

LVLMM-Aware Multimodal Retrieval for RAG-Based Medical Diagnosis with General-Purpose Models

Anonymous ACL submission

Abstract

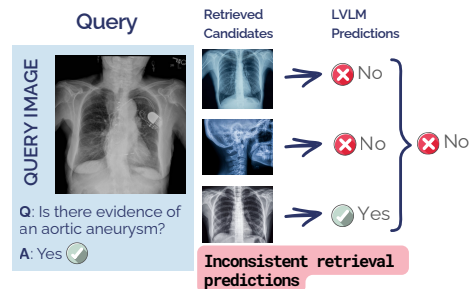
Retrieving visual and textual information from medical literature and hospital records can enhance diagnostic accuracy for clinical image interpretation. However, multimodal retrieval-augmented diagnosis is highly challenging. We explore a lightweight mechanism for enhancing diagnostic performance of retrieval-augmented LVLMMs. We train an LVLMM-aware multimodal retriever, such that the retriever learns to return images and texts that guide the LVLMM toward correct predictions. In our low-resource setting, we perform only lightweight fine-tuning with small amounts of data, and use only general-purpose backbone models, achieving competitive results in clinical classification and VQA tasks compared to medically pre-trained models with extensive training. In a novel analysis, we highlight a previously unexplored class of errors that we term inconsistent retrieval predictions: cases where different top-retrieved images yield different predictions for the same target. We find that these cases are challenging for all models, even for non-retrieval models, and that our retrieval optimization mechanism significantly improves these cases over standard RAG. However, our analysis also sheds light on gaps in the ability of LVLMMs to utilize retrieved information for clinical predictions.¹

1 Introduction

Inferring diagnoses from medical imagery is a fundamental part of clinical decision-making. Large Vision Language Models (LVLMMs) have been widely explored for medical diagnosis (Thawakar et al., 2024; Wu et al., 2023; Zhang et al., 2023b; Moor et al., 2023; Li et al., 2023). To improve the performance of LVLMMs in the medical field, retrieval augmentation (RAG) has been adopted and has shown promising results, providing more accurate and also explainable methods (He et al., 2024a; Xia et al., 2024b,a; Wu et al., 2025).

¹Code and models available at [redacted for anonymity].

Fine-Tuned Multimodal RAG



CLARE: Clinical LVLMM - Aware Retrieval

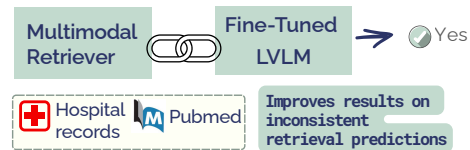


Figure 1: We optimize a multimodal retriever and a Large Vision-Language Model for medical tasks. We achieve competitive results without resource-intensive medical pre-training and significantly improve performance on challenging cases where different retrieved images lead to inconsistent retrieval prediction.

In this work, we explore a lightweight fine-tuning approach using a *general-purpose* LVLMM with a *general-purpose* multimodal retriever for medical diagnosis tasks. Unlike standard multimodal RAG, we train an LVLMM-aware multimodal retriever to find knowledge—clinical images, captions, and reports from both medical literature and hospital records—that leads to correct LVLMM predictions. Poorly optimized retrieval mechanisms can mislead models (Yoran et al., 2023; Sun et al., 2024). As shown in Figure 2, a RAG model incorrectly classifies a benign ultrasound image when conditioned on a cancerous retrieved image, whereas our LVLMM-aware retrieval optimization leads to the correct predicted diagnosis.

Our method involves sequential multimodal training with a dual-head retriever architecture and

059 a customized retrieval loss and visual question
060 answering (VQA) training recipe. Our method
061 achieves competitive results in medical image clas-
062 sification and VQA. Importantly, our focus in this
063 work is not to chase SOTA results but to shed light
064 on several interesting observations. One, that a
065 simple and effective LVLM-aware retrieval opti-
066 mization mechanism can lead to a substantial boost
067 in results over current multimodal RAG methods
068 and should thus be more widely adopted for med-
069 ical diagnosis; beyond downstream performance,
070 our analysis also shows improved relevance of re-
071 trieval candidates. Second, that this mechanism
072 achieves competitive results by using models with
073 no medical pre-training (for neither the LVLM nor
074 the retriever) and only lightweight fine-tuning. This
075 resonates with recent findings showing general-
076 purpose LVLMs can rival their medical counter-
077 parts (Jeong et al., 2024) and has potentially impor-
078 tant implications considering the resource-intensive
079 nature of pre-training processes.

080 Third, previous related work in the general do-
081 main that tuned both retrieval and generation mod-
082 els required large-scale pre-training followed by
083 few-shot task-specific adaptation (Izacard et al.,
084 2023; Hu et al., 2023; Lin et al., 2024). In con-
085 trast, we conduct our optimization *directly on*
086 *downstream tasks* with no pretraining and only
087 lightweight fine-tuning (as few as 546 samples).
088 This provides first evidence about the feasibility
089 and utility of lightweight data-efficient generation-
090 aware retrieval optimization directly on down-
091 stream tasks as opposed to in the pre-training set-
092 ting. We are also the first to explore and show
093 the utility of such optimization in the multimodal
094 medical domain.

095 Fourth and importantly, we conduct a novel anal-
096 ysis that finds that our method helps in particular
097 to address a class of errors we term *inconsistent*
098 *retrieval predictions*. Inconsistent retrieval predic-
099 tions are cases in which for a given patient image,
100 each retrieved image leads the model to make dif-
101 ferent predictions. These cases commonly occur
102 across our experiments and are substantially more
103 difficult for models, both retrieval-augmented and
104 non-retrieval models. This instability with respect
105 to different retrieval candidates leads to high predic-
106 tion entropy/uncertainty across classes, degrading
107 not only overall RAG results but also their reli-
108 ability (Lambert et al., 2022). As we show, our re-
109 trieval mechanism significantly mitigates this issue
110 and achieves large improvement over standard fine-

111 tuned RAG on these cases. We further discover that
112 for a large proportion of these inconsistent cases,
113 at least one retrieved candidate enables the model
114 to predict the correct answer in an oracle setting
115 in which we provide the model with the correct
116 retrieved candidate. This suggests useful informa-
117 tion is being retrieved but is often lost among less
118 helpful candidates. To address this, we explored
119 using the powerful o3 model (OpenAI, 2025) as
120 a reranker to identify predictive candidates. We
121 found it did not close the performance gap to the or-
122 acle and was comparable or inferior to our method,
123 leaving significant room for future improvement.

