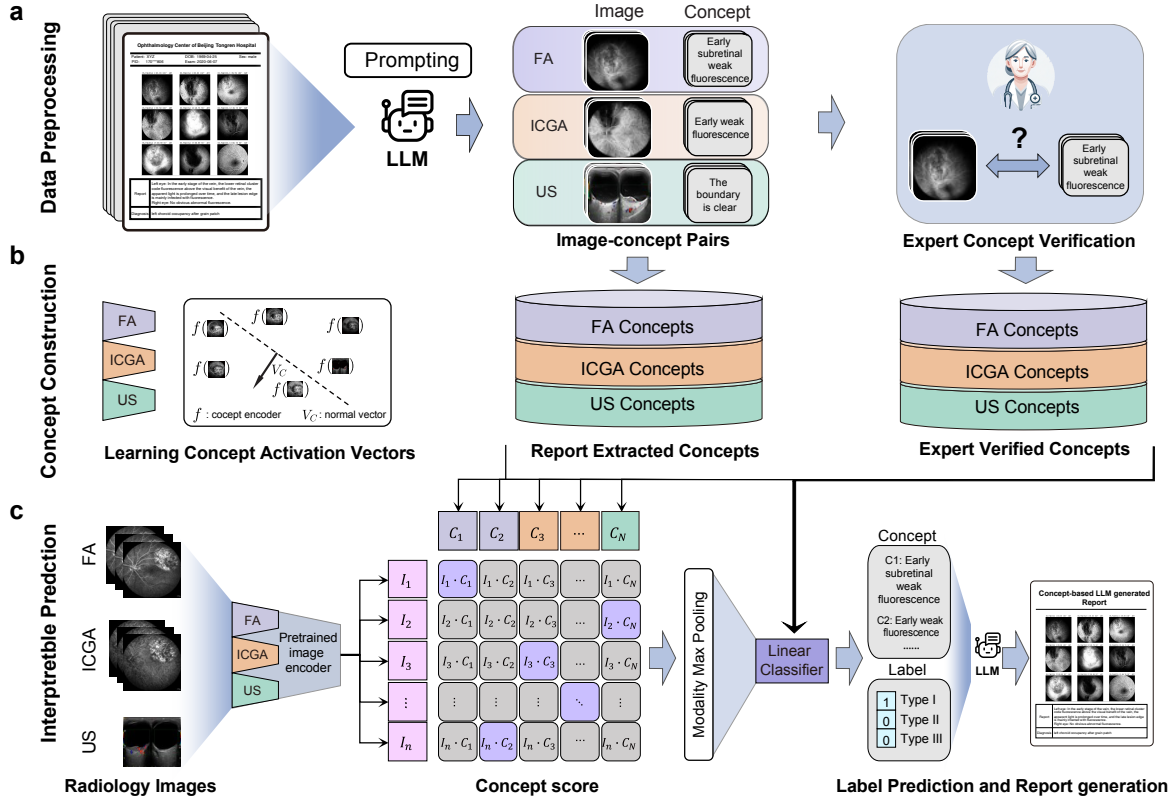
**Keywords:** Uveal Melanoma, Computer-aided Diagnosis, Interpretable Machine Learning, Multi-modality Classification, Concept Bottleneck Model, Rare Disease Diagnostics

## 33 Introduction

34 Recent advancements in machine learning and deep neural networks have accelerated the development  
35 of computer-aided diagnosis (CAD) methods in the past decade [1]. For common diseases amenable to  
36 automated diagnoses with large publicly available datasets, deep learning-based models have performed com-  
37 parably to radiologists across a variety of diagnostic tasks. Example applications include analyzing chest  
38 X-rays (CXRs) [2, 3], fundus photography for automated retinopathy screening [4], and brain MRIs for  
39 tumor and stroke lesion quantification [5, 6]. However, for many diseases, especially rare diseases, there is  
40 often a lack of high-quality datasets to power such learning-based approaches due to the prohibitive cost  
41 of professional annotations and potential incompatibilities between clinical and research protocols. In addi-  
42 tion, high-stakes diagnostic tasks for diseases associated with high patient morbidity and mortality crucially  
43 require *interpretable* machine-generated predictions for easy downstream verification by specialist clinicians  
44 with relevant domain expertise [7–9]. Addressing such challenges is therefore crucial for developing CAD  
45 technologies for rare disease identification.

46 In the workup of rare oncologic diseases, physicians use multimodal imaging biomarkers collected from  
47 different acquisition methods to provide a comprehensive diagnosis for patients. The development of an  
48 applicable CAD model should also employ such a multi-modal pipeline. Recent work proposes implementing  
49 these pipelines by aligning medical imaging data to text-based context descriptions [10, 11], which leverages  
50 the representation power of foundation models such as Contrastive Language-Image Pre-training (CLIP) [12]  
51 and generative Large Language Models (LLMs) [13]. While these methods may help augment the clinical  
52 workflows in many common diseases, their applications to rare disease diagnostic tasks are largely unexplored.  
53 Notably, when attempting to diagnose rare diseases with artificial intelligence (AI) in particular, there is a  
54 need for well-curated data and feasible tools to generate explainable biomarkers aligned with experienced  
55 specialists. High-quality data with labeled image-text pairs are necessary to train sophisticated AI models  
56 with sufficient diagnostic performance. Because of the low prevalence of rare diseases, the analytical tools  
57 also need to produce interpretable explanations as well as accurate predictions to comprehensively facilitate  
58 the clinical management [14–16]. These practical constraints for diagnosing rare oncologic diseases call for a  
59 machine-learning paradigm distinguished from existing performance-focused approaches.

60 In this work, we discuss our approach on engineering machine learning model architectures specifically  
61 designed for rare disease diagnosis. We focus on the diagnosis of uveal melanoma, a rare cancer originating  
62 from the iris, ciliary body, choroid, or other components of the uveal tract in the eye [17, 18]. While cases  
63 of uveal melanoma are rare, with an estimated incidence of 5.1 per million in the United States, the long-  
64 term prognosis is poor due to the high risk of metastasis at the time of diagnosis [19]. Uveal melanomas are



**Fig. 1: Overview of the MMCBM workflow.** (a) Utilizing a large language model (LLM), concept banks are formulated by extracting image-concept pairs from comprehensive medical reports. Senior experts help examine the faithfulness of the image-concept pairs and make corresponding modifications. (b) Based on such pairs, we construct the concept bank by learning concept activation vectors. (c) The model’s output stage takes a series of images spanning 1 to 3 modalities. A pre-trained image encoder is employed to convert these images into tokenized features. Subsequent calculations produce concept scores. The model then delivers an explainable prediction, spotlighting the diagnostic evidence. Moreover, it crafts an interpretative report, enhancing the transparency of the diagnostic process.

65 frequently missed in routine clinical workups due to their low prevalence in the general population. As a result,  
 66 few clinicians are well-trained in their diagnosis and clinical management [20]. Initial diagnosis requires a  
 67 detailed fundoscopic examination with an expert clinician followed by additional advanced imaging techniques  
 68 such as ocular ultrasound (US), fluorescein angiography (FA), and indocyanine green angiography (ICGA)  
 69 for confirmation and prognostication [21–24]. Domain-specific expert physicians specializing in managing  
 70 uveal melanomas are few and far between, further complicating diagnostic workup [23]. To overcome these  
 71 challenges, we aim to build a computer-aided system to differentiate between choroidal melanoma, metastatic  
 72 carcinoma, and hemangioma—, all occurring in the choroid of the fundus and often appearing as solitary  
 73 tumors. These diseases may have similar symptoms in the early stages and overlapping imaging features [25].  
 74 Given the poor prognosis associated with uveal melanomas and the consequent need for timely diagnosis

75 and treatment, it is crucial to have high confidence in a diagnosis of choroidal neoplasias prior to definitive  
76 intervention.

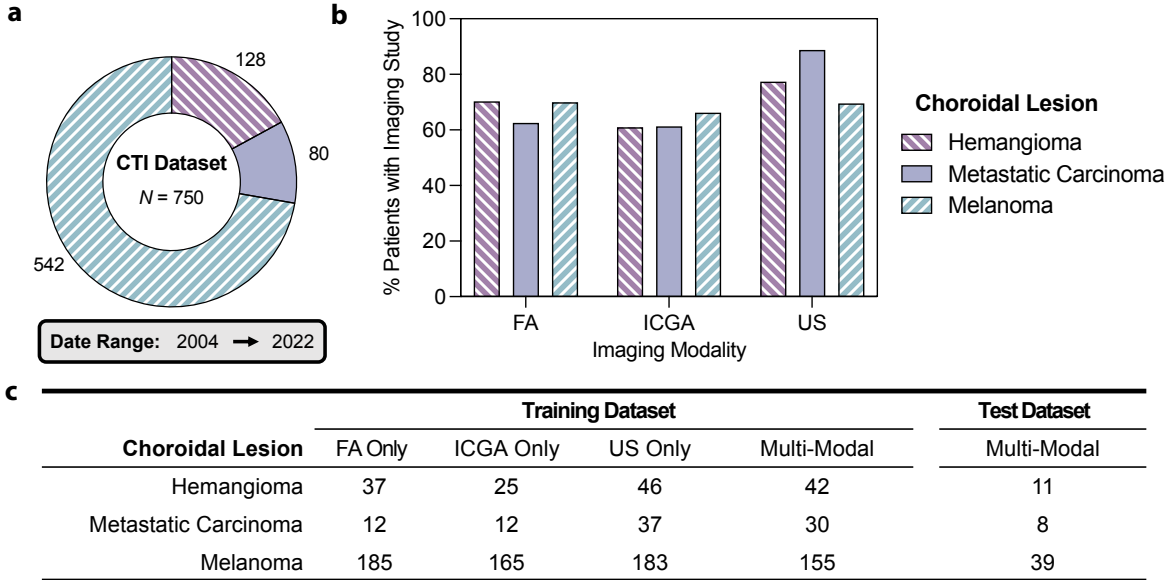
77 To establish such a pipeline for the automated interpretable diagnosis of choroidal neoplasias, we need  
78 to address key challenges in data curation and model training and verification. First, to enable our work,  
79 we collect the first well-curated multimodal dataset of uveal oncologic pathologies to train classifier models  
80 that accurately differentiate choroidal melanomas from other clinically similar diseases. To our knowledge,  
81 this dataset is the largest one for choroidal melanomas in the world. We then use this dataset to develop  
82 the Multimodal Medical Concept Bottleneck Model (MMCBM), a domain knowledge-enhanced model that  
83 predicts interpretable classifications from patient data. MMCBM supports a human-in-the-loop mechanism  
84 to learn from feedback provided by domain experts. We find that MMCBMs not only provide accurate  
85 classifications, but also offer interpretable concepts that explain its reasoning process. The concepts align  
86 well with senior doctors and provide substantial assistance for the junior ones to more accurately diagnose  
87 choroidal neoplasias. Our methodology leverages the extensive knowledge in clinical reports to offer a pathway  
88 towards building interpretable models for diagnosing rare diseases.

## 89 **Results**

### 90 **Dataset Description**

91 To support the development of interpretable models for diagnosing choroidal tumors, we built the Choroid  
92 Tri-Modal Imaging (CTI) dataset, a anonymized, multimodal, and annotated collection of medical images  
93 from Beijing Tongren Hospital (2004-2022) encompassing Fluorescence Angiography (FA), Indocyanine  
94 Green Angiography (ICGA), and Ocular Ultrasound (US) images. Construction of this dataset was approved  
95 by the Ethics Committee of Beijing Tongren Hospital. CTI includes images from patients diagnosed with  
96 benign hemangioma, secondary metastatic carcinoma to the eye, or primary choroidal melanoma. The CTI  
97 dataset (Fig. 2) consists of 542 patients with choroidal melanoma (FA: 379, ICGA: 359, US: 377), 128  
98 patients with choroidal hemangioma (FA: 90, ICGA: 78, US: 99), and 80 patients with choroidal metastatic  
99 carcinoma (FA: 50, ICGA: 49, US: 71). The numbers indicate the quantity of imaging studies for each  
100 specific imaging modality. Note that not every patient has images across all modalities. We refer to the  
101 subset where patients have all three modalities as Multi-Modal (MM) data and reserve 20% of this MM  
102 data as a hold-out test set. In the MM data training split, 97 patients have anonymized reports for all three  
103 modalities, describing the radiological features observed in the images.

104



**Fig. 2: Statistics of the CTI Dataset.** (a) The CTI dataset is composed of 750 patients: 542 with melanoma, 128 with hemangioma, and 80 with metastatic carcinoma, collected from 2004 to 2022. (b) Proportions of patients with hemangioma, metastatic carcinoma, and melanoma imaged by Fluorescein Angiography (FA), Indocyanine Green Angiography (ICGA), and Ultrasound (US). (c) Split of imaging studies in the training and test datasets across various imaging modalities: 20% of the Multi-Modal data (MM), representing patients imaged with all three modalities, is set aside for testing. The remaining 80% of MM and all non-MM data are allocated for training using 5-fold cross-validation.

## 105 Baseline Blackbox Model

106 We first sought to build baseline black-box machine learning models. Taking inspiration from recent success  
 107 in natural image processing [26], our baseline black-box model is composed of three separate modality-  
 108 specific encoders trained to encode corresponding imaging study inputs into intermediate lower-dimensional  
 109 representations. The encoder output (or *outputs*, if multiple imaging studies of different modalities are  
 110 available for a given patient) is then passed to an attention pooling block [27] and subsequent dense layer to  
 111 yield the final classification prediction. We refer to this model architecture as the *Pre-Trained Multimodal*  
 112 *Classifier*. Our baseline model performs accurately across different input image modalities, validating the  
 113 feasibility of deep-learning models for this clinical problem. Using FA imaging studies alone, the pre-trained  
 114 classifier achieved an  $F_1$  score of 78.3% (95% CI: 74.0 - 81.7%); using ICGA studies alone, it achieved an  
 115  $F_1$  score of 85.9% (95% CI: 83.7 - 88.2%); and using US studies alone, it achieved an  $F_1$  score of 72.1%  
 116 (95% CI: 67.1 - 76.7%). When using all three imaging studies together, the baseline classifier attained  
 117 an  $F_1$  score of 89.2% (95% CI: 87.9 - 90.6%). Additional classification are included in **Supplementary**  
 118 **Table A1**. Our results show that using a multimodal inputs lead to models that are more accurate than  
 119 those leveraging any individual imaging study as input alone. However, while the Pre-Trained Multimodal  
 120 Classifier demonstrates impressive performance, it is impossible to interrogate the model’s predictions for

121 human experts to interpret —a key limitation of existing black-box approaches.