124 **Our contributions** are three-fold:

- 125 • We perform LVLM-aware retrieval optimiza- 125
126 tion for a multimodal retriever and an LVLM 126
127 on medical classification and VQA tasks (in- 127
128 cluding text generation VQA). We use only 128
129 general-purpose backbones without medical 129
130 pretraining, and show that by retrieving rel- 130
131 evant medical evidence and training the re- 131
132 triever to guide the LVLM, we close much of 132
133 the gap to expensive medically pre-trained 133
134 models and achieve superior performance 134
135 compared to general-purpose RAG methods. 135
- 136 • Unlike previous general-domain methods that 136
137 required pre-training, we conduct lightweight 137
138 data-efficient optimization of a retriever and 138
139 LVLM directly on downstream tasks. 139
- 140 • We identify *inconsistent retrieval predictions*, 140
141 challenging cases where retrieved images lead 141
142 to different predictions for the same query. 142
143 Our method significantly improves perfor- 143
144 mance on these cases over standard RAG. 144

145 2 Related Work

146 **Retrieval-Augmented LLM-Aware Optimiza-** 146
147 **tion** has been explored primarily in unimodal (text 147
148 only) NLP with work such as ATLAS (Izacard 148
149 et al., 2023), which conducted large-scale pretrain- 149
150 ing of both reader and retriever models and eval- 150
151 uated in zero-shot and few-shot scenarios in the 151
152 general domain. REVEAL (Hu et al., 2023) ex- 152
153 tended this to the multimodal setting in the gen- 153
154 eral domain with extensive pre-training; REVEAL 154
155 used an encoder-decoder architecture with a T5 155
156 generator—to our knowledge, since REVEAL no 156
157 work has explored optimizing multimodal retriev- 157
158 ers and generators with causal decoder models and 158

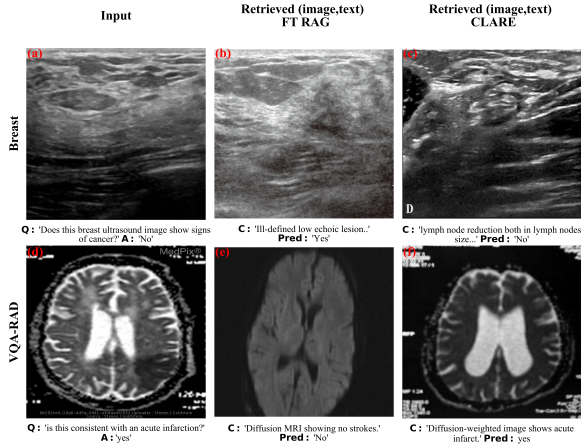


Figure 2: LVLMM-aware multimodal retrieval finetuning impacts inconsistent retrieval predictions. After LVLMM-aware multimodal retrieval finetuning, retrieved images show greater alignment with query image labels. In Breast, a cancer-free ultrasound query (a) initially retrieved a lesion image (b), however, after finetuning, the retrieved image is less directly related to the wrong label (c). In VQA-RAD, retrieval shifted from an unrelated medical condition (e) to an image depicting the condition of the query image.

modern VLMMs. Shi et al. (2023) proposed training textual retrievers based on reader performance, influencing subsequent retrieval alignment LLM-aware optimization methods. Lin et al. (2024) introduced another LLM-aware optimization variant and proposed first training the reader with retrieval augmentation, then training the retriever. Siriwardhana et al. (2023) demonstrated LLM-aware optimization’s effectiveness for domain adaptation. Our work differs by being the first to explore generator-aware retrieval optimization directly for downstream tasks, with minimal fine-tuning, as opposed to requiring large-scale pre-training. Prior work, such as ATLAS and REVEAL, follows a two-step approach—first conducting large-scale pre-training, then evaluating in few-shot or zero-shot settings, while we use a lightweight one-step approach instead. Our method includes a new loss and new training methodologies. Our extensive novel analyses are the first to shed light on inconsistent predictions and methods to help address them.

Multimodal Retrieval Augmentation in Medical Applications. Multimodal retrieval augmentation began with encoder-based architectures (Yuan et al., 2023), which integrated retrieved text and images with query modalities. With LVLMMs, He et al. (2024a) proposed retrieving labeled examples dur-

ing inference. Xia et al. (2024b) improved LVLMM RAG with contrastive learning and context selection strategies. Xia et al. (2024a) proposed domain-aware retrieval for diverse medical data, while Wu et al. (2025) combined retrieval augmentation with knowledge graphs. All current multimodal medical RAG methods are not LVLMM-aware, and the retriever is trained disjointly from the LVLMM.

3 Method

Overview. In this section, we present our methodology, termed CLARE (Clinical LVLMM-Aware Retrieval). Our task involves medical image classification and visual question answering. CLARE comprises two main components: a multimodal retriever and a reader. Given an input patient image and a diagnostic question, the multimodal retriever identifies relevant medical knowledge—images and their associated captions or hospital reports—that provide predictive information. The reader then analyzes the retrieved candidates along with the question to generate an answer. We optimize the reader to better leverage the retrieved information and the multimodal retriever to enhance its ability to select informative candidates for the reader. Figure 3 illustrates our framework.

Unlike previous work (He et al., 2024a; Xia et al., 2024b,a) we do not use medical pretraining, and instead we use readily available general LVLMMs (Pixtral and Qwen2-vl) and a general retriever, jinaclip (Xiao et al., 2024). The selected LVLMMs can process multiple images and texts, allowing us to incorporate both the retrieved image and the retrieved text. We also compare to Med-Flamingo (Moor et al., 2023), a medical LVLMM which is also able to process multiple images, unlike other medical LVLMMs available at the time of conducting our experiments. We now elaborate on the details.

3.1 Reader Fine-Tuning

We first fine-tune the reader with retrieval augmentation. This step aims to achieve two main goals: improving the reader’s performance on the dataset and teaching the reader to effectively utilize retrieved (image, text) pairs in context.

For an image-question tuple (i_d, q_d) where $d \in \{1, \dots, D\}$ from dataset D , we retrieve K relevant tuples of images and texts (reports or captions) (i_k, t_k) where $k \in \{1, \dots, K\}$ using our multimodal retriever module. To simplify notation, we denote $z_k = (i_k, t_k)$ as the k -th retrieved image-

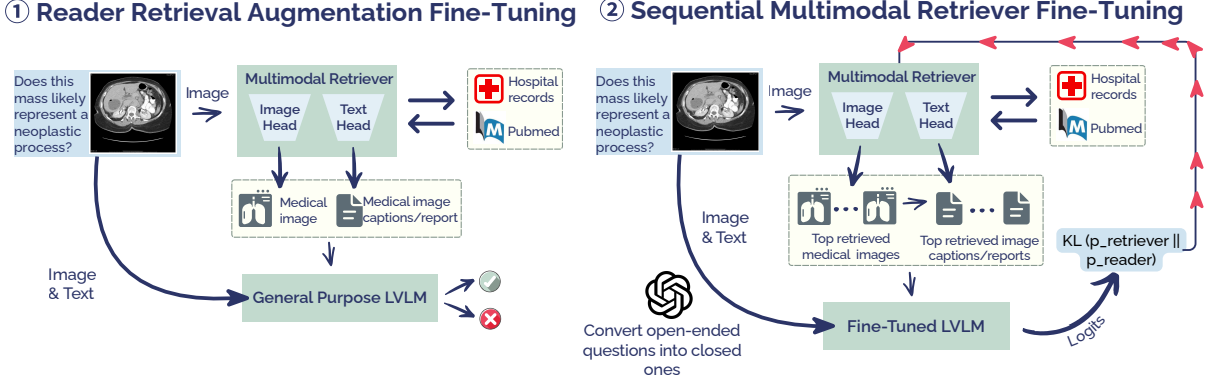


Figure 3: The two-phase CLARE training. The LVLM is first trained with a frozen retriever on augmented prompts containing retrieved image–caption/report pairs. With the LVLM frozen, the dual-head multimodal retriever is optimized using a KL divergence computed over a selected subset of the model’s logits. An LLM converts open-ended questions into closed form (during training only), which improves results.

text pair and $q_d = (i_d, q_d)$ as the query pair for example d . These retrieved pairs are prepended to the query pair to create augmented inputs. The augmented input consists of the retrieved context followed by the target question-image pair: $z_k \circ q_d$ for $k \in \{1, \dots, K\}$. We use supervised fine-tuning on the augmented input, minimizing the negative log-likelihood to predict the answer a_d :

$$\mathcal{L}(\theta) = - \sum_{d=1}^D \sum_{k=1}^K \log p_{\theta}(a_d | z_k \circ q_d)$$

where θ represents the parameters of the LVLM, and \circ denotes the concatenation operation.