122

### 123 **Trustworthy Interpretable Framework: MMCBMs**

124 The lack of interpretability in the baseline pre-trained classification model is a common trait of many  
125 modern AI tools. To address the need for trustworthiness in medical diagnostics, we sought to engineer a  
126 framework with interpretability baked into the model design. Our approach, referred to as the Multimodal  
127 Medical Concept Bottleneck Models (MMCBMs), is a task-agnostic framework designed for high-stakes  
128 applications where human-in-the-loop or subsequent human verification is critical. Our key insight is to  
129 leverage prior knowledge from domain experts to align the intermediate representations of input images by  
130 the model as representations that human experts can easily understand. In this way, model predictions can  
131 be easily interpreted as activations and linear combinations of these representations. These representations  
132 can be visual patterns or findings that clinical experts consider the evidence for making diagnoses and are  
133 used as educational guidelines in a natural language format. We refer to these representations as *concepts*.

134

### 135 **Concept Construction and Grounding**

136 Using medical reports as the knowledge database, we prompt GPT-4 [28] to extract concepts from reports  
137 and construct a bank of concepts containing phrases related to imaging findings of choroidal tumors. For  
138 instance, a description in a fluorescein angiography (FA) report states, “In the venous phase, a clustered  
139 hypofluorescence under the subretinal can be seen in the temporal part of the macula. Fluorescence increases  
140 with time, and lesions are dominated by fluorescent staining at the late stage.” The extracted concepts for  
141 this FA study include “Clustered Hypofluorescence During Venous Phase”, “Globally Increasing Fluores-  
142 cence Intensity”, and “Late-Stage Staining.” After extracting concepts from the reports of 97 patients, we use  
143 GPT-4 to aggregate semantically similar concepts, ensuring each concept’s uniqueness and relevance. The  
144 final concept bank consists of 47 concepts for FA, 30 for ICGA, and 26 for US, with an average of 3 concepts  
145 for FA, 2 for ICGA, and 5 for US per patient. The comprehensive list of all  $N = 103$  concepts is presented  
146 in **Supplementary Table A2**. To validate that the concepts extracted by the LLM accurately represented  
147 real-world clinical reasoning, two senior ophthalmologists specializing in diagnosing and managing choroidal  
148 tumors at Beijing Tongren Hospital were asked to verify and amend the concepts. Quantitatively, the initial  
149 concept bank constructed by GPT-4 was assessed to be reasonable and relevant, requiring only minor modi-  
150 fications: 5 concepts were removed, and 8 new ones were added to the FA category, 4 to the ICGA category,  
151 and no changes to the US category.

152 To ground concepts as feature embeddings, we employed support vector machines (SVMs) for concept-  
153 level binary classification. We used image representations from a pre-trained model as input and binary  
154 derived from the concept construction process as labels. Images associated with assigned concepts were  
155 used as positive samples and all other images were used as negative samples. The classification hyperplane  
156 vector from each SVM serves as the concept’s representation, which we refer to as concept activation vec-  
157 tors (CAVs) [29]. Subsequently, in MMCBM, an image is projected into the space of concepts to estimate  
158 the input image alignment with any given modality-specific concept. The alignment scores are then used  
159 as input into a linear classifier to predict the relative probabilities of each of the three targeted choroidal  
160 diseases. **Fig. 3a** shows this process and shows the top- $k$  concepts derived from concept scores to explain  
161 the model’s predictions.

162

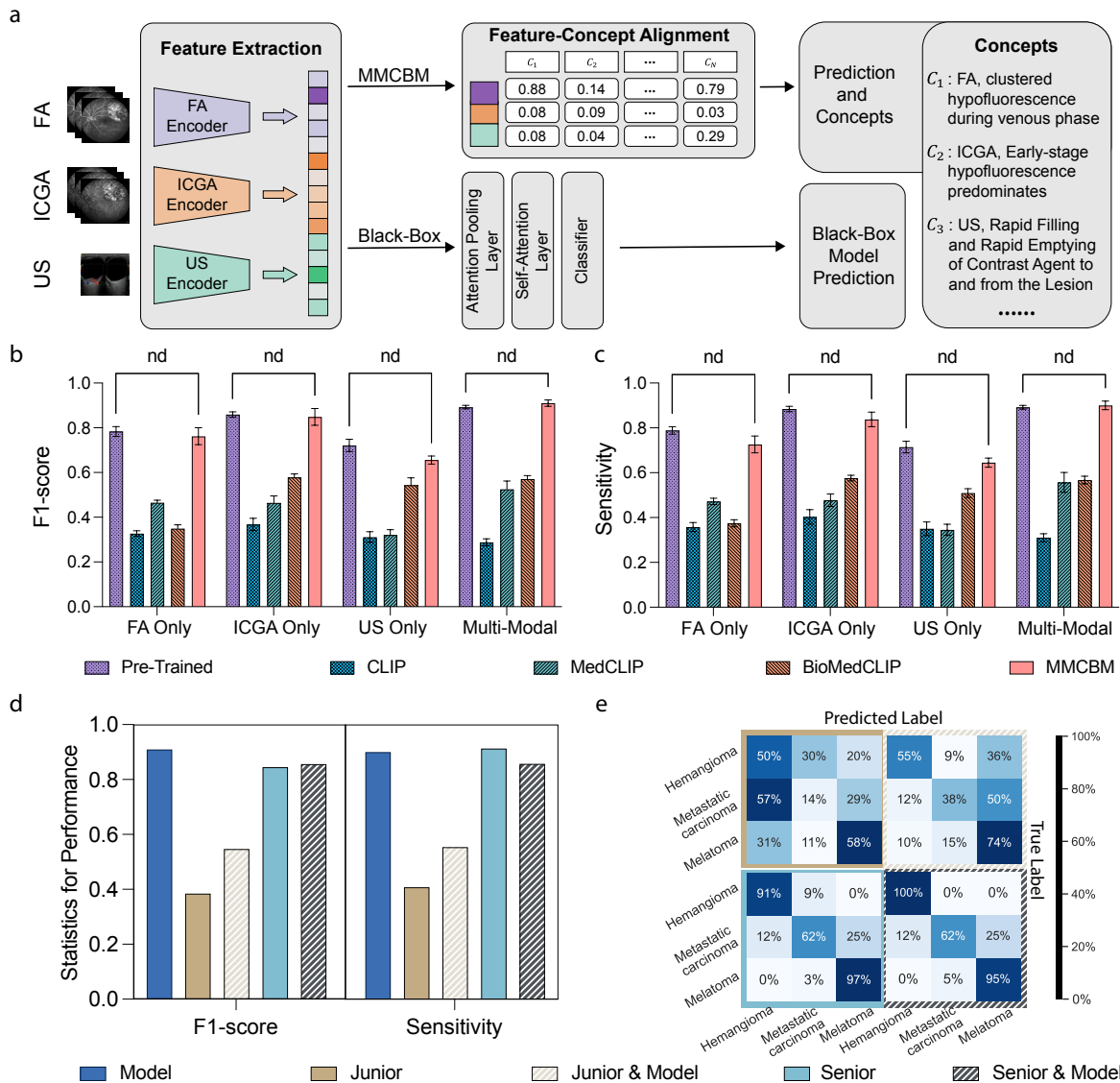
### 163 **Noninferior Accuracy of MMCBM to Black-Box Model**

164 A common critique of interpretable machine learning models is that enforcing priors on the model, such as  
165 requiring input images to align with concept activation vectors, is equivalent to adding additional regulariza-  
166 tion to the hypothesis space [30]. Such constraints may adversely impact the performance of trained models  
167 [31, 32]. To this end, we sought to evaluate the classification performance of our MMCBM model against  
168 the black-box pre-trained multimodal classifier baseline (**Fig. 3**). On the MM testing dataset, MMCBM  
169 achieved an overall classification  $F_1$  score of 91.0% (95% CI: 88.2 - 93.4%), which is comparable with the  
170 performance of the baseline black-box model (89.2%; 95% CI: 87.9 - 90.6%). These results show that the  
171 MMCBM framework is non-inferior to the black-box pre-trained model. We attribute this improvement to  
172 the fact that by adding interpretable regularization, the framework mitigates the issue of class imbalance  
173 in the data. Additionally, comparing classifier performance across unimodal imaging inputs revealed no  
174 statistically significant differences in classification metrics (**Table A1**). This indicates that our MMCBM  
175 framework matches the performance of black-box approaches in automating the diagnosis of rare choroidal  
176 tumors according to clinically relevant metrics.

177

### 178 **Integration of MMCBM in Clinical Workflows**

179 We have shown that the MMCBM effectively leverages prior knowledge from domain experts to represent  
180 input data aligned with interpretable concepts. However, it remains unknown whether our framework  
181 can provide real-world utility in augmenting existing clinical workflows. To investigate the applications of  
182 MMCBM in clinical practice, we recruited the help of 8 doctors from Beijing Tongren Hospital: 2 senior  
183 ophthalmologists specializing in the diagnosis and management of choroidal melanomas, and 6 resident



**Fig. 3: Multimodal Medical Concept Bottleneck Model (MMCBM).** Black-box models, such as the pre-trained classifier, learn directly from the encoded image features and output a single model prediction without any insight as to how the prediction was computed. In contrast, the MMCBM shown in (a) instead represents encoded features by their alignment with key medical concepts derived from domain experts. This allows MMCBM to return not only its prediction but also the top- $k$  activated concepts that best describe the input images, giving insight into how the model arrived at its prediction. Comparing both the classification (b) accuracy and (c) sensitivity of the models, there is no statistically significant difference between black-box models and MMCBM across all sets of imaging inputs. MMCBM concepts also outperform features derived from CLIP-based models, highlighting the importance of sourcing prior knowledge from domain experts. (d) Performance Benchmark with Human Evaluators: A comparison of our model’s performance against junior and senior doctors. After presenting them with the model’s predicted concepts, they conducted a subsequent assessment, enabling us to document and compare performance metrics. (e) Confusion Matrix for Human Evaluators with and without Concepts. The matrices correspond to Junior, Junior & Model, Senior, and Senior & Model groups.



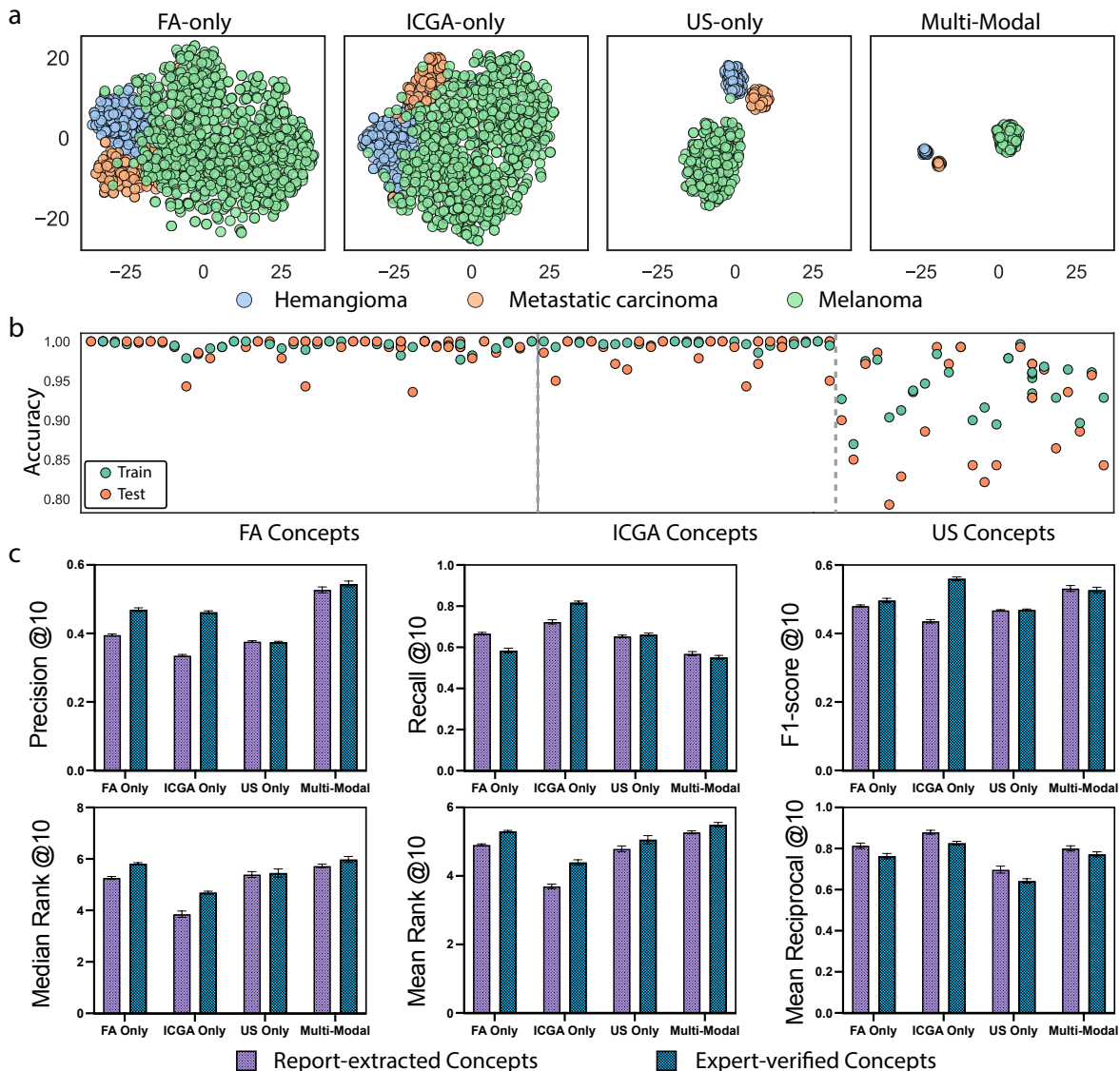
184 ophthalmologists in training. We assessed the diagnostic performance of ophthalmologists alone against  
185 ophthalmologists with our trained MMCBM model. In order to avoid the memorization of the seen cases,  
186 the time lapse between the experiments w./w.o. the generated concepts is 2 months. The ophthalmologists  
187 leveraging our MMCBM model for diagnostic workflow augmentation have access to the top-10 activated  
188 concepts from the MMCBM concept bank and can adjust the confidence scores of the concepts based on  
189 their judgment. This human-in-the-loop interactive feature improves the practical utility of MMCBM in  
190 clinical decision-making, fostering a more collaborative and accurate diagnostic process. For the 6 junior  
191 ophthalmologists, the average accuracy is 51.9%, precision 40.5%, recall 40.9%, and  $F_1$  score 38.5%; with  
192 the aid of our MMCBM model, their accuracy improves to 65.5%, precision to 54.3%, recall to 55.5%, and  
193  $F_1$  score to 54.7% (**Fig. 3[d-e]**). The 2 senior ophthalmologists demonstrate a high diagnostic accuracy at  
194 baseline of 91.4%, precision of 85.8%, recall of 83.6%, and  $F_1$  score of 84.6%. When augmented with the  
195 model’s predictions, their performance remains relatively unchanged with an accuracy of 91.4%, precision  
196 of 86%, recall of 85.8%, and  $F_1$  score of 85.7%. In particular, the use of MMCBM improves junior doctors’  
197 performance by 42% on the  $F_1$  score. These results not only validate the quality and precision of the pre-  
198 dicted concepts of our MMCBM model but also highlight our model’s ability to serve educational purposes  
199 by improving the diagnostic accuracy of less experienced doctors for complex and rare diseases.