3.2 Sequential Retriever Fine-Tuning

Our dual-head multimodal retriever has a text-based head for retrieving relevant (image, text) pairs based on textual similarity to the query, and an image-based head for retrieving pairs based on visual similarity. To train this retriever, we adopt a sequential training strategy, first optimizing the text retrieval head followed by the image retrieval head. During this process, the reader model remains frozen, focusing solely on improving the retriever’s embedding space. Similarly to Izacard et al. (2023) we minimize the KL divergence between the LVLM’s posterior distribution over retrieved pairs and the retriever’s distribution, however we do so directly on downstream tasks with only lightweight fine-tuning in the data-efficient regime without pre-training, and we customize the loss and the training recipe for open-ended questions leading to empirical gains.

Let $z_k = (\mathbf{i}_k, \mathbf{t}_k)$ denote the k -th retrieved image-text pair and $q = (\mathbf{i}, \mathbf{q})$ denote the query pair.

The posterior distribution reflects the model’s confidence in predicting the correct answer \mathbf{a} given a retrieved pair z_k and query q : $p_k \propto p_{\text{LVLM}}(\mathbf{a} | z_k \circ q)$, where $p_{\text{LVLM}}(\mathbf{a} | z_k \circ q)$ is the probability assigned by the LVLM to \mathbf{a} , and \circ denotes prepending the retrieved chunk to the query.

We sharpen the LVLM’s output distribution by restricting it to only the relevant class tokens from the benchmark’s possible answers. Specifically, we extract the model’s logits, select only the relevant tokens for each class, and apply a softmax to obtain a more discriminative distribution. We denote this class-restricted distribution as p_{LVLM_C} where $C = \{c_1, c_2, c_3, \dots\}$ is the set of chosen class tokens. The normalized posterior p_k is formulated as:

$$p_k = \frac{\exp(\log p_{\text{LVLM}_C}(\mathbf{a} | z_k \circ q))}{\sum_{i=1}^K \exp(\log p_{\text{LVLM}_C}(\mathbf{a} | z_i \circ q))},$$

where K is the total number of retrieved pairs.

The retriever’s distribution is defined by:

$$p_{\text{RETR}}(z | q) = \frac{\exp(s(z, q)/\tau)}{\sum_{k=1}^K \exp(s(z_k, q)/\tau)},$$

where $s(z, q)$ is the similarity score between the pair and the query, and temperature τ controls distribution sharpness. For the text retrieval head, the score is computed using the dot product between the index caption/report embeddings and the query image embedding, while for the image retrieval head, we use the index image embeddings.

For our loss we compute the KL-divergence between retriever and reader $\text{KL}(p_{\text{LVLM}_C} || p_{\text{RETR}})$:

$$\sum_{k=1}^K p_{\text{LVLM}_C}(z_k) \log \left(\frac{p_{\text{LVLM}_C}(z_k)}{p_{\text{RETR}}(z_k)} \right)$$

Finally, for visual question answering (VQA), we use the o3 model to convert open-ended questions into closed-ended ones (yes/no). See prompt in Appendix E and analysis of approach contribution in Appendix D.2. We apply this only during training², without modifying questions at inference time. Reformatted questions will be made publicly available. This reformulation improves results, akin to restricting the LVLM output distribution.

3.3 Inference

Following Shi et al. (2023) for a given query pair $q = (\mathbf{i}, \mathbf{q})$, we retrieve the top- K relevant chunks $z_k = (\mathbf{i}_k, \mathbf{t}_k)$ using the image as the query. Each chunk is prepended to the question, and the LVLM computes predictions for the augmented prompts in parallel. The final output probability is:

$$p_{\text{LVLM}}(a | q) = \sum_{k=1}^K p_{\text{LVLM}}(a | z_k \circ q) \cdot p_{\text{R}}(z_k | q),$$

where $p_{\text{R}}(z_k | q) = \frac{\exp(s(q, z_k))}{\sum_{j=1}^K \exp(s(q, z_j))}$, and $s(q, z_k)$ is the similarity score between the query and the retrieved chunk.

4 Experiments

4.1 Experimental Setup

Our datasets include real-world hospital datasets (BRSET (Nakayama et al., 2024) and VinDr-PCXR (Pham et al., 2022)) alongside a variety of classification and visual question answering datasets. We focus on a low-resource data-efficient setting (training sets ranging from 546 to 7007 samples in classification, 1790-19,700 in VQA). Medical image annotation is a resource-intensive task, demanding expert annotators whose availability is limited and expensive. Clinical AI also often has poor generalization across heterogeneous medical centers and patient populations, often requiring that models be trained for the specific context in which they are deployed (Casey, 2022). Such constraints often restrict researchers to datasets comprising only a few thousand samples for a given study. Retrieval augmentation is well-suited for this setting, as it

²Run once per question; the total cost is only a few dollars.

is known to benefit low-data regimes the most by leveraging external knowledge.

We consider binary, multi-class, and multi-label classification. These include BreastMNIST (“Breast”; breast ultrasound imaging, binary) (Al-Dhabyani et al., 2020), DermaMNIST (“Derma”; pigmented skin lesion images, multi-class) (Tschandl et al., 2018), RetinaMNIST (“Retina”; retinal fundus images, multi-class) (Liu et al., 2022), and two challenging real-world multi-label datasets: VinDr-PCXR (chest X-rays, 15 labels) (Pham et al., 2022) and BRSET (ophthalmology, 14 labels) (Nakayama et al., 2024). For VQA, we use widely adopted benchmarks: VQA-RAD (Lau et al., 2018), SLAKE-English (Liu et al., 2021), and PathVQA (He et al., 2020). Full details are provided in Appendix A.1.

Retrieval augmentation is supported by an external index constructed from PubMed and medical records: PMC-OA (Lin et al., 2023), MIMIC-CXR (Johnson et al., 2019a), and ROCO (Rückert et al., 2024). Full descriptions are available in Appendix A.2. Our experiments leverage two backbone models—Pixtral (12B) (Agrawal et al., 2024) and Qwen2-vl (7B parameters) (Wang et al., 2024)—with jina-clip-v1 (Xiao et al., 2024) serving as the general-purpose retriever’s visual head, embedding both the texts and images of the retrieval index. All runs were performed on a single L40s GPU. More details in Appendix B.

Baselines. We compare to several RAG baselines. RAD (He et al., 2024a), a retrieval-based approach for classification, where the label of the most similar training image is used as context; we use RAD with both Qwen2-VL and Pixtral. MMed-RAG (Xia et al., 2024a), a recent state-of-the-art multimodal RAG framework in which the retriever is optimized independently from the LVLM and is not aware of the LVLM during training. We evaluated MMed-RAG using backbones LLaVA-Med and LLaVA, both paired with a medically pre-trained retriever (CLIP). We also adopt a standard fine-tuned RAG baseline in which the retriever is fixed and only the reader is fine-tuned, corresponding to the first phase of our approach. Finally, we compare our retriever loss with the base Perplexity Distillation Loss (PDist) (Izacard et al., 2023), integrated into CLARE in place of our customized loss. Implementation details are in Appendix C. To compare to medically pre-trained LVLMs, we include competitive baselines: BiomedGPT (Luo et al., 2023); three LLaVA-Med (Li et al., 2023)

(a)

Backbone (Size)	Model	Breast		Derma		Retina		VinDr-PCXR		BRSET		Mean	
		ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1
LLaVA (7B)	MMed-RAG	.85	.84	.75	.30	.63	.46	.55	.11	.42	.30	.64	.40
Qwen2-vl (7B)	Reader	.83	.77	.68	.27	.60	.44	.48	.08	.34	.23	.59	.36
	RAD	.84	.79	.43	.34	.55	.40	.57	.09	.40	.25	.56	.37
	FT RAG	.85	.82	.71	.42	.62	.48	.55	.09	.48	.27	.64	.42
	CLARE + PDist	.86	.83	.73	.47	.63	.49	.57	.09	.47	.29	.65	.43
	CLARE	.87	.84	.76	.50	.65	.50	.57	.14	.49	.37	.67	.47
	CLARE _{oracle}	.87	.84	.85	.64	.69	.51	.64	.14	.52	.41	.71	.51
Pixtral (12B)	Reader	.82	.77	.75	.52	.55	.44	.48	.08	.44	.33	.61	.43
	RAD	.85	.80	.75	.47	.56	.41	.48	.09	.47	.30	.62	.41
	FT RAG	.88	.85	.79	.60	.57	.47	.49	.09	.47	.33	.64	.47
	CLARE + PDist	.88	.84	.80	.62	.60	.47	.50	.09	.45	.35	.65	.47
	CLARE	.90	.87	.80	.62	.60	.51	.56	.14	.51	.37	.67	.50
	CLARE _{oracle}	.93	.90	.83	.67	.63	.54	.67	.15	.57	.41	.73	.53

(b)

Backbone (Size)	Model	VQA-RAD		SLAKE		PathVQA		Mean	
		Closed	Open	Closed	Open	Closed	Open	Closed	Open
LLaVA (7B)	MMed-RAG	.74	.39	.87	.81	.90	.31	.84	.50
Qwen2-vl (7B)	Reader	.73	.41	.84	.80	.87	.25	.81	.49
	FT RAG	.76	.45	.88	.81	.91	.33	.85	.53
	CLARE + PDist	.77	.45	.89	.81	.92	.32	.86	.53
	CLARE	.79	.48	.90	.84	.93	.38	.87	.57
	CLARE _{oracle}	.86	-	.92	-	.95	-	.91	-
Pixtral (12B)	Reader	.72	.38	.83	.80	.87	.26	.81	.48
	FT RAG	.74	.41	.88	.81	.88	.31	.83	.51
	CLARE + PDist	.75	.41	.89	.81	.89	.32	.84	.51
	CLARE	.76	.45	.90	.84	.90	.36	.85	.55
	CLARE _{oracle}	.88	-	.92	-	.95	-	.92	-

Table 1: Results for classification (a) and VQA (b) with non-medically pretrained LVLm backbones (see Tables 3, 4 for comparisons to medical pre-training; MMed-RAG’s retrieval component is medically pre-trained). CLARE consistently leads to the best results. For VQA, we evaluate closed questions using the exact match metric and open questions using token recall. Using an oracle reranker (for classification and closed VQA; §4.2) demonstrates that even larger improvements are achievable with CLARE’s top-retrieved images in some datasets.

	MMed-RAG	LLaVaMed variants	MedDr variants	GSCo	MedVInT variants
Medical Pre-Training Size	600K	600K	255K	255K	177K

Table 2: Medical pre-training sizes of existing LVLms.

variants; two LVLm variants based on PMC CLIP (MedVInT-TE and MedVInT-TD) (Zhang et al., 2023b); and three InternVL(Chen et al., 2024) variants: MedDr, MedDr + RAD (He et al., 2024a), GSCo (He et al., 2024b). Details in Appendix C.

4.2 Results

Table 1a and Table 1b show the effectiveness of CLARE on five medical classification and three

VQA benchmarks. We compare with methods of similar computational cost; all LVLms were not extensively pre-trained on medical data, whereas the MMed-RAG retriever was. CLARE consistently outperforms MMed-RAG, RAD, the fine-tuned RAG baseline, CLARE + PDist loss, and the Reader baseline (LVLm without retrieval). In addition, we evaluated our model in an oracle scenario. Instead of the final fusion of logits (Section 3.3), we applied an oracle that, given the model responses for different retrieved images, selects the correct answer if it exists.³ Our findings in Tables 1a and 1b show that the correct answer exists in

³For open-ended questions we do not conduct this analysis as there is no simple boundary between correct/incorrect.

Model	Breast	Derma	VinDr	BRSET
MMed-RAG _{LLaVA-Med}	.89	.79	.11	.33
MedVInT-TE	.88	.78	-	-
MedVInT-TD	.90	.80	-	-
MedDr _{InternVL}	.72	-	.08	.08
MedDr + RAD _{InternVL}	.88	-	-	-
GSCo _{InternVL}	.93	-	.09	.33
CLARE _{Qwen2-vl}	.87	.76	.14	.37
CLARE _{Pixtral}	.90	.80	.14	.37

Table 3: Comparison to prior reported results of medical pre-trained models, for binary and multiclass classification (Breast and Derma; accuracy), and for multi-label classification (VinDr-PCXR and BRSET; F1).

Model	SLAKE		PathVQA	
	Closed	Open	Closed	Open
MMed-RAG _{LLaVA-Med}	.89	.84	.92	.39
BiomedGPT-B	.90	.85	.88	.28
LLaVA-Med _{LLaVA}	.83	.85	.91	.38
LLaVA-Med _{Vicuna}	.85	.83	.92	.39
LLaVA-Med _{BioMedCLIP}	.87	.87	.91	.39
CLARE _{Qwen2-vl}	.90	.84 (.84)	.93	.38 (.37)
CLARE _{Pixtral}	.90	.84 (.82)	.90	.36 (.35)

Table 4: Comparison to medical pre-trained models for visual question answering. We report the token recall metric, reported for LLaVA-Med variants and token F1 reported for BiomedGPT (in red).

a large proportion of cases, yielding superior results. This shows that in these cases, CLARE is able to surface in its top retrieved images an image which could lead to a correct prediction, however the information gets lost in the process of prediction fusion with other retrieved images. Simply taking the label with highest confidence or label with highest average confidence, underperforms. This motivates us to explore the o3 multimodal reasoning model to detect this image (Section 5).

Benchmarking versus medically pretrained LLMs for classification. In Table 3, we include previously reported results for leading medically pretrained LLM methods (detailed in Section 4.1). CLARE demonstrates results competitive with models that underwent extensive medical pretraining (pretraining sizes in Table 2). CLARE with the Pixtral backbone matches or exceeds all models except GSCo on the Breast dataset. To our knowledge, no medically pretrained LLMs have been

evaluated on the Retina benchmark, and we are not aware of any prior LLM results on our other classification benchmarks. For binary and multi-class classification, we report only accuracy (ACC), as prior works did not report F1 scores. Conversely, for multi-label classification, we report only F1 scores. We also evaluated a medical baseline, MedFlamingo, which generally underperformed and showed limited benefit from retrieval augmentation. We attribute this to its relatively older LLaMA backbone (Touvron et al., 2023). Additionally, we experimented with using BiomedCLIP (Johnson et al., 2019b) as the retriever for our general-purpose LLMs, which exhibited a similar trend. Full results in Appendix D.4.

Benchmarking against pretrained LLMs for visual question answering. In Table 4, we compare our approach with medically pretrained LLM methods (detailed in Section 4.1) for VQA. CLARE achieves competitive performance relative to extensively pretrained models. For closed-question tasks, CLARE outperforms or matches existing models. For open-ended question answering, CLARE with the Qwen2-vl backbone matches LLaVA-Med_{Vicuna} on the SLAKE dataset and LLaVA-Med_{LLaVA} on the PathVQA dataset. Note that we do not report results on VQA-RAD, since our method was trained using an internal split of training and validation, whereas VQA-RAD provides only official training and test splits.