200

### 201 **Comparison between MMCBM and Alternative Feature Embedding Methods**

202 Given the recent progression of cross-modality foundation models, it may be possible to leverage existing  
203 feature embedding models trained on extensive corpora of medical information to represent input ocular  
204 imaging data and concepts. This approach might offer greater generalizability and require less effort than  
205 our MMCBM setup. To evaluate this alternative framework, we compared our concept embedding procedure  
206 and image feature extraction with those using Contrastive Language-Image Pre-training (CLIP) [12] and  
207 its biomedical variants, including MedCLIP [33] and BioMedCLIP [34], which are specifically fine-tuned  
208 for medical data. Briefly, MedCLIP was fine-tuned on multiple Chest X-ray datasets, while BioMedCLIP  
209 underwent fine-tuning on 15 million figure-caption pairs extracted from biomedical research articles in  
210 PubMed Central. Our results suggest that all assessed CLIP-based frameworks perform significantly worse  
211 than our CAV-based feature extraction method used in our MMCBM framework (**Fig. 3[b-c]**). As expected,  
212 methods fine-tuned on specialized medical datasets—such as MedCLIP and BioMedCLIP—outperform  
213 the generic CLIP model as feature extractors for choroidal disease diagnosis using both multimodal and  
214 unimodal image inputs (**Fig. 3[b-c]**, MedCLIP: 52.5% (95% CI: 47.2 - 59.6%), BioMedCLIP: 57.2% (95%  
215 CI: 54.7 - 59.6%), CLIP: 28.8% (95% CI: 26.2 - 31.3%)). The analysis of unimodal input results and

216 additional classification metrics further aligns with these findings. Specifically, embedding model inputs  
 217 with expertise-curated knowledge significantly outperforms the use of general domain knowledge. These  
 218 observations highlight the necessity for fine-tuning and domain-specific adaptation or embedding images  
 219 and texts in medical applications. Furthermore, they affirm the efficacy of our MMCBM as a viable and  
 220 effective means to achieve model interpretability without compromising algorithmic performance.  
 221



**Fig. 4: Comparative Human Evaluation and Model Insights.** (a) Embedding Visualizations via t-SNE: This offers a graphical representation of embeddings from the trio of pretrained encoders. Notably, the fused MM embeddings are processed through the attention-pooling mechanism. (b) Accuracy of SVMs in generating concept banks using Concept Activation Vectors (CAVs). (c) Metrics of predicted Top- $k$  concepts on test dataset with  $k = 10$ . This evaluation includes precision@ $k$ , recall@ $k$ , and F1@ $k$ , as well as mean rank@ $k$ , median rank@ $k$ , and mean reciprocal rank@ $k$ .

## 222 Evaluation of Image-Concept Alignment

223 Our MMCBM demonstrates classification performance on par with state-of-the-art black-box models and  
224 offers interpretable insights into final model outputs. We have also shown that the quality of model inter-  
225 pretability depends on the quality of (1) the prior knowledge used to construct the MMCBM concept bank;  
226 (2) the image and concept embedding functions, and (3) image-concept alignment. We sought to evaluate  
227 our model’s interpretability according to these three aspects.

228 First, we evaluated the MMCBM feature representations and their accuracy in describing input images.  
229 Model representations for each of FA, ICGA, and US imaging studies were computed by the respective  
230 MMCBM encoders before leveraging t-SNE [35] dimensionality reduction techniques to visualize the complex  
231 feature landscapes from our multimodal dataset (**Fig. 4a**). We observe distinct clusters corresponding to,  
232 hemangioma, metastatic carcinoma, and melanoma, indicating effective class separation by the MMCBM  
233 encoders. Qualitatively, the clusters corresponding to multimodal data inputs appear more cohesive and  
234 less dispersed, suggesting that integrating multi-modal inputs may improve the separability of the different  
235 class representations in this representation space. This enhanced clustering density may contribute to the  
236 improved discriminative performance of our multimodal MMCBM models in contrast to models with only  
237 unimodal inputs accessible.

238 Next, we evaluated the quality of the MMCBM concept representation and image-concept alignment  
239 by examining the accuracy of the SVM classifiers employed in generating concept vectors for each medical  
240 concept. A high SVM accuracy score indicates a concept’s representational effectiveness and consistent  
241 presence across the dataset. According to this metric, FA and ICGA concepts achieve high accuracy across  
242 the board (**Fig. 4b**), with accuracy on test data exceeding 90% for all concepts. This suggests that concepts  
243 derived from FA and ICGA are well-represented and aligned with the input images. In contrast, though  
244 less accurate, the accuracy scores for US-based concepts are still higher than 80% for all concepts. This  
245 suggests that classifying diseases from ultrasound images alone may be more challenging. Specific details of  
246 the individual concepts and their corresponding accuracies are detailed in **Supplementary Table A2**.

247 To further assess the quality of MMCBM concept-based interpretability, we examined how well the model  
248 concepts align with ophthalmologist annotations. We selected the top- $k$  concepts predicted by MMCBM  
249 for each patient in the multimodal testing dataset. In **Fig. 4c**, we quantify our model’s alignment with  
250 expert annotations according to key performance metrics: Precision@ $k$ , Recall@ $k$ , F1@ $k$ , Median-Rank@ $k$ ,  
251 Mean-Rank@ $k$ , and Mean-Reciprocal-Rank-(MRR)@ $k$ , with  $k = 10$ . We compared two setups of concept  
252 banks: the report-extracted and the expert-verified. We found that report-extracted concepts achieved Pre-  
253 cision@10 = 0.53 and Recall@10 = 0.57, similar to expert-verified concepts (Precision@10 = 0.54, Recall@10

254 = 0.55). It is worth noting that expert-verified concepts yielded better alignment with expert annotations,  
255 suggesting that human intervention in the verification process improves the concept bank’s ability to  
256 capture domain knowledge. Our analysis demonstrates that the MMCBM model concepts extracted from  
257 reports closely match the performance of expert-verified annotations across various metrics. This suggests  
258 that report-extracted concepts achieve interpretability comparable to expert-verified concepts, negating the  
259 need for time-intensive expert annotation while effectively capturing the salient clinical features of interest  
260 to ophthalmologists.

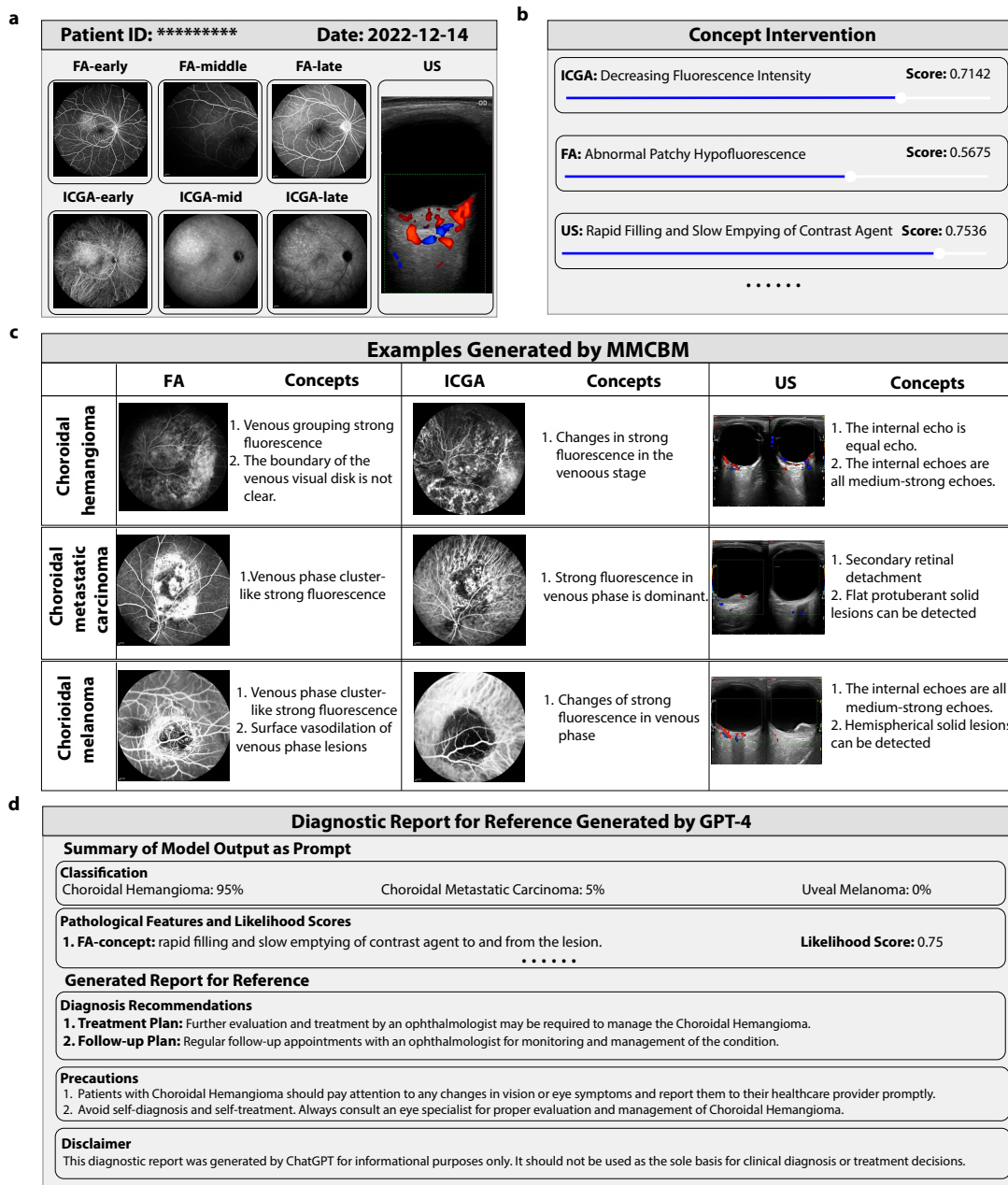
261

## 262 **Demonstration of Human-Model Interaction**

263 To exemplify a practical engineered system for enabling human-model interactions, we make available our  
264 website (<https://mmcbm.liuy.site>) used for this user-based study with the eight ophthalmologists. Our  
265 website provides a user-friendly online interface for concept bank verification and predication evaluation.  
266 The annotation system allows ophthalmologists to upload de-identified images and annotate them with clin-  
267 ically meaningful concepts (**Fig. 5a**) or verify images along with report extracted concepts. The prediction  
268 system can accept FA, ICGA, and/or US images, and use them to output imaging concepts with confidence  
269 scores (**Fig. 5b**). In instances where MMCBM may produce erroneous concept predictions, clinicians can  
270 easily adjust the confidence scores of individual concepts within the user interface. Such adjustments can  
271 refine and correct model predictions to better align them with clinical findings that may be otherwise  
272 inaccessible to the model. This feature of human intervention significantly improves the practical utility of  
273 MMCBMs in clinical decision-making, fostering a more collaborative and accurate diagnostic process. The  
274 **Fig. 5c** displays several examples generated by MMCBM, including a curated selection of representative  
275 cases processed by the model. Finally, given the model outputs, a basic diagnostic report can be generated  
276 by leveraging LLMs to interpret the MMCBM outputs and concept activations (**Fig. 5d**). The generative  
277 model highlights the top- $k$  activated concepts before presenting the final generated diagnostic report. The  
278 report generation prompt example is included in **Supplementary Fig. A6**.

279

280 In summary, our results highlight the MMCBM model as a promising tool for clinical decision support.  
281 While the model’s predictions are accurate on their own, they are most effective when combined with  
282 human expertise, offering the most comprehensive diagnostic performance and underscoring the potential of  
283 AI-assisted diagnostics.



**Fig. 5: Example Human Interactive Interface.** We offer a website to facilitate the user interactive study with ophthalmologists and our trained MMCBM model. **(a)** Image Display Panel: as FA and ICGA imaging span various time frames, ophthalmologists pinpoint images from early, medium, and late phases for accurate classification. **(b)** Interventions interface on concept bottleneck: a panel that allows adjustment of the concept scores to refine the final prediction. **(c)** Visual Emphasis on Bottlenecks: a curated selection of representative cases processed by the model, highlighting the top-k concepts prioritized by their attention scores in the weight matrix displayed across three distinct tumor classes. **(d)** Diagnostic Reporting in Action: an example of a diagnostic report formulated by ChatGPT during the testing phase. The input to ChatGPT includes the predicted top-k concepts combined with patient-specific details, highlighting the model’s capability to produce interpretable diagnoses.

## 284 Discussion

285 In this work, we establish the Multimodal Medical Concept Bottleneck Model (MMCBM) as a novel approach  
286 for the interpretable diagnosis of rare choroid tumors. To facilitate the application of advanced machine  
287 learning techniques, we initially tackled the significant challenge of scarce comprehensive training data by  
288 curating the Choroidal Tri-Modal Imaging clinical dataset. This dataset, which includes image data of Flu-  
289 orescein Angiography (FA), Indocyanine Green Angiography (ICGA), and Ultrasound (US) with associated  
290 radiology reports, to our knowledge, is the largest dataset of choroidal melanoma. Based on this dataset,  
291 our MMCBM maintains the accuracy of prior “black-box” models and introduces interpretability through  
292 the concept bottleneck model. Furthermore, by incorporating the explainable MMCBM into the diagnostic  
293 workflow, our model significantly enhances the performance of junior ophthalmologists.

294 Unlike traditional methods for explainable AI, which often rely on saliency maps [36–38] to highlight  
295 important spatial attributes, our approach is inherently finding-based and therefore aligns more closely with  
296 clinical practice by mimicking the diagnostic thought process used by domain experts. Clinicians identify a  
297 range of descriptive visual features, including textual elements, contrast, shape, and dynamic changes, that  
298 extend beyond pixel values alone. Traditional approaches to incorporating this additional information and  
299 prior knowledge typically require either expensive labeling or sophisticated network infrastructure designs to  
300 integrate clinical insights, thereby limiting the generalized utility of explainable AI tools. Our introduction of  
301 “concepts” addresses this gap by providing human-comprehensible descriptions that facilitate intervention in  
302 the diagnostic process. This yields a twofold benefit: it simplifies the alignment between domain knowledge  
303 in clinical practice and the representational power of neural networks, and it proves immensely beneficial for  
304 junior doctors, who may lack experience in finding identification and risk over-reliance on AI outputs [39].

305 Moreover, recent advancements in vision and natural language processing, such as Large Language Mod-  
306 els (LLMs) and Contrastive Language–Image Pre-training (CLIP), have paved new pathways for research  
307 into interpretable diagnostic systems. However, for rare diseases like choroidal melanoma, the scarcity of  
308 paired image-text knowledge on the internet presents a significant challenge to the reliability of these mod-  
309 els’ reasoning capabilities, as evidenced in Fig. 3. While professional annotation of high-quality data can  
310 mitigate this issue, further data access and expertise challenges remain [33], especially for rare diseases. Our  
311 concept-based multimodal model circumvents these challenges by utilizing LLMs to process texts without  
312 necessitating detailed labeling of image features. The model’s predictive and interpretive power stems from  
313 integrating the pre-trained model with the extracted relationship between reports and images. This approach  
314 mitigates the data scarcity issue for rare diseases in recent foundation models, avoiding the need for extensive  
315 labeling efforts in medical AI preparation, thus making the design extendable to other rare diseases.

316 In the realism of AI-aid medical diagnosis, particularly for the detection and intervention of serious  
317 diseases like the choroid neoplasias we considered in the current work, ethical considerations are of critical  
318 importance [40]. Our methodology, which enables human-in-the-loop feedback, helps address this issue by  
319 aligning human expertise with AI diagnosis. Specifically, by actively involving domain experts in the training  
320 and validation phases of AI model development, we not only ensure that the AI’s diagnostic concepts are  
321 vetted by experienced clinicians but also provide feasible constraints of the degree of AI intervention. This  
322 reduces the risk of hallucinations that could arise from sole reliance on AI. This approach may foster trust  
323 among clinicians and patients in AI-assisted medical decisions. The inclusion of HITL integration in our AI  
324 models aligns with ethical guidelines for AI in healthcare, emphasizing the safeguarding of patient dignity  
325 and privacy. As we advance the frontiers of medical AI, it is crucial to maintain a balanced synergy between  
326 technological innovation and ethical responsibility, ensuring that AI serves as a supportive tool rather than  
327 a replacement for the nuanced judgment of medical professionals.

328 While our results presented in the study are promising, it nonetheless has limitations. Firstly, the multi-  
329 modal data can be noisy due to inconsistencies in image acquisition and labeling. Our filtering pipeline  
330 improves data quality, as confirmed by human expert evaluation, but it still requires careful oversight to  
331 ensure that the MMCBM concepts are well aligned with prior knowledge. Secondly, as with any application  
332 of machine learning in healthcare, the clinical implementation of such models requires rigorous valida-  
333 tion through prospective studies and randomized clinical trials. Collaboration with regulatory bodies will  
334 ensure these diagnostic tools meet safety, efficacy, and equity standards. We hope to explore the capacity of  
335 MMCBMs and other interpretable models to meet these standards in future work.

336 In summary, the development of MMCBMs marks a significant advancement toward achieving inter-  
337 pretable and reliable diagnoses within the healthcare domain. As efforts to refine and incorporate these  
338 models into clinical workflows progress, it is imperative to carefully consider the ethical and regulatory  
339 dimensions to ensure that these innovations enhance patient outcomes without compromising the standards  
340 of care or jeopardizing patient safety. This work delineates a promising avenue for applying artificial intelli-  
341 gence in the nuanced and critical field of diagnosing rare diseases, offering a blueprint for future explorations  
342 in this vital area of medical research.

## 343 **Materials and Methods**

344 **Dataset Collection and Ethics Statement.** The patient data in the CTI dataset were collected at Beijing  
345 Tongren Hospital from January 2004 to December 2022 (Approval No. TRECKY2018-056-GZ(2022)-07). To  
346 our knowledge, it is the largest clinical database containing multimodal data from patients with choroidal

347 melanoma and other closely related ocular pathologies. This extensive database contains diagnostic and  
348 pathological data of patients with choroidal diseases. The database includes a total of 925 cases, which  
349 comprise 161 cases of choroidal hemangioma, 82 cases of choroidal metastatic carcinoma, and 682 cases of  
350 choroidal melanoma. The image collection includes three types of radiological images: fluorescein angiogra-  
351 phy (FA), indocyanine green angiography (ICGA), and Doppler ultrasound images (US). Each patient has  
352 one or more modalities of images. The FA and ICGA images, being time-series, were captured from three  
353 angles: 30, 55, and 102 degrees. The US images include two types: B-mode ultrasound and color Doppler  
354 ultrasound. Medical professionals have thoroughly reviewed the data-cleaning process to ensure its integrity  
355 and clinical relevance. For FA and ICGA modalities, we ignored the shooting angle and categorized the FA  
356 and ICGA images into three periods—early, middle, and late—in alignment with existing clinical diagnostic  
357 recommendations. The time frames for these periods are as follows: ICGA (Early: less than 5 minutes; Mid-  
358 dle: between 5 and 20 minutes; and Late: at least 20 minutes) and FA (Early: less than 5 minutes, Middle:  
359 between 5 and 10 minutes, Late: at least 10 minutes). We selected binocular color Doppler images contain-  
360 ing blood flow information for the US modality. Finally, the cleaned dataset includes a total of 750 cases,  
361 which comprises 128 cases of choroidal hemangioma, 80 cases of choroidal metastatic carcinoma, and 542  
362 cases of choroidal melanoma. There are 53 patients with choroidal hemangioma, 38 patients with choroidal  
363 metastatic carcinoma, and 194 patients with choroidal melanoma with all three imaging modalities, which  
364 we refer to as multi-modal (MM) data. Additionally, 97 cases have clinical diagnostic reports that describe  
365 the radiological features observed in the FA, ICGA, and US images. Informed consent was obtained from all  
366 patients whose anonymized and de-identified data is included in the dataset. Per the Declaration of Helsinki  
367 2000, the collecting organization obtained written informed consent from the patients.

368

369 **Data Splitting.** To optimize data utilization and establish reliable evaluation indicators, we initially allo-  
370 cated 20% of patients with all three imaging studies as the test set and performed 5-fold cross-validation at  
371 the patient level on the remaining data. Specifically, the remaining data is split into five folds based on each  
372 pathology and modality. Data augmentation techniques were applied during training, including random hor-  
373 izontal flipping, random rotating, and random zooming. To build the multi-modal concept banks, we used  
374 97 diagnosis reports, comprising 39 cases of choroidal hemangioma, 18 of choroidal metastatic carcinoma,  
375 and 40 of choroidal melanoma. Each report included three modal images and prompted GPT-4 to extract  
376 relevant medical concepts from reports. The prompts are detailed in **Supplementary Fig. A5**, and the  
377 extracted concepts are in **Supplementary Table A2**.

378



379 **Model Training.** Consider a training dataset  $\mathcal{D}_{train} = \{(\mathbf{x}, \mathbf{r}, y)\}$  comprising image-report pairs, where  
 380  $\mathbf{x} \in \mathcal{X}$  represents a fundus image (of any imaging modality),  $\mathbf{r} \in \mathcal{R}$  is the clinical patient report collected  
 381 by doctors,  $y \in \mathcal{Y} := \{\text{hemangioma, carcinoma, melanoma}\}$  is the corresponding disease label. We utilize  
 382 GPT-4 to analyze the reports and extract relevant concepts, represented as a function  $\text{LLM} : \mathcal{R} \rightarrow \mathcal{C}$   
 383 where  $\mathcal{C}$  is the space of concepts. We can then prompt GPT-4 to combine concepts with the same semantic  
 384 meaning, resulting in a compressed representation of  $N$  concepts  $\mathcal{C} = \{c_1, c_2, \dots, c_N\}$ . Using a pre-trained  
 385 multi-modality backbone  $\phi : \mathcal{X} \rightarrow \mathcal{Z}$  capable of mapping different modality images into a shared feature  
 386 space, we can generate bottleneck embeddings to establish a concept bank, denoted as  $\mathcal{Z}_{\mathcal{C}} \in \mathbb{R}^{N \times d}$ , where  
 387  $N$  is the number of concepts and  $d$  the size of the embedding space of  $\phi$ . Row  $i$  of the two-dimensional  
 388 matrix  $\mathcal{Z}_{\mathcal{C}}$  represents the learned representation of the  $i$ th concept  $c_i$  obtained through Concept Activation  
 389 Vectors (CAVs) [29]. MMCBM generates a prediction  $\hat{y} = g(\text{sim}(\phi(x), \mathcal{Z}_{\mathcal{C}}))$ . The function  $\text{sim} : \mathbb{R}^d \rightarrow \mathbb{R}^N$   
 390 computes the concept scores by calculating the similarities between image features and each element of the  
 391 concept bank  $\mathcal{Z}_{\mathcal{C}}$ . The function  $g : \mathbb{R}^N \rightarrow \mathcal{Y}$  predicts the final label based on the concept scores, serving as  
 392 an interpretable predictor. To learn the MMCBM, we solve the following problem:

$$\min_g \mathbb{E}_{(\mathbf{x}, c, y) \sim \mathcal{D}} \mathcal{L}[g(\text{sim}(\phi(x), \mathcal{Z}_{\mathcal{C}})), y] \quad (1)$$

393 where  $\phi(x)$  is the projection to the concept space and  $\mathcal{L}$  is the cross-entropy loss. To ensure that final  
 394 prediction  $\hat{y}$  can be easily derived from input  $\text{sim}(\phi(x), \mathcal{Z}_{\mathcal{C}})$ , we model  $g$  as a linear classifier.

395

396 **Evaluation of Model Performance.** Using a 5-fold cross-validation framework, we report the macro-  
 397 averaged metrics accuracy, precision, recall, and F1 score, which considers both precision and recall while  
 398 addressing potential class imbalances. In addition to these traditional classification metrics, we also focused  
 399 on interpretability metrics such as Precision@ $k$ , Recall@ $k$ , Mean Rank@ $k$ , and Median Rank@ $k$ . Precision@ $k$   
 400 measures how many of the top- $k$  identified concepts were right compared with the annotated ground truth.  
 401 Recall@ $k$  evaluates the ratio of correct concepts in the first  $k$  predictions to all correct concepts for the  
 402 patient.  $F_1@k$  is the harmonic mean of Precision@ $k$  and Recall@ $k$ . Mean Rank@ $k$  and Median Rank@ $k$   
 403 indicate the average ranking position of the correct concept; lower scores are better.

## 404 Data and Code Availability

405 Due to legal restrictions on the patients' privacy information, the complete raw dataset cannot be made  
406 publicly available. However, we will release the processed data, codes, and trained model to enable repro-  
407 ducibility of the results. The code is currently available at <https://github.com/ly1998117/MMCBM> and will  
408 be moved to a permanent address upon acceptance.