5 Analysis

CLARE boosts empirically challenging cases. We observe that the performance gap is substantially larger for predictions classified as inconsistent retrieval predictions (Table 6 shows the proportion of such cases). We define an inconsistent retrieval prediction as one in which the model predicts different labels for retrieved candidates given the same query (see Figure 1), i.e., if the model’s predictions vary across retrieved candidates. Conversely, if the model predicts the same label for all retrieved candidates, the instance is consistent.

Consistent and inconsistent prediction sets are established after the initial reader training stage, and their performance is assessed following LLM-aware multimodal retrieval fine-tuning. Detailed results are presented in Table 5. Overall, CLARE improves performance on inconsistent retrieval predictions by +0.12 in accuracy and +0.13 in F1 score when using the Qwen2-vl backbone, and by +0.09

Backbone	Model	Breast		Derma		Retina		VinDr		BRSET		VQA-RAD	SLAKE	Mean	Mean	
		ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	Closed Acc	Closed Acc	ACC	F1	
Qwen2-vl	Inconsistent	Reader	.40	.40	.50	.25	.35	.31	.41	.07	.33	.21	.43	.70	.40	.25
		RAG	.40	.40	.52	.36	.40	.28	.46	.07	.49	.25	.33	.84	.45	.27
		CLARE	.80	.80	.62	.46	.52	.29	.48	.16	.50	.36	.47	.86	.58	.41
	Consistent	Reader	.84	.77	.83	.29	.63	.45	.56	.09	.34	.22	.77	.85	.64	.36
		RAG	.86	.82	.90	.49	.66	.50	.66	.10	.43	.23	.82	.89	.70	.43
		CLARE	.88	.84	.91	.50	.66	.55	.66	.10	.41	.37	.82	.91	.70	.47
Pixtral	Inconsistent	Reader	.60	.59	.45	.33	.40	.32	.38	.07	.46	.28	.66	.64	.46	.32
		RAG	.72	.71	.44	.45	.35	.37	.38	.08	.39	.26	.72	.65	.46	.37
		CLARE	.84	.84	.50	.50	.42	.46	.49	.13	.45	.35	.77	.68	.54	.46
	Consistent	Reader	.86	.80	.80	.56	.56	.44	.66	.08	.61	.37	.74	.85	.70	.45
		RAG	.91	.87	.84	.67	.62	.50	.71	.08	.64	.38	.75	.91	.74	.50
		CLARE	.91	.87	.85	.67	.62	.53	.71	.10	.64	.42	.75	.91	.75	.52

Table 5: Analysis of LVLM-aware multimodal retrieval fine-tuning effect: for inconsistent retrieval predictions, we can see that applying LVLM-aware multimodal retrieval fine-tuning boosts performance significantly for all datasets. For consistent retrieval predictions, the results are slightly improved. Note that inconsistent retrievals are challenging to all models, even non-retrieval ones.

Model	Breast	Derma	Retina	VinDr	BRSET	VQARAD	SLAKE
Qwen2-vl	3%	51%	15%	50%	93%	12%	16%
Pixtral	16%	13%	16%	66%	70%	7%	8%

Table 6: Proportion inconsistent predictions cases, per dataset, for CLARE based on Qwen2-vl and Pixtral.

in both accuracy and F1 score with the Pixtral backbone, while also offering a slight improvement on the consistent. We also observed that inconsistent retrieval predictions are empirically more challenging for the reader, FT RAG and CLARE models, with a 10–20 point performance gap.

Performance with reranking. Results for inconsistent predictions remain relatively low (Table 5). However, in the oracle analysis reported above, we found that CLARE can sometimes retrieve images that lead to correct predictions, together with other images that lead to wrong predictions. We thus explore whether a state-of-the-art multimodal model, when used as an optional reranker, could close the gap toward oracle performance. Specifically, we investigated whether the o3 reasoning model (OpenAI, 2025) could select a retrieved image that leads to correct prediction. We provided o3 with an image to classify, along with four retrieved images+captions, and tasked it with identifying the one containing the most predictive information. Our results show that o3’s performance generally surpasses simply taking the image with highest confidence, but o3 is generally inferior compared to fusing logits in CLARE. We conclude that reranking in this setting remains a significant challenge

even for a frontier model. Results in Appendix D.3.

Retrieval relevance and additional analyses.

We explore the utility of our method in retrieving more relevant candidates. In our setting where no retrieval ground-truth labels are available, we use GPT-5.2 as a judge to compare relevance of retrieved candidates before and after the LVLM-aware retriever fine-tuning. Our method improves candidate relevance with an average 17% win rate, as detailed in Appendix D. We further evaluate the contribution of each component in Appendix D.2 and present robustness ablations in Appendix D.5.

6 Conclusion

We demonstrated that LVLM-aware multimodal retrieval achieves competitive medical diagnosis performance through lightweight, data-efficient fine-tuning, without medical pre-training. Future work may explore whether this approach generalizes as a cost-effective alternative to domain-specific pre-training when task-specific data and domain knowledge bases are available but large-scale pre-training is prohibitive. In addition, we identified and substantially mitigated inconsistent retrieval predictions, where different retrieved candidates lead to conflicting diagnoses. Our oracle analysis revealed a considerable performance gap between actual results and what was theoretically achievable using the retrieved images, providing a foundation for future research on closing this gap. Future work may also look into the prevalence and impact of inconsistent retrieval predictions more broadly.

535 Limitations

536 Scope of evaluation. Our evaluation focuses on
537 classification and visual question-answering tasks,
538 including both open questions (a text-generation
539 setting) and closed questions (closed-form answer
540 selection). We did not evaluate our method on
541 report generation, which requires different capa-
542 bilities (longer-form generation, clinical writing
543 conventions) and would benefit from dedicated in-
544 vestigation. Additionally, while our method shows
545 consistent improvements across diverse imaging
546 modalities (ultrasound, fundus photography, X-ray,
547 histopathology, dermatoscopy, CT) and tasks (vi-
548 sual question answering, malignant lesion classifi-
549 cation, diabetic retinopathy severity grading), eval-
550 uation on additional specialized modalities (e.g.,
551 PET, Microscopy, OCT) and tasks (e.g., organ clas-
552 sification, retinal OCT disease classification) would
553 further validate generalizability.

554 References

555 Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna,
556 Baptiste Bout, Devendra Chaplot, Jessica Chud-
557 novsky, Diogo Costa, Baudouin De Monicault,
558 Saurabh Garg, Theophile Gervet, and 1 others. 2024.
559 Pixtral 12b. *arXiv preprint arXiv:2410.07073*.

560 Walid Al-Dhabyani, Mohammed Gomaa, Hussien
561 Khaled, and Aly Fahmy. 2020. Dataset of breast
562 ultrasound images. *Data in brief*, 28:104863.

563 Ross Casey. 2022. [\[link\]](#).

564 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo
565 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,
566 Xizhou Zhu, Lewei Lu, and 1 others. 2024. Internvl:
567 Scaling up vision foundation models and aligning
568 for generic visual-linguistic tasks. In *Proceedings of*
569 *the IEEE/CVF conference on computer vision and*
570 *pattern recognition*, pages 24185–24198.

571 Sunan He, Yuxiang Nie, Zhixuan Chen, Zhiyuan Cai,
572 Hongmei Wang, Shu Yang, and Hao Chen. 2024a.
573 Meddr: Diagnosis-guided bootstrapping for large-
574 scale medical vision-language learning. *CoRR*.

575 Sunan He, Yuxiang Nie, Hongmei Wang, Shu Yang, Yi-
576 hui Wang, Zhiyuan Cai, Zhixuan Chen, Yingxue Xu,
577 Luyang Luo, Huiling Xiang, and 1 others. 2024b.
578 Gsco: Towards generalizable ai in medicine via
579 generalist-specialist collaboration. *arXiv preprint*
580 *arXiv:2404.15127*.