## 409 Acknowledgements

410 S.G. is supported by NSFC Key Program 62236009, Shenzhen Fundamental Research Program (General  
411 Program) JCYJ 20210324140807019, NSFC General Program 61876032, and Key Laboratory of Data Intel-  
412 ligence and Cognitive Computing, Longhua District, Shenzhen. X.Y. and WB.W. are supported by NSFC  
413 82220108017, 82141128, 82002883, the Capital Health Research and Development of Special 2024-1-2052,  
414 Science & Technology Project of Beijing Municipal Science & Technology Commission Z201100005520045,  
415 Sanming Project of Medicine in Shenzhen No. SZSM202311018. M.S.Y. is supported by NIH T32 EB009384.

## 416 References

- 417 [1] Chan, H.-P., Hadjiiski, L. M. & Samala, R. K. Computer-aided diagnosis in the era of deep learning.  
418 *Medical physics* **47**, e218–e227 (2020).
- 419 [2] Johnson, A. E. *et al.* MIMIC-CXR, a de-identified publicly available database of chest radiographs with  
420 free-text reports. *Scientific data* **6**, 317 (2019).
- 421 [3] Lin, M. *et al.* Improving model fairness in image-based computer-aided diagnosis. *Nature Communica-*  
422 *tions* **14**, 6261 (2023).
- 423 [4] Gao, M. *et al.* Discriminative ensemble meta-learning with co-regularization for rare fundus diseases  
424 diagnosis. *Medical Image Analysis* **89**, 102884 (2023).
- 425 [5] Menze, B. H. *et al.* The multimodal brain tumor image segmentation benchmark (brats). *IEEE*  
426 *transactions on medical imaging* **34**, 1993–2024 (2014).
- 427 [6] Liew, S.-L. *et al.* A large, open source dataset of stroke anatomical brain images and manual lesion  
428 segmentations. *Scientific data* **5**, 1–11 (2018).

- 429 [7] Decherchi, S., Pedrini, E., Mordenti, M., Cavalli, A. & Sangiorgi, L. Opportunities and challenges for  
430 machine learning in rare diseases. *Front Med (Lausanne)* **8**, 747612 (2021).
- 431 [8] Molnar, M. J. & Molnar, V. Ai-based tools for the diagnosis and treatment of rare neurological disorders.  
432 *Nature Reviews Neurology* **19**, 455–456 (2023).
- 433 [9] Richens, J. G., Lee, C. M. & Johri, S. Improving the accuracy of medical diagnosis with causal machine  
434 learning. *Nature communications* **11**, 3923 (2020).
- 435 [10] Huang, Z., Bianchi, F., Yuksekgonul, M., Montine, T. J. & Zou, J. A visual–language foundation model  
436 for pathology image analysis using medical twitter. *Nature medicine* **29**, 2307–2316 (2023).
- 437 [11] Zhang, X., Wu, C., Zhang, Y., Xie, W. & Wang, Y. Knowledge-enhanced visual-language pre-training  
438 on chest radiology images. *Nature Communications* **14**, 4542 (2023).
- 439 [12] Radford, A. *et al.* Learning transferable visual models from natural language supervision. *International*  
440 *conference on machine learning (PMLR)* 8748–8763 (2021).
- 441 [13] Singhal, K. *et al.* Large language models encode clinical knowledge. *Nature* **620**, 172–180 (2023).
- 442 [14] Chae, A. *et al.* Strategies for implementing machine learning algorithms in the clinical practice of  
443 radiology. *Radiology* **310**, e223170 (2024).
- 444 [15] Rajpurkar, P., Chen, E., Banerjee, O. & Topol, E. J. Ai in health and medicine. *Nature medicine* **28**,  
445 31–38 (2022).
- 446 [16] Varoquaux, G. & Cheplygina, V. Machine learning for medical imaging: methodological failures and  
447 recommendations for the future. *NPJ digital medicine* **5**, 48 (2022).
- 448 [17] Jager, M. J. *et al.* Uveal melanoma. *Nature reviews Disease primers* **6**, 24 (2020).
- 449 [18] Shields, C. L. *et al.* Metastatic tumours to the eye. review of metastasis to the iris, ciliary body, choroid,  
450 retina, optic disc, vitreous, and/or lens capsule. *Eye* **37**, 809–814 (2023).
- 451 [19] Singh, A. D., Turell, M. E. & Topham, A. K. Uveal melanoma: trends in incidence, treatment, and  
452 survival. *Ophthalmology* **118**, 1881–1885 (2011).
- 453 [20] Kaliki, S., Shields, C. L. & Shields, J. A. Uveal melanoma: estimating prognosis. *Indian journal of*  
454 *ophthalmology* **63**, 93 (2015).

- 455 [21] Egan, K. M., Seddon, J. M., Glynn, R. J., Gragoudas, E. S. & Albert, D. M. Epidemiologic aspects of  
456 uveal melanoma. *Survey of ophthalmology* **32**, 239–251 (1988).
- 457 [22] Augsburger, J. J. & Gamel, J. W. Clinical prognostic factors in patients with posterior uveal malignant  
458 melanoma. *Cancer* **66**, 1596–1600 (1990).
- 459 [23] Carvajal, R. D. *et al.* Metastatic disease from uveal melanoma: treatment options and future prospects.  
460 *British Journal of Ophthalmology* (2016).
- 461 [24] Khoja, L. *et al.* Meta-analysis in metastatic uveal melanoma to determine progression free and overall  
462 survival benchmarks: an international rare cancers initiative (irci) ocular melanoma study. *Annals of*  
463 *Oncology* **30**, 1370–1380 (2019).
- 464 [25] Mathis, T. *et al.* New concepts in the diagnosis and management of choroidal metastases. *Progress in*  
465 *retinal and eye research* **68**, 144–176 (2019).
- 466 [26] Tan, M. & Le, Q. Efficientnet: Rethinking model scaling for convolutional neural networks. *International*  
467 *conference on machine learning (PMLR)* 6105–6114 (2019).
- 468 [27] Safari, P., India, M. & Hernando, J. Self-attention encoding and pooling for speaker recognition. *arXiv*  
469 *preprint arXiv:2008.01077* (2020).
- 470 [28] Achiam, J. *et al.* Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- 471 [29] Kim, B. *et al.* Interpretability beyond feature attribution: Quantitative testing with concept activation  
472 vectors (tcav). *International conference on machine learning (PMLR)* 2668–2677 (2018).
- 473 [30] Koh, P. W. *et al.* Concept bottleneck models. *International conference on machine learning (PMLR)*  
474 5338–5348 (2020).
- 475 [31] Yuksekgonul, M., Wang, M. & Zou, J. Post-hoc concept bottleneck models. *International conference*  
476 *on learning representations (ICLR)* (2023).
- 477 [32] Yang, Y. *et al.* Language in a bottle: Language model guided concept bottlenecks for interpretable image  
478 classification. *In Proc. IEEE Conference on Computer Vision and Pattern Recognition* 19187–19197  
479 (2023).

- 480 [33] Wang, Z., Wu, Z., Agarwal, D. & Sun, J. Medclip: Contrastive learning from unpaired medical images  
481 and text. *arXiv preprint arXiv:2210.10163* (2022).
- 482 [34] Zhang, S. *et al.* Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv*  
483 *preprint arXiv:2303.00915* (2023).
- 484 [35] Van der Maaten, L. & Hinton, G. Visualizing data using t-sne. *Journal of machine learning research* **9**  
485 (2008).
- 486 [36] Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-based localization.  
487 *In Proc. of the IEEE international conference on computer vision* 618–626 (2017).
- 488 [37] Wu, Y. *et al.* Interpretable identification of interstitial lung disease (ild) associated findings from ct.  
489 *Medical Image Computing and Computer Assisted Intervention* 560–569 (2020).
- 490 [38] Tjoa, E. & Guan, C. A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE*  
491 *transactions on neural networks and learning systems* **32**, 4793–4813 (2020).
- 492 [39] Kostick-Quenet, K. M. & Gerke, S. Ai in the hands of imperfect users. *npj Digital Medicine* **5**, 197  
493 (2022).
- 494 [40] Char, D. S., Abràmoff, M. D. & Feudtner, C. Identifying ethical considerations for machine learning  
495 healthcare applications. *The American Journal of Bioethics* **20**, 7–17 (2020).

## 496 Appendix A

### 497 A.1 Supplementary Results

#### 498 A.1.1 Varying the Image Encoder Size

499 To further explore the impact of encoders’ size on the efficacy of pretrained models and MMCBMs, we under-  
500 took a study incorporating three variants of EfficientNet encoders: `efficientnet-b0`, `efficientnet-b1`, and  
501 `efficientnet-b2`—each integrated separately within the respective modalities. **Figure A2** shows the rela-  
502 tionship between encoder size and pivotal performance metrics, depicted across three detailed comparative  
503 bar plots which focus on two key aspects: classification performance and model interpretability.

504 **Figure A2[a-b]** show the accuracy, precision, recall, and  $F_1$  scores of the resulting models. Our results  
505 demonstrate high performance across the board with negligible differences between the encoder sizes, as  
506 evidenced by the closely clustered bars and their associated error margins.

507 Regarding the interpretability of MMCBMs, **Figure A2c** shows a consistent performance trend across  
508 the different encoders, without any marked disparity in the results. This section of the figure examines the  
509  $\text{precision}@k$ ,  $\text{recall}@k$ , and  $F_1@k$  metrics, alongside mean  $\text{rank}@k$ , median  $\text{rank}@k$ , and mean reciprocal  
510  $\text{rank}@k$  for predicted Top- $k$  concepts.

511 Drawing from these insights, we selected `efficientnet-b0` as the encoder for our MMCBMs. This deci-  
512 sion is underpinned by the encoder’s ability to deliver high classification accuracy and robust interpretability,  
513 while also ensuring a more compact model with fewer parameters. `efficientnet-b0` strikes a balance  
514 between maintaining performance standards and optimizing computational efficiency, which is particularly  
515 advantageous in scenarios where resource constraints are a key consideration.

#### 516 A.1.2 Ablation Study for the Concept Bank

517 To further explore the efficacy of the Multi-modal Medical Concept Bottleneck Model (MMCBM), an ablation  
518 study was conducted, focusing on two critical aspects of the model: the concept extraction from diagnostic  
519 reports and the verification of these concepts. Our primary aim was to compare the impact of varying the  
520 number of reports and concepts on the performance of two distinct concept banks: *the Report-Extracted*  
521 *Concept Bank* and *the Expert-verified Concept Bank*. Such a comparison is intended to shed light on the  
522 relative influence in the overall functioning of the MMCBM.

523 **Report-Extracted Concept Bank:** In our study, we first prompt the LLM to extract concepts from a set  
524 of 97 diagnostic reports. Then those image-concept pairs are used directly to generate a multi-modal concept  
525 bank with CAVs.

526 **Expert-verified Concept Bank:** Subsequently, the concepts derived from the LLM were meticulously  
527 reviewed by two experienced ophthalmologists. Their role was to verify the accuracy of the concept extraction,  
528 remap the relationships between images and concepts, and rectify any inaccuracies found in the initial concept  
529 set. This refined collection of concepts, validated and enhanced by medical expertise, was termed the Expert-  
530 verified Concept Bank. To streamline this verification and correction process, we developed a specialized  
531 web interface (**Fig. 5**). This interface empowered the ophthalmologists to efficiently identify and correct  
532 erroneous concepts, delete irrelevant or incorrect entries, and introduce additional concepts as necessary.

533 **The impact of varying the number of concepts.** The first part of the ablation study assessed the perfor-  
534 mance impact of varying the number of concepts extracted from these reports. For both the Report-Extracted  
535 Concept Bank and Expert-verified Concept Bank, a noticeable trend in the FA and ICGA modalities  
536 showed performance improvements as the number of concepts increased, plateauing at higher concept counts.  
537 Specifically, for concept counts below 60, the US modality demonstrated marked improvements. Conversely,  
538 exceeding 60 concepts led to a notable decline in performance, suggesting an influx of potentially irrelevant  
539 or ‘bad concepts’, as indicated on our graphs. Interestingly, for the Report-Extracted Concept Bank, this  
540 trend was also apparent in the MM modality when the concept counts around 80. However, for Expert-  
541 verified Concept Bank, there is no exhibit significant decrease in performance, underscoring the importance  
542 of concept relevance and quality in the bank.

543 **The impact of varying the number of reports.** In the second part of our ablation study, shown in  
544 **Figure A4**, we varied the number of reports from 7 to 97. Contrary to our expectations, the performance  
545 metrics remained relatively stable across this range, thereby supporting the conclusion that the model’s  
546 effectiveness is not substantially influenced by the quantity of reports. This observed stability in performance  
547 metrics, even with a limited number of reports, may be attributed to the comprehensive nature of the  
548 concepts contained within these reports. It appears that even a smaller set of templated reports encompasses  
549 a sufficient range of concepts. This suggests that the key determinants for model performance lie not in the  
550 sheer quantity of reports, but rather in the richness and relevance of the concepts they encompass. Notably,  
551 despite the general stability in performance across both concept banks, the Expert-verified Concept Bank  
552 consistently exhibited more robust and superior performance. This finding further emphasize the critical role  
553 of concept validity in enhancing model accuracy.

## 554 **A.2 Supplementary Methods**

555 **Interpretable Predictor.** We trained a linear layer, denoted as the interpretable predictor  $g$ , to learn the  
556 preferences (weights) of specific categories for certain concepts. Then, we linearly combined these weights

557 and used a multi-modality concept score  $\mathbf{C}_{score}$  as the only input to obtain the final prediction. Here,  $\mathbf{C}_{score}$   
 558 represents the similarities between image features and each row within the concept bank  $\mathcal{Z}_{\mathcal{C}}$ . To measure the  
 559 attention of the predictor on an image’s concept score, we applied a sigmoid activation function to the weight  
 560 and performed an element-wise multiplication between this activated weight and concept score to produce  
 561 the attention matrix  $W_{atten}$ . Formally, considering  $W$  as a learnable matrix and  $\sigma$  as the sigmoid activation  
 562 function, the prediction  $\hat{\mathcal{Y}}$  was determined by  $\hat{\mathcal{Y}} = \text{argmax}(\sum W_{atten})$ , where  $W_{atten} = \mathbf{C}_{score} \odot \sigma(W)$  and  
 563  $\odot$  represents the element-wise multiplication.