581 Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and
582 Pengtao Xie. 2020. Pathvqa: 30000+ questions for
583 medical visual question answering. *arXiv preprint*
584 *arXiv:2003.10286*.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
Weizhu Chen. 2021. *Lora: Low-rank adaptation of*
large language models. *Preprint*, arXiv:2106.09685.

Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-
Wei Chang, Yizhou Sun, Cordelia Schmid, David A
Ross, and Alireza Fathi. 2023. Reveal: Retrieval-
augmented visual-language pre-training with multi-
source multimodal knowledge memory. In *Proceed-*
ings of the IEEE/CVF conference on computer vision
and pattern recognition, pages 23369–23379.

Gautier Izacard and Edouard Grave. 2020. Leverag-
ing passage retrieval with generative models for
open domain question answering. *arXiv preprint*
arXiv:2007.01282.

Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas
Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-
Yu, Armand Joulin, Sebastian Riedel, and Edouard
Grave. 2023. Atlas: Few-shot learning with retrieval
augmented language models. *Journal of Machine*
Learning Research, 24(251):1–43.

Daniel P Jeong, Pranav Mani, Saurabh Garg, Zachary C
Lipton, and Michael Oberst. 2024. The lim-
ited impact of medical adaptation of large lan-
guage and vision-language models. *arXiv preprint*
arXiv:2411.08870.

Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz,
Nathaniel R Greenbaum, Matthew P Lungren, Chih-
ying Deng, Roger G Mark, and Steven Horng.
2019a. Mimic-cxr, a de-identified publicly available
database of chest radiographs with free-text reports.
Scientific data, 6(1):317.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019b.
Billion-scale similarity search with GPUs. *IEEE*
Transactions on Big Data, 7(3):535–547.

Benjamin Lambert, Florence Forbes, Alan Tucholka,
Senan Doyle, Harmonie Dehaene, and Michel Dojat.
2022. Trustworthy clinical ai solutions: a unified
review of uncertainty quantification in deep learning
models for medical image analysis. *arXiv preprint*
arXiv:2210.03736.

Jason J Lau, Soumya Gayen, Asma Ben Abacha, and
Dina Demner-Fushman. 2018. A dataset of clini-
cally generated visual questions and answers about
radiology images. *Scientific data*, 5(1):1–10.

Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto
Usuyama, Haotian Liu, Jianwei Yang, Tristan Nau-
mann, Hoifung Poon, and Jianfeng Gao. 2023. Llava-
med: Training a large language-and-vision assistant
for biomedicine in one day. *Advances in Neural In-*
formation Processing Systems, 36:28541–28564.

Weixiong Lin, Ziheng Zhao, Xiaoman Zhang, Chaoyi
Wu, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023.
Pmc-clip: Contrastive language-image pre-training
using biomedical documents. In *International Con-*
ference on Medical Image Computing and Computer-
Assisted Intervention, pages 525–536. Springer.

642	Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Richard James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, and 1 others. 2024. Ra-dit: Retrieval-augmented dual instruction tuning. In <i>The Twelfth International Conference on Learning Representations</i> .	699
643		700
644		701
645		702
646		
647		
648	Bo Liu, Li-Ming Zhan, Li Xu, Lin Ma, Yan Yang, and Xiao-Ming Wu. 2021. Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering. In <i>2021 IEEE 18th international symposium on biomedical imaging (ISBI)</i> , pages 1650–1654. IEEE.	
649		
650		
651		
652		
653		
654	Ruhan Liu, Xiangning Wang, Qiang Wu, Ling Dai, Xi Fang, Tao Yan, Jaemin Son, Shiqi Tang, Jiang Li, Zijian Gao, and 1 others. 2022. Deepdrid: Diabetic retinopathy—grading and image quality estimation challenge. <i>Patterns</i> , 3(6).	
655		
656		
657		
658		
659	Yizhen Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. 2023. Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. <i>arXiv preprint arXiv:2308.09442</i> .	
660		
661		
662		
663		
664	Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakkka, Eduardo Pontes Reis, and Pranav Rajpurkar. 2023. Med-flamingo: a multimodal medical few-shot learner. In <i>Machine Learning for Health (MLAH)</i> , pages 353–367. PMLR.	
665		
666		
667		
668		
669		
670	LF Nakayama, M Goncalves, L Zago Ribeiro, H Santos, D Ferraz, F Malerbi, and 1 others. 2024. A brazilian multilabel ophthalmological dataset (brset). 2023. URL: https://physionet.org/content/brazilian-ophthalmological/1.0.0/ [accessed 2024-08-14].	
671		
672		
673		
674		
675	OpenAI. 2025. Openai o3 and o4-mini system card .	
676		
677	H Hieu Pham, T Thanh Tran, and Ha Quy Nguyen. 2022. Vindr-pcxr: An open, large-scale pediatric chest x-ray dataset for interpretation of common thoracic diseases. <i>PhysioNet (version 1.0. 0)</i> , 10(2).	
678		
679		
680	Johannes Rückert, Louise Bloch, Raphael Brüngel, Ahmad Idrissi-Yaghir, Henning Schäfer, Cynthia S Schmidt, Sven Koitka, Obioma Pelka, Asma Ben Abacha, Alba G. Seco de Herrera, and 1 others. 2024. Rocov2: Radiology objects in context version 2, an updated multimodal image dataset. <i>Scientific Data</i> , 11(1):688.	
681		
682		
683		
684		
685		
686		
687	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Replug: Retrieval-augmented black-box language models. <i>arXiv preprint arXiv:2301.12652</i> .	
688		
689		
690		
691		
692	Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering. <i>Transactions of the Association for Computational Linguistics</i> , 11:1–17.	
693		
694		
695		
696		
697		
698		
	Jiashuo Sun, Jihai Zhang, Yucheng Zhou, Zhaochen Su, Xiaoye Qu, and Yu Cheng. 2024. Surf: Teaching large vision-language models to selectively utilize retrieved information. In <i>EMNLP</i> .	699
		700
		701
		702
	Omkar Chakradhar Thawakar, Abdelrahman M Shaker, Sahal Shaji Mullappilly, Hisham Cholakkal, Rao Muhammad Anwer, Salman Khan, Jorma Laaksonen, and Fahad Khan. 2024. Xraygpt: Chest radiographs summarization using large medical vision-language models. In <i>Proceedings of the 23rd workshop on biomedical natural language processing</i> , pages 440–448.	703
		704
		705
		706
		707
		708
		709
		710
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	711
		712
		713
		714
		715
		716
	Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. 2018. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. <i>Scientific data</i> , 5(1):1–9.	717
		718
		719
		720
	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	721
		722
		723
		724
		725
		726
	Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023. Towards generalist foundation model for radiology by leveraging web-scale 2d&3d medical data. <i>arXiv preprint arXiv:2308.02463</i> .	727
		728
		729
		730
	Yinan Wu, Yuming Lu, Yan Zhou, Yifan Ding, Jingping Liu, and Tong Ruan. 2025. Mkgf: A multi-modal knowledge graph based rag framework to enhance lvlms for medical visual question answering. <i>Neurocomputing</i> , page 129999.	731
		732
		733
		734
		735
	Peng Xia, Kangyu Zhu, Haoran Li, Tianze Wang, Weijia Shi, Sheng Wang, Linjun Zhang, James Zou, and Huaxiu Yao. 2024a. Mmed-rag: Versatile multimodal rag system for medical vision language models. <i>arXiv preprint arXiv:2410.13085</i> .	736
		737
		738
		739
		740
	Peng Xia, Kangyu Zhu, Haoran Li, Hongtu Zhu, Yun Li, Gang Li, Linjun Zhang, and Huaxiu Yao. 2024b. Rule: Reliable multimodal rag for factuality in medical vision language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1081–1093.	741
		742
		743
		744
		745
		746
	Han Xiao, Georgios Mastrapas, and Bo Wang. 2024. Jina clip: Your clip model is also your text retriever. In <i>Multi-modal Foundation Model meets Embodied AI Workshop@ ICML2024</i> .	747
		748
		749
		750
	Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. 2023. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. <i>Scientific Data</i> , 10(1):41.	751
		752
		753
		754
		755

Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. Making retrieval-augmented language models robust to irrelevant context. *arXiv preprint arXiv:2310.01558*.

Zheng Yuan, Qiao Jin, Chuanqi Tan, Zhengyun Zhao, Hongyi Yuan, Fei Huang, and Songfang Huang. 2023. Ramm: Retrieval-augmented biomedical visual question answering with multi-modal pre-training. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 547–556.

Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, and 1 others. 2023a. Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. *arXiv preprint arXiv:2303.00915*.

Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023b. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. *Llamafactory: Unified efficient fine-tuning of 100+ language models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

A Dataset

A.1 Dataset for Evaluation

Our datasets include real-world hospital datasets (BRSET (Nakayama et al., 2024) and VinDr-PCXR (Pham et al., 2022)) alongside a variety of classification and visual question answering datasets. We focus on a low-resource data-efficient setting (training sets ranging from 546 to 7007 samples in classification, 1790-19,700 in VQA). Medical image annotation is a resource-intensive task, demanding expert annotators whose availability is limited and expensive. Clinical AI also often has poor generalization across heterogeneous medical centers and patient populations, often requiring that models be trained for the specific context in which they are deployed (Casey, 2022). Such constraints often restrict researchers to datasets comprising only a few thousand samples for a given study. Retrieval augmentation is well-suited for this setting, as it is known to benefit low-data regimes the most by leveraging external knowledge to compensate for sparse training samples. For the classification tasks, we used the following datasets:

BreastMNIST (nickname: Breast) (Al-Dhabyani et al., 2020) is a binary classification dataset of breast ultrasound. Following the official split, we use 546 for training, 78 for validation and 156 for testing.

RetinaMNIST (nickname: Retina) (Liu et al., 2022) is a multi-label classification dataset of retina Fundus Camera. Following the official split, we use 1,080 for training, 120 for validation and 400 for testing.

DermaMNIST (nickname: Derma) (Tschandl et al., 2018) is a multi-class classification dataset of common pigmented skin lesions. Following the official split, we use 7,007 for training, 1,003 for validation and 2,005 for testing. The dataset contains clinical images from 7 different diagnostic categories: actinic keratoses, basal cell carcinoma, benign keratosis, dermatofibroma, melanoma, melanocytic nevi, and vascular lesions.

VinDr-PCXR (Pham et al., 2022) is a large, open pediatric chest X-ray dataset collected in Vietnam (2020–2021) containing 9,125 posteroanterior scans from patients under 10 years old. It provides both lesion-level bounding-box annotations for 36 findings and image-level labels for multi-label of 15 diagnoses, curated by experienced radiologists. We follow the official split, which includes 7,728 training and 1,397 test images, with an additional internal split of the training data into 85% for training and 15% for validation.

BRSET (Nakayama et al., 2024) is a Brazilian multilabel ophthalmological dataset comprising retinal fundus images annotated with 14 distinct pathological findings. The dataset, composed of 16,266 images, presents a challenging real-world scenario for multilabel classification. There is no publicly available split; therefore, we apply an internal split for train, validation, and test sets (70% train, 10% validation, and 20% test).

Breast, Derma, and Retina taken from the large-scale MNIST-like collection of standardized biomedical images, including 12 datasets for 2D and 6 datasets for 3D. We especially used MedMNIST+ (Yang et al., 2023), which is a higher resolution extension of the original MedMNIST. We used the highest resolution of 224×224 .

For medical question answering, we have three well-known datasets:

VQA-RAD (Lau et al., 2018) is a clinician-annotated visual question answering dataset focused on radiology images (e.g., X-ray, CT, MRI). It pairs each image with both open-ended and

(a)

Backbone (Size)	Model	Breast		Derma		Retina		VinDr-PCXR		BRSET		Mean	
		ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1
Qwen-vl (7B)	FT RAG	.85	.82	.71	.42	.62	.48	.55	.09	.48	.27	.64	.42
	Text Retriever Head Only	.86	.83	.74	.46	.62	.48	.56	.13	.48	.32	.65	.44
	Image Retriever Head Only	.87	.83	.75	.48	.64	.49	.56	.13	.48	.30	.66	.45
	CLARE	.87	.84	.76	.50	.65	.50	.57	.14	.49	.37	.67	.47
Pixtral (12B)	FT RAG	.88	.85	.79	.60	.57	.47	.49	.09	.47	.33	.64	.47
	Text Retriever Head Only	.89	.86	.79	.61	.58	.48	.55	.13	.50	.35	.66	.49
	Text Retriever Head Only	.90	.86	.79	.60	.59	.48	.54	.13	.49	.35	.66	.48
	CLARE	.90	.87	.80	.62	.60	.51	.56	.14	.51	.37	.67	.50

(b)

Backbone (Size)	Model	VQA-RAD		SLAKE		PathVQA		Mean	
		Closed	Open	Closed	Open	Closed	Open	Closed	Open
Qwen2-vl (7B)	FT RAG	.76	.45	.88	.81	.91	.33	.85	.53
	Text Retriever Head Only	.77	.46	.89	.82	.93	.35	.86	.54
	Image Retriever Head Only	.78	.47	.89	.82	.92	.34	.86	.54
	Closed Question Only	.78	.47	.90	.83	.93	.37	.87	.56
	CLARE	.79	.48	.90	.84	.93	.38	.87	.57
Pixtral (12B)	FT RAG	.74	.41	.88	.81	.88	.31	.83	.51
	Text Retriever Head Only	.76	.44	.90	.83	.89	.33	.85	.53
	Image Retriever Head Only	.76	.44	.89	.82	.90	.32	.85	.53
	Closed Question Only	.76	.45	.90	.84	.90	.36	.85	.55
	CLARE	.78	.47	.90	.83	.93	.37	.87	.56

Table 7: Ablation results for classification (a) and visual question answering (b). We demonstrate the necessity of each stage in our training paradigm. The first stage, FT RAG, trains only the reader. Subsequently, training the text retriever improves performance, and finally, completing the process with the image retriever (CLARE) yields the best results. We also evaluate training only on closed questions—without applying the o3 model to convert open questions into closed ones—and observe a gain in open-question performance in this setting

yes/no questions. Following the official split, we use 1,753 questions for training and 453 for testing. In addition, we further partition the training set into 85% for training and 15% for validation.

PathVQA (He et al., 2020) is a pathology-focused medical VQA dataset built from textbooks and the PEIR digital library, comprising 4,289 images and 32,632 question–answer pairs spanning both open-ended and yes/no questions. We follow the official split provided by the authors: 19,700 for training, 6,260 for validation, and 6,720 for testing.

SLAKE (Liu et al., 2021) is a bilingual medical VQA dataset. We use the English subset, *SLAKE-English*, which comprises 642 radiology images and 7.03k English QA pairs. We follow the official train/validation/test split of 4.92k/1.05k/1.06k QA pairs, and include all question types, covering both open-ended and closed-ended formats.

A.2 Dataset for Index

For construction of the Index we used three large datasets of (image, text) pairs: PMC-OA (Lin et al., 2023), ROCO (Rückert et al., 2024) and MIMIC-CXR (Johnson et al., 2019a):

PMC-OA: is a large-scale dataset that contains 1.65M image-text pairs. The figures and captions from PubMed Central, 2,478,267 available papers are covered.