564 **Model Inference** The Concept Bottleneck Model involves the initial transformation or mapping of the  
 565 input data into a representation comprising a set of concepts, which are subsequently utilized for prediction.  
 566 During the reasoning stage, the CBM model’s interpretability is primarily demonstrated by its ability to  
 567 provide clear explanations of how each concept contributes to the predicted results.

568 **Report Generation.** Medical report generation (MRG) is a task that involves automatically generating a  
 569 descriptive narrative report in the medical domain based on a given medical image. By utilizing our proposed  
 570 interpretable predictor and predicted concept activation scores, we can extract multiple concepts associated  
 571 with an image. This can be done by selecting either the top-k concepts or those that exceed a predetermined  
 572 threshold. It is important to note that clinical diagnosis reports follow a structured format that includes  
 573 components such as patient information, medical details, diagnosis, treatment recommendations, and other  
 574 relevant information. Consequently, it becomes feasible to generate standardized clinical diagnosis reports  
 575 by effectively combining the predicted concepts and prompting a Large Language Model (LLM) to produce  
 576 comprehensive reports.

577 **Test-time intervention.** A salient difference between CBMs and standard models is that a practitioner  
 578 utilizing a CBM model can interact with it by intervening on concept predictions. This ability to intervene  
 579 on concept bottleneck models enables human users to have richer interactions with them. This kind of test-  
 580 time intervention can be particularly useful in high-stakes settings like medicine. For example, “correcting”  
 581 the model by replacing the  $j$ -th concept value  $\hat{c}_j$  with the ground truth value  $c_j$ , and then updating the  
 582 prediction  $\hat{\mathcal{Y}}$  after this replacement. We can qualitatively see the contribution of each concept by removing  
 583 the concept and seeing the changes in the corresponding prediction’s output  $\hat{\mathcal{Y}} = g(\hat{\mathbf{C}}_{score})$ , where  $\hat{\mathbf{C}}_{score} =$   
 584  $(\hat{c}_{\{1, \dots, N\} \setminus j}, c_j)$ .

### 585 A.3 Implementation Details

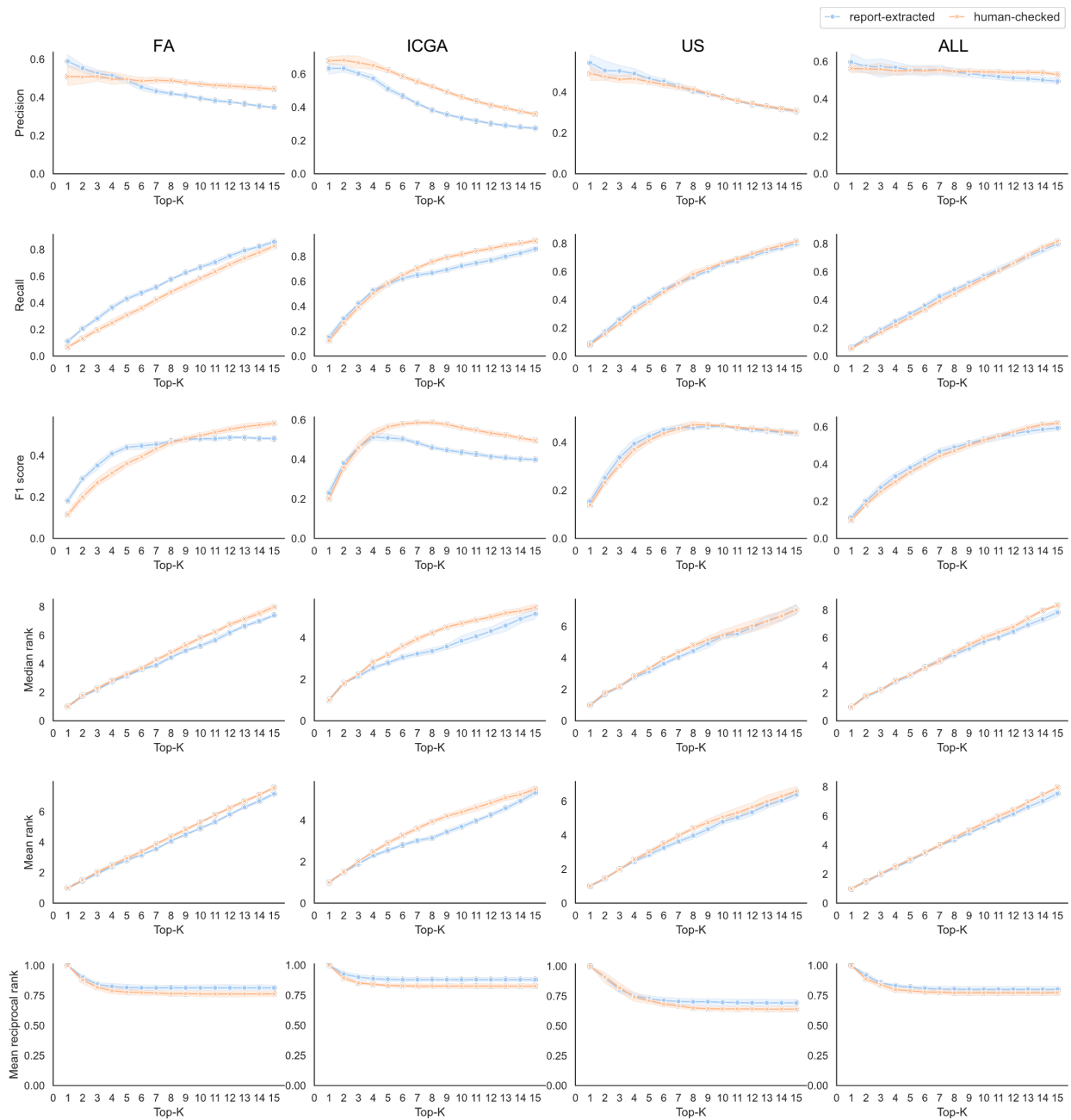
586 In the initial assessment of our dataset, we encountered several quality issues. These included discrepancies  
 587 between diagnostic reports and classification labels, images that were entirely black, blurred images with



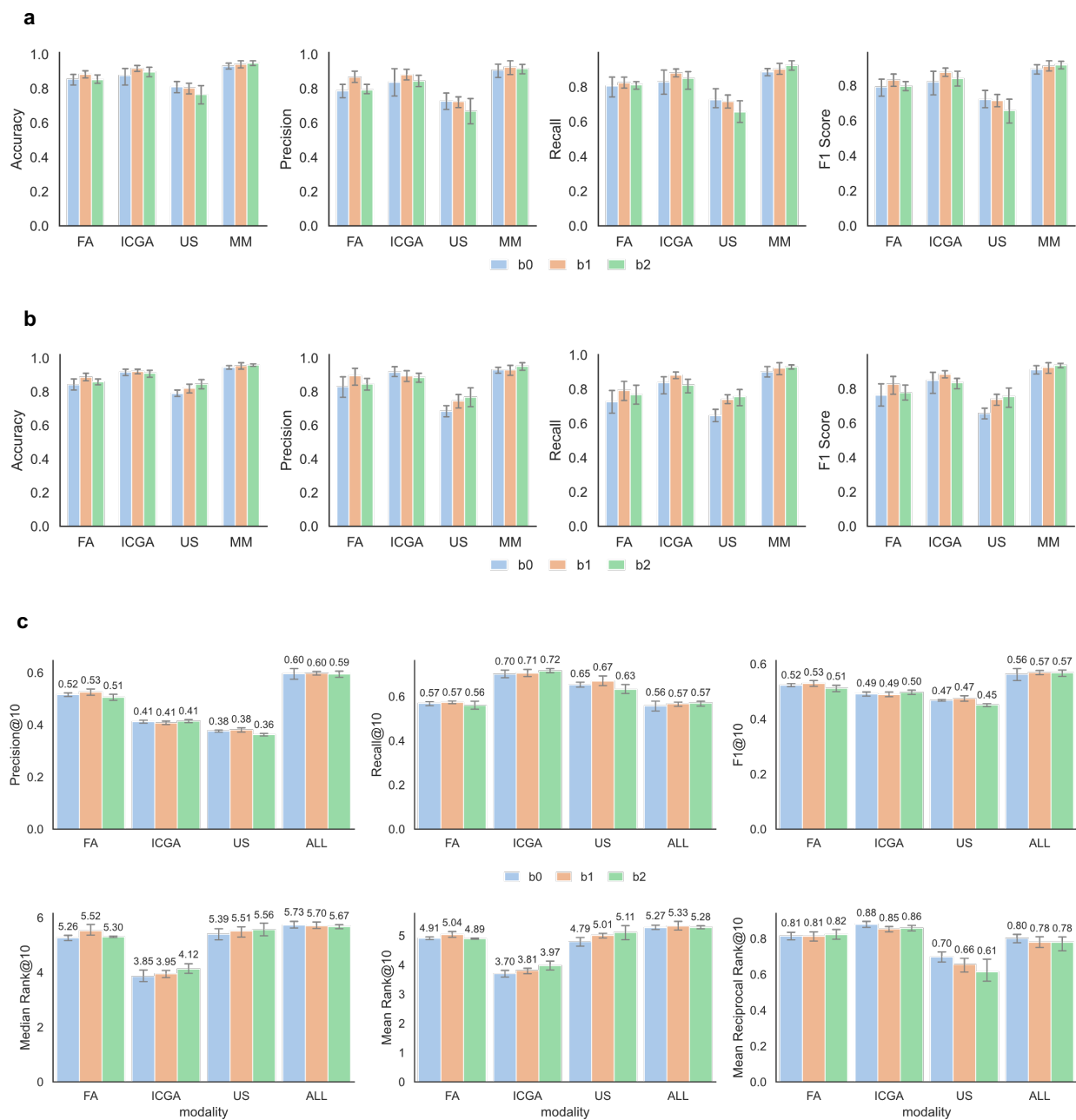
588 noise, and instances of data being saved twice. To address these, we implemented a rigorous data-cleaning  
589 process. Firstly, we removed all modalities containing entirely black images. Secondly, we eliminated images  
590 of poor quality, which included those with excessive blurring and noise. Thirdly, we deleted any duplicate  
591 images to ensure each data point was unique.

592 In our approach, we employ an early stopping mechanism to efficiently mitigate overfitting and save  
593 training time. Additionally, we perform standardization and normalization on the input images. To train  
594 the network, we set the initial learning rate to  $10^{-4}$ , and the weight decay factor to  $10^{-2}$  for the Adam  
595 optimizer. The network is trained for 200 epochs, with a batch size of 8. All experiments were implemented  
596 using the NVIDIA GeForce TiTAN XP GPU and a RAM capacity of 12 GB.

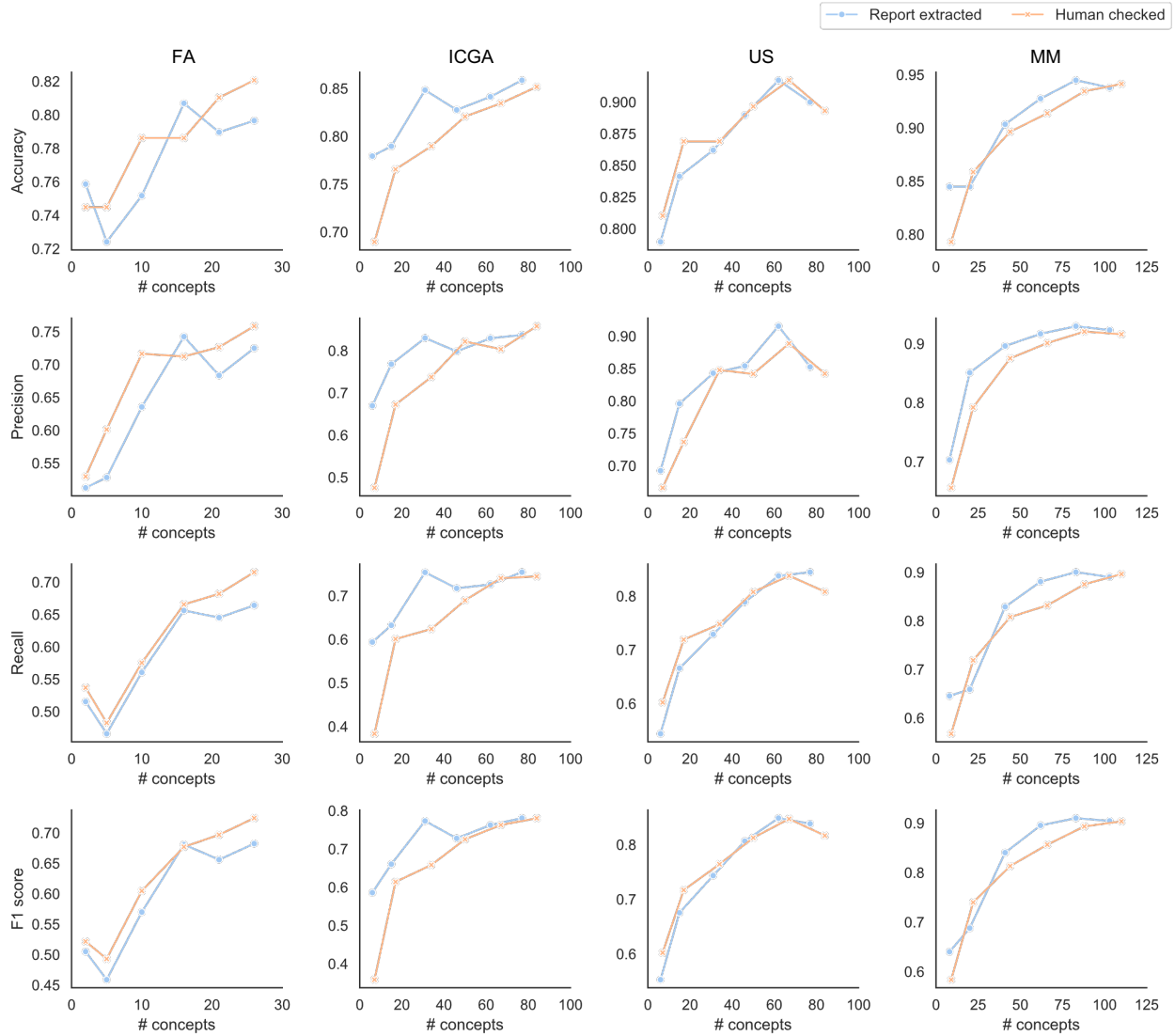
597 The MMCBMs underwent training and evaluation under similar conditions to previous models, with the  
598 notable difference of a learning rate set at  $10^{-3}$ .



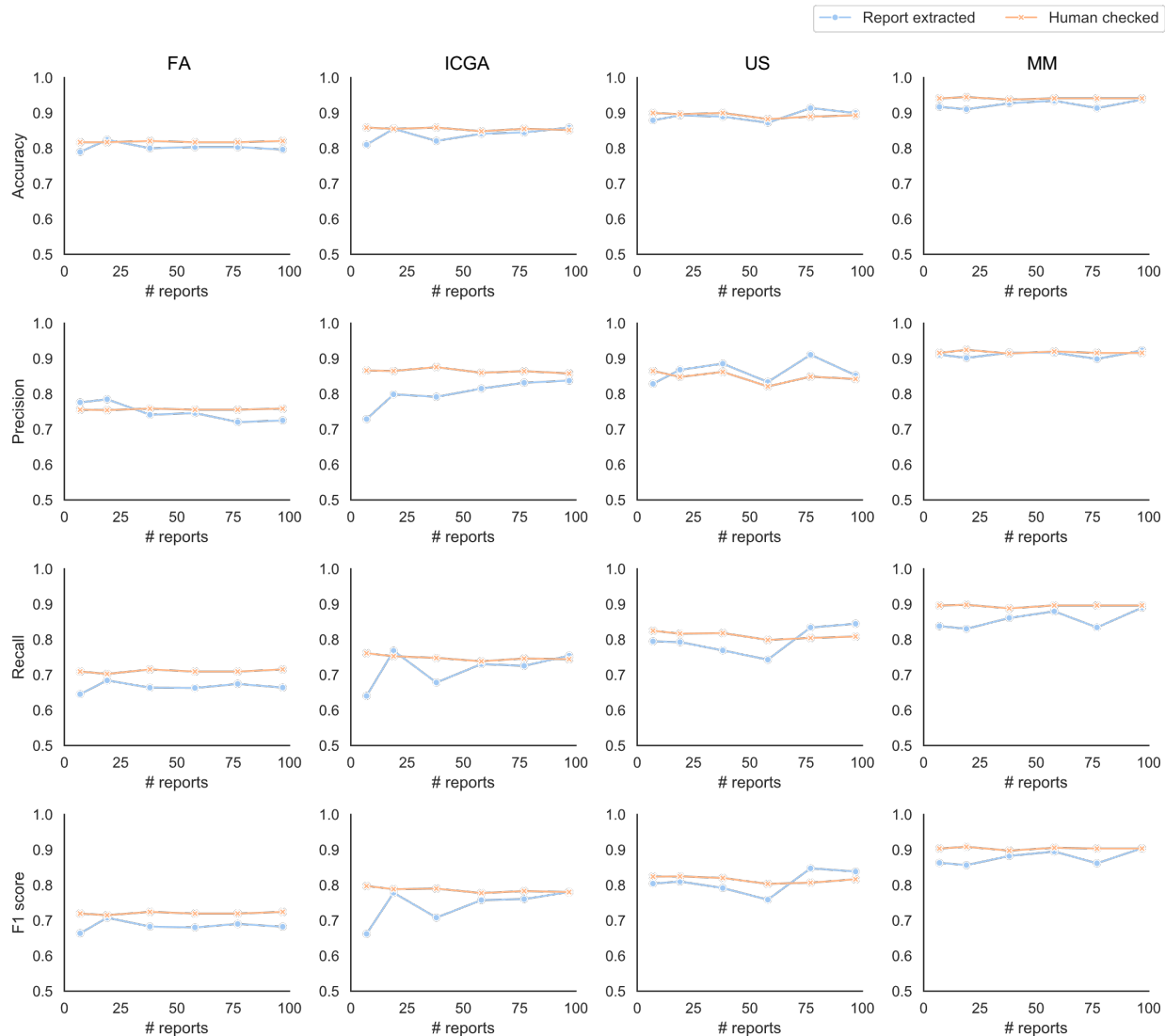
**Fig. A1: Extended details of comparative human evaluation for two concept bank generation methods across multiple modalities.** Each row displays key retrieval metrics: precision@ $k$ , recall@ $k$ , and  $F_1$ @ $k$ , as well as median rank@ $k$ , mean rank@ $k$ , and mean reciprocal rank@ $k$ . Each column corresponds to various data modalities: FA, ICGA, US, and ALL. The ‘ALL’ category represents the aggregation of top- $k$  concepts derived from FA, ICGA and US.



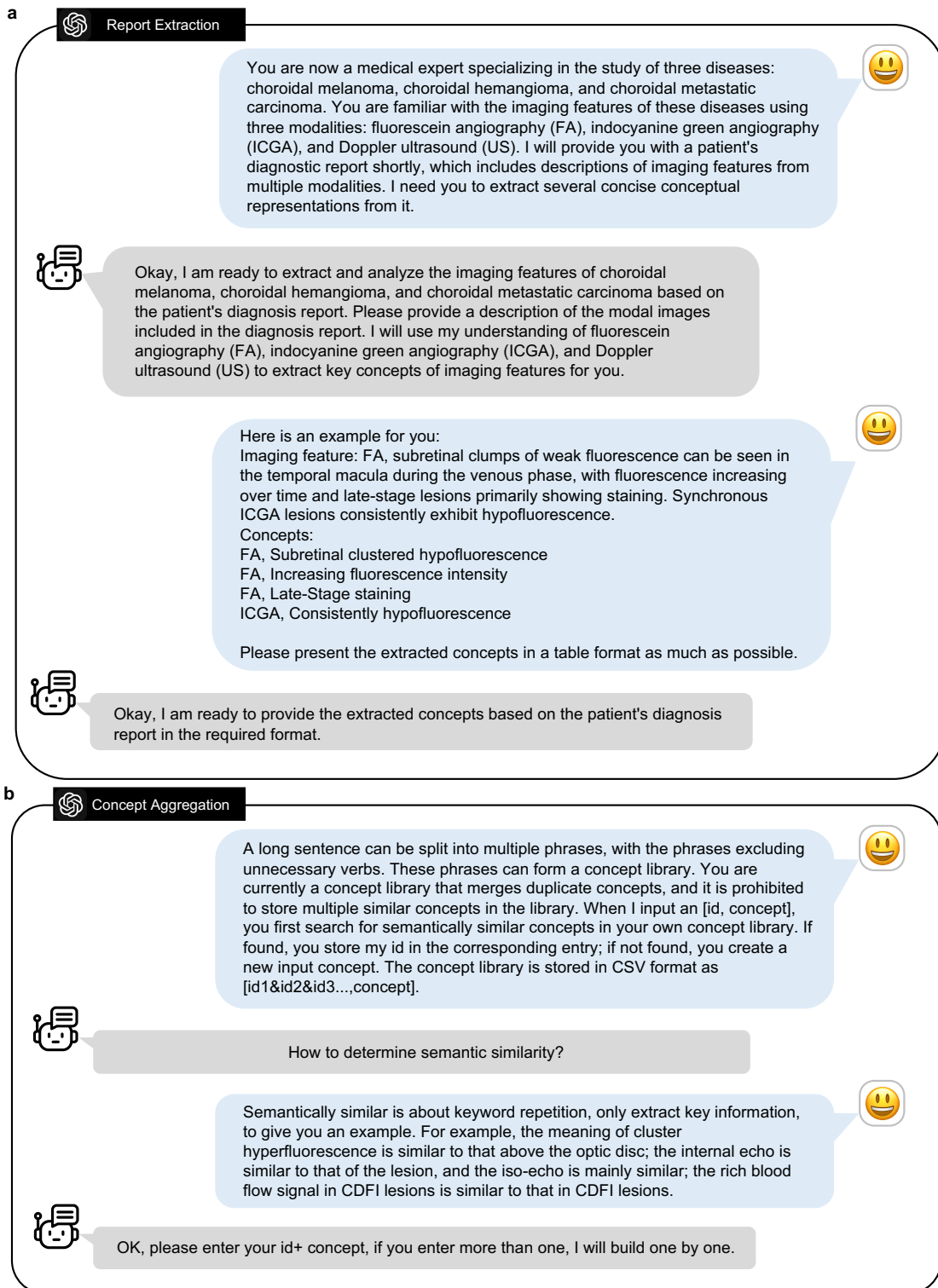
**Fig. A2: Comparing different image encoder sizes.** (a) Comparative bar plot illustrating classification performance metrics of the pre-trained model employing different encoders (b0, b1, b2): accuracy, precision, recall, and  $F_1$  scores. Notably, the best-performing model on the validation dataset was chosen for evaluation. Metrics are represented as mean values, with error bars indicating the 95% confidence interval based on 5-fold cross-validation on the test dataset. (b) A similar comparative bar chart is shown for the performance metrics of MMCBMs, illustrating how varying encoder sizes impact their classification effectiveness. (c) Metrics of predicted Top- $k$  concepts of MMCBMs with different encoders (b0, b1, b2) on test dataset with  $k = 10$ . This evaluation includes precision@ $k$ , recall@ $k$ , and F1@ $k$ , as well as mean rank@ $k$ , median rank@ $k$ , and mean reciprocal rank@ $k$ .



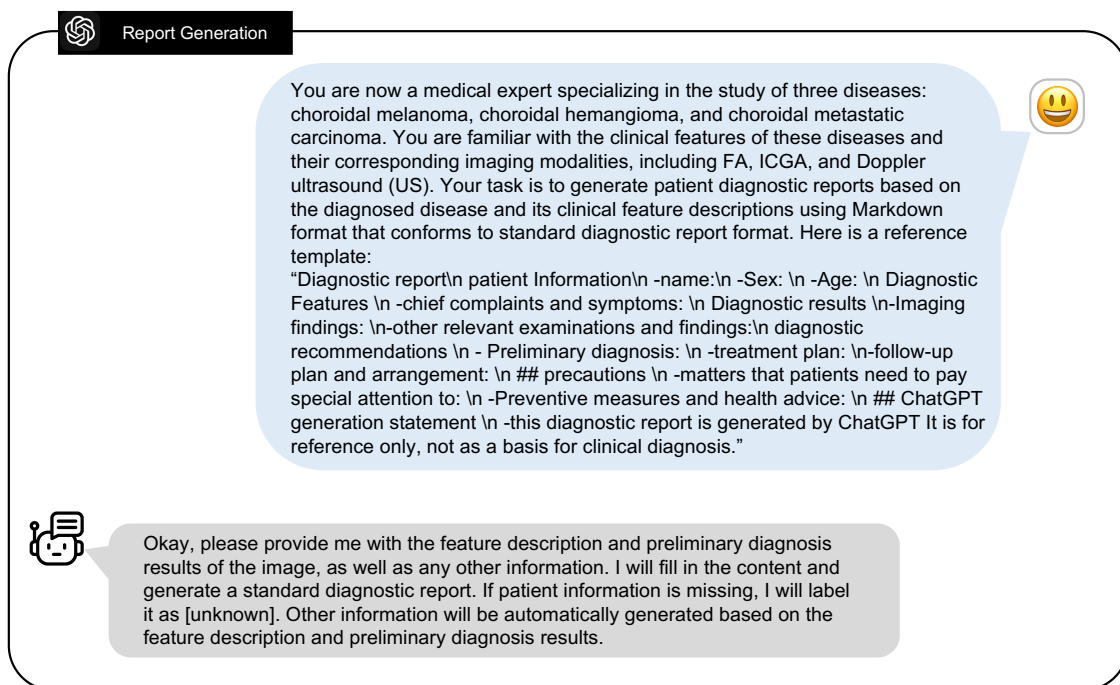
**Fig. A3: Ablation Study Highlighting Dependency on Concept Bank Size.** The concept bank, devised leveraging reports from 100 patients with tri-modality images, forms the basis of this examination. **(a)** Report-extracted concept bank. The initial set of figures (first row) showcases the impact of Report Volume: This set of graphs portrays how adjusting the total number of reports impacts the model’s performance metrics—precision, accuracy, recall, and  $F_1$  score. The subsequent set of figures (second row) explores the effect of Concept Count: By altering the number of disease-associated concepts, we analyze how the model’s performance metrics fluctuate. This study assists in understanding the optimal number of concepts required for reliable diagnosis. **(b)** Expert-verified concept bank. We maintains identical experimental settings to **(a)** for direct comparability.



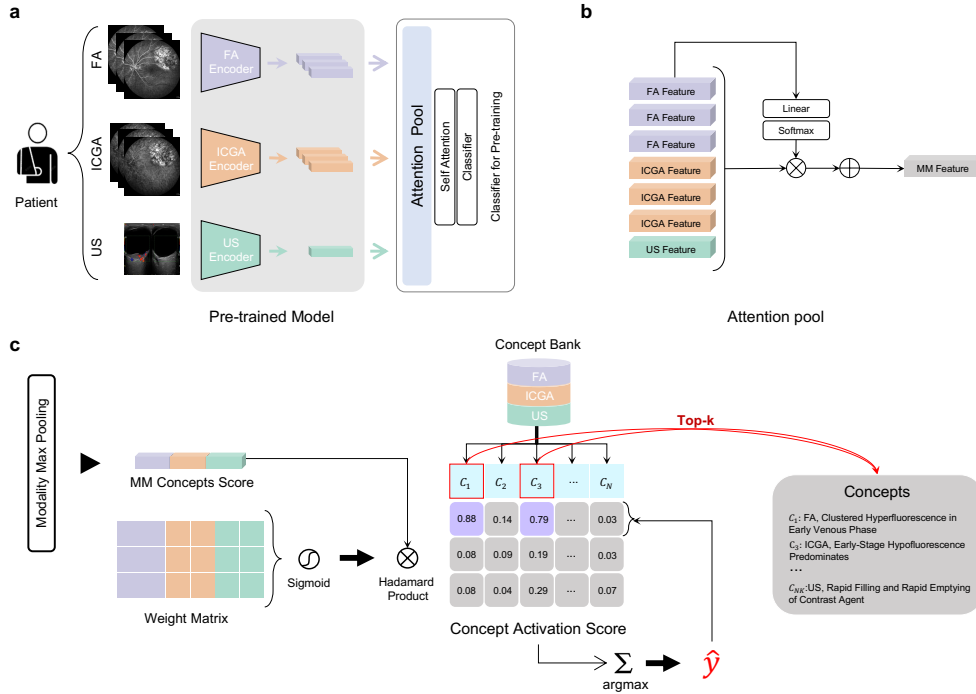
**Fig. A4: Ablation Study on Number of Reports.** The concept bank, devised leveraging reports from 100 patients with tri-modality images, forms the basis of this examination. **(a)** Report-extracted concept bank. The initial set of figures (first row) showcases the impact of Report Volume: This set of graphs portrays how adjusting the total number of reports impacts the model’s performance metrics—precision, accuracy, recall, and  $F_1$  score. The subsequent set of figures (second row) explores the effect of Concept Count: By altering the number of disease-associated concepts, we analyze how the model’s performance metrics fluctuate. This study assists in understanding the optimal number of concepts required for reliable diagnosis. **(b)** Expert-verified concept bank. We maintain identical experimental settings to **(a)** for direct comparability.



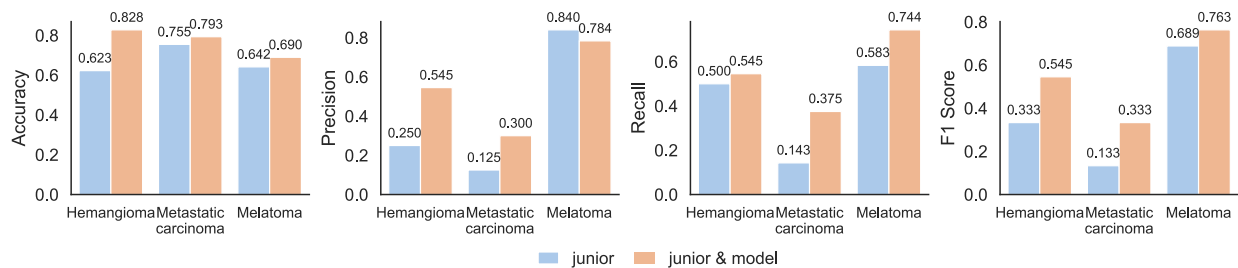
**Fig. A5: Detailed GPT-4 Prompts.** (a) Concept Extraction Prompts. This part details the prompts used for extracting concepts from diagnostic reports, forming a unique concept library. (b) Concepts Aggregation Prompts. This part describes the prompts for aggregating semantically similar concepts, ensuring each concept's uniqueness.



**Fig. A6: Detailed GPT-4 Report Generation Prompts.** The prompt used for generating patient diagnostic reports based on extracted concepts and diagnosed diseases are shown.



**Fig. A7: Details of the MMCBM Framework Module Design.** (a) Illustration of the multi-modal pre-trained prediction model, where images from distinct modalities are transformed into features through dedicated encoders and merged by an attention-pooling block to determine the final classification prediction. (b) Feature Fusion Process via the Attention Pool Module. (c) Linear Layer Mechanism: a trainable weight matrix is optimized to predict tumor classes based on concept scores, with the concept activation score—derived from the Hadamard Product—elucidating the relationship between the current images, representing individual patients, and the respective concepts to support prediction.



**Fig. A8: Details of the Performance Benchmark with Junior Evaluators.**



**Table A1:** Comparative Performance of Pre-trained Classifier and MMCBM on the hold-out test datasets.

Model	Modal	Test Dataset Metrics with 95% Confidence Interval			
		Accuracy	Precision	Recall	F1-score
Pre-trained	FA	84.8% (82.1 - 86.9%)	78.6% (73.8 – 82.6%)	78.9% (76.0 – 81.8%)	78.3% (74.0 – 81.7%)
	ICGA	90.3% (88.6 - 92.4%)	85.5% (82.8 – 88.0%)	88.4% (86.4 – 90.6%)	85.9% (83.7 – 88.2%)
	US	81.4% (77.2 – 84.5%)	73.5% (68.7 – 78.0%)	71.5% (66.4 – 75.6%)	72.1% (67.1 – 76.7%)
	MM	92.8% (91.7 – 93.8%)	89.8% (87.9 – 92.1%)	89.2% (88.0 – 90.9%)	89.2% (87.9 – 90.6%)
MMCBM	FA	84.1% (81.0 - 87.6%)	82.9% (76.4 – 89.3%)	72.6% (66.0 – 79.0%)	76.2% (69.9 – 82.5%)
	ICGA	91.7% (89.7 - 93.4%)	91.5% (89.0 – 95.0%)	83.8% (77.2 – 87.2%)	84.9% (77.3 – 89.4%)
	US	79.0% (76.9 – 81.0%)	68.3% (65.0 – 71.8%)	64.5% (61.1 – 68.2%)	65.6% (62.5 – 68.7%)
	MM	94.5% (93.4 – 95.5%)	93.0% (90.9 – 94.5%)	90.0% (86.9 – 93.4%)	91.0% (88.2 – 93.4%)

**Table A2:** List of report-extracted concepts and the corresponding accuracy of SVMs.

Concept	Train	Test	Concept	Train	Test
FA, Patchy ("Stippled") Hypofluorescence	1.00	0.99	ICGA, Rich Vascular Network within the Lesion	0.99	0.97
FA, Abnormal Vasodilation at the Lesion Surface	0.98	0.98	ICGA, "Double Circulation" Pattern	1.00	1.00
FA, Rich Vascular Network within the Lesion	0.98	0.99	ICGA, Constant Hypofluorescence	1.00	0.94
FA, Pinpoint Hyperfluorescence	0.99	0.94	ICGA, Pinpoint Hyperfluorescence	0.99	1.00
FA, Pinpoint Hyperfluorescence in Arteriovenous Phase	1.00	1.00	ICGA, Large Fluorescence Changes in Venous Phase	1.00	1.00
FA, Increasing Fluorescence Intensity During the Arteriovenous Phase	1.00	0.99	ICGA, Late-Stage Hyperfluorescence Predominates	1.00	0.99
FA, Large Fluorescence Changes in Arterial Phase	1.00	1.00	ICGA, Clustered Venous Phase Hyperfluorescence	1.00	1.00
FA, Hyperfluorescence Primarily During Arterial Phase	1.00	1.00	ICGA, Large Fluorescence Changes	0.99	0.99
FA, Clustered Hyperfluorescence in Arterial Phase	1.00	0.99	ICGA, Predominant Hyperfluorescence	1.00	1.00
FA, Rich Vascular Network in Early Arterial Phase	1.00	0.99	ICGA, Dynamic Hypofluorescence Changes	0.99	1.00
FA, Clustered Hyperfluorescence During Early Arterial Phase	1.00	1.00	ICGA, Predominant Hypofluorescence	1.00	0.95
FA, Abnormal Vasodilation on the Lesion Surface During Venous Phase	0.99	0.99	ICGA, Decreasing Fluorescence Intensity	1.00	1.00
FA, Rich Vascular Network in Venous Phase	1.00	1.00	ICGA, Gradually Increasing Fluorescence Intensity Over Time	1.00	1.00
FA, Abnormal Pinpoint Hyperfluorescence During Venous Phase	0.99	0.98	ICGA, Clustered Hyperfluorescence	1.00	0.99
FA, Hemispherical Solid Lesions During Venous Phase	1.00	1.00	ICGA, Clustered Hypofluorescence	0.99	0.95
FA, Large Fluorescence Changes in Venous Phase	1.00	1.00	ICGA, Late-Stage 'Double Circulation' Pattern in Lesions	1.00	1.00
FA, Hyperfluorescence Primarily During Venous Phase	1.00	1.00	ICGA, Constant Late-Stage Fluorescence	1.00	1.00
FA, Hypofluorescence Primarily During Venous Phase	1.00	1.00	ICGA, Late-Stage Leakage, Mild	1.00	0.99
FA, Blurred Optic Disc Margins in Venous Phase	0.99	0.99	ICGA, Late-Stage Decreasing Fluorescence Intensity	1.00	0.98
FA, Increasing Fluorescence Intensity During the Venous Phase	1.00	1.00	ICGA, Late-Stage "Washout" Phenomenon	1.00	1.00
FA, Clustered Hyperfluorescence During Venous Phase	1.00	1.00	ICGA, Late-Stage Staining	1.00	1.00
FA, Clustered Hypofluorescence During Venous Phase	0.99	0.99	ICGA, Inferior Retinal Elevation	1.00	1.00
FA, Inferior Retinal Elevation During Venous Phase	1.00	1.00	ICGA, Obscured Hypofluorescence	0.99	0.99
FA, Venous Phase Leakage	1.00	1.00	ICGA, Staining	1.00	1.00
FA, Obscured Venous-Phase Hypofluorescence	0.99	0.99	ICGA, Early-Stage Patchy ("Stippled") Hypofluorescence	1.00	1.00
FA, Early Pinpoint Hyperfluorescence in Venous Phase	1.00	1.00	ICGA, Early-Stage Rich Vascular Network	1.00	1.00
FA, Diffuse Hyperfluorescence in Early Venous Phase	1.00	1.00	ICGA, Early-Stage Pinpoint Hyperfluorescence	1.00	0.99
FA, Hyperfluorescence Primarily During Early Venous Phase	1.00	1.00	ICGA, Early-Stage Hyperfluorescence Predominates	1.00	0.97
FA, Clustered Hyperfluorescence in Early Venous Phase	1.00	1.00	ICGA, Early-Stage Hypofluorescence Predominates	1.00	0.96
FA, Clustered Hypofluorescence in Early Venous Phase	1.00	1.00	ICGA, Early-Stage Clustered Hyperfluorescence	1.00	1.00
FA, Blurred Optic Disc Margins	0.99	1.00	US, Abnormal Blood Flow in Lesion Using Doppler	0.87	0.85
FA, Gradually Increasing Fluorescence Intensity Over Time	0.98	0.99	US, Clear and Regular Lesion Borders	0.89	0.84
FA, Clustered Hyperfluorescence	1.00	1.00	US, Irregular Lesion Borders	1.00	0.99
FA, Constant Late-Stage Fluorescence	1.00	1.00	US, Band-shaped Echo Visualized on Lesion Surface	0.93	0.90
FA, Late-Stage Diffuse Hyperfluorescence	1.00	1.00	US, Rapid Filling and Rapid Emptying of Contrast Agent to and from the Lesion	0.93	0.86
FA, Primarily Late-Stage Hyperfluorescence	1.00	1.00	US, Rapid Filling and Slow Emptying of Contrast Agent to and from the Lesion	0.97	0.96
FA, Late-Stage Leakage, Mild	1.00	0.99	US, Vitreous Opacity	0.96	0.96
FA, Late-Stage Hypofluorescence Predominates	1.00	1.00	US, Intravitreal Band Connected to Primary Lesion	0.92	0.82
FA, Late-Stage Decreasing Fluorescence Intensity	1.00	1.00	US, Point- or Band- shaped Hypoechoic Intravitreal Lesion(s) Connected to the Retina	0.96	0.94
FA, Late-Stage Increasing Fluorescence Intensity	1.00	1.00	US, Point- or Band- shaped Hypoechoic Intravitreal Lesion(s) Not Connected to the Retina	0.93	0.84
FA, Late-Stage Leakage	0.99	0.94	US, Solid Lesion with Irregular Borders	0.94	0.94
FA, Late-Stage Staining	1.00	1.00	US, Negative Imaging Findings with Movement on Dynamic Imaging	0.98	0.98
FA, Inferior Retinal Elevation	0.98	0.99	US, Positive Posterior Scleral Concavity	0.99	0.99
FA, Leakage	0.99	0.98	US, Fluttering with Movement on Dynamic Imaging	0.96	0.97
FA, Staining	0.99	0.99	US, Secondary Retinal Detachment	0.90	0.84
FA, Early-Stage Hyperfluorescence Predominates	0.99	0.99	US, Hemispherical Solid Lesions	0.90	0.79
FA, Early-Stage Hypofluorescence Predominates	0.98	0.94	US, Flat and Raised Solid Lesions	0.98	0.99
			US, Abnormal Choroidal Concavity	0.98	0.97
			US, Predominantly Internally Isoechoic Imaging Findings Present	0.93	0.93
			US, Internally Isoechoic and Hypoechoic Imaging Findings Present	0.96	0.97
			US, Internally Isoechoic and Hyperechoic Imaging Findings Present	0.95	0.97
			US, Predominantly Hypoechoic Imaging Findings Present	0.91	0.83
			US, Isoechoic and Hypoechoic Imaging Findings Present	0.96	0.93
			US, Irregular, Solid, Nodular Lesions Near the Optic Disc	0.98	0.99
			US, Ultrasonographic Hollowing	0.90	0.89
			US, No Ultrasonographic Hollowing or Abnormal Findings in the Posterior Fossa Choroid Plexus	0.95	0.89