ROCO: is an image-caption dataset collected from PubMed. It filters out all the compound or non-radiological images, and consists of 81K samples.

MIMIC-CXR: is the largest chest X-ray dataset, containing 377,110 samples (image-report pairs). Each image is paired with a clinical report describing findings from doctors.

Backbone	Retrieval type	Breast		Derma		Retina		VQA-RAD		Mean	
		ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1
Flip a coin		.46	.30	.14	.10	.17	.16	.50		.35	.31
Pixtral	WO Q img	.75	.42	.66	.18	.47	.33	.67		.67	.49
	Reader only	.82	.66	.75	.53	.56	.44	.72		.74	.65
	WO Retrieval	.84	.65	.78	.55	.58	.48	.72		.74	.65
	Random Retrieval	.84	.69	.78	.56	.56	.45	.70		.74	.65
Qwen2-vl	WO Q img	.81	.52	.56	.20	.47	.20	.70		.68	.48
	Reader only	.82	.63	.68	.27	.61	.44	.74		.74	.59
	WO Retrieval	.86	.70	.70	.33	.62	.42	.76		.76	.62
	Random Retrieval	.85	.67	.73	.39	.63	.42	.77		.76	.62

Table 8: We evaluate our models in three scenarios: (1) the model is not given the query image and must rely solely on the retrieved context for prediction; (2) the model is not given any retrieval; and (3) the model is given a retrieved context that is randomly chosen and may mislead it. We use two baselines: a random classifier (uniformly selects a class) and a reader-only model (trained without retrieval). We evaluate on Breast, Derma, Retina, and the closed-question subset of VQA-RAD.

Backbone	Model	Breast		Derma		Retina	
		ACC	F1	ACC	F1	ACC	F1
Qwen2-vl	Top-1	.50	.60	.32	.25	.38	.32
	Top-1 logits	.60	.58	.43	.35	.37	.38
	o3	.77	.80	.42	.32	.35	.33
	o3 multi-image	.77	.80	.43	.33	.43	.35
	CLARE	.77	.80	.57	.43	.51	.47
Pixtral	Top-1	.72	.75	.46	.48	.29	.33
	Top-1 logits	.64	.63	.40	.40	.29	.32
	o3	.73	.73	.42	.37	.25	.30
	o3 multi-image	.73	.73	.43	.42	.34	.37
	CLARE	.76	.77	.45	.47	.34	.37

Table 9: Study of o3 as a reranker. We evaluate o3’s ability to rerank the inconsistent predictions, selecting which of the retrieved (image, caption) pairs has the most predictive information for the query image. We compared the performance of choosing the highest similarity retrieved (image, caption) pair, choosing the highest confidence (top logits), using only the caption for reranking, and using our current strategy of fusing logits from all predictions. The metrics used are accuracy and F1 score.

B Implementation Details

B.1 Retrieval Augmentation

As described earlier, for each image–question pair, we retrieve r image–text pairs, with $r = 4$. We use the input image as the query and Jina-CLIP (Xiao et al., 2024) as the retriever head in our multimodal retriever. The index is constructed with FAISS (Johnson et al., 2019b) over MIMIC-CXR, PMC-OA, and ROCO, which are fully described

in Section A.2 in Appendix A. To embed images, we use the Jina-CLIP visual head (the same model used for retrieval); to embed captions/reports, we use the Jina-CLIP text head. The index is stored in float16 due to storage constraints. We also experimented with BiomedCLIP (Zhang et al., 2023a).

B.2 Reader Fine-Tuning

We use Pixtral (12B) (Agrawal et al., 2024) and Qwen2-vl (7B) (Wang et al., 2024) as base models. We train the models with and without retrieval augmentation to assess their effect. The retrieval-augmented prompt is described in Appendix E. Note, we also tried Med-Flamingo (9B) (Moor et al., 2023) as a base model.

We fine-tune Pixtral and Qwen2-vl using LLaMA-Factory (Zheng et al., 2024) on a single NVIDIA L40 (48 GB) GPU for 10 epochs. We use a learning rate of 2×10^{-5} and apply Low-Rank Adaptation (LoRA) (Hu et al., 2021) for parameter-efficient fine-tuning. To all models (including baseline models), we conduct a grid search over batch size (2, 4, 6) and whether to freeze the vision head, selecting the best configuration by validation performance. For Med-Flamingo, we fine-tune using our codebase on a single NVIDIA L40 (48 GB) GPU for 10 epochs. Following the Med-Flamingo paper, the language model and image encoder are frozen, and only the Gated Cross-Attention layers and the Perceiver Resampler are optimized for stable and efficient learning. We use a learning rate of 2×10^{-5} .

935 **B.3 LVLM-Aware Multimodal Retrieval** 936 **Fine-Tuning**

937 We train the retrieval model to fetch relevant con-
938 text while keeping the reader frozen. We use Jina-
939 CLIP’s visual head as the base model for the multi-
940 modal retriever. The model is trained on a single
941 NVIDIA L40 (48 GB) GPU for 100 epochs with
942 a learning rate of 2×10^{-5} , freezing all layers ex-
943 cept the last ten. We set the number of retrieved
944 candidates per query to 50. In our experiments, re-
945 trieving more candidates further improved retriever
946 performance.

947 **B.4 Evaluation Metrics**

948 **Classification.** We report Accuracy (ACC) and
949 Macro F1:

$$950 \text{ACC} = \frac{1}{N} \sum_{n=1}^N \mathbf{1}[\hat{y}_n = y_n],$$

$$951 \text{MacroF1} = \frac{1}{C} \sum_{c=1}^C \frac{2 \text{Prec}_c \text{Rec}_c}{\text{Prec}_c + \text{Rec}_c},$$

952 where Prec_c and Rec_c are precision and recall
953 computed per class c , C is the number of classes,
954 and N is the number of examples.

955 **Visual Question Answering (VQA).** For open-
956 ended questions, we compute token-level Recall
957 and F1 between the predicted token set \hat{A} and the
958 reference token set A :

$$959 \text{Prec} = \frac{|\hat{A} \cap A|}{|\hat{A}|}, \quad \text{Rec} = \frac{|\hat{A} \cap A|}{|A|},$$

$$960 \text{F1} = \frac{2 \text{Prec} \cdot \text{Rec}}{\text{Prec} + \text{Rec}}.$$

961 For closed-ended questions, we use Exact Match.

962 **C Baseline Methods**

963 **C.1 RAG baseline**

964 We evaluate our method against several represen-
965 tative RAG-based and multimodal retrieval ap-
966 proaches to ensure a fair and comprehensive com-
967 parison.

968 **RAD** (He et al., 2024a) is a retrieval-based
969 method designed for classification tasks. During
970 inference, the most similar image from the training
971 set is retrieved, and its supervised label is directly
972 used as the prediction. We integrate RAD with both

of our backbone models (Qwen2-VL and Pixtral)
under the same data splits and retrieval settings.

973 **MMed-RAG** (Xia et al., 2024a) represents a re-
974 cent state-of-the-art multimodal RAG approach in
975 which the retriever and LVLM are trained inde-
976 pendently. We follow the official GitHub imple-
977 mentation, first training the CLIP module and then
978 performing DPO training using two backbone mod-
979 els: LLaVA-Med and LLaVA.

982 **Fusion-in-Decoder (FiD) Pipeline** (Izacard and
983 Grave, 2020) serves as a standard RAG-style base-
984 line. In this setting, the retriever is frozen, and only
985 the reader is fine-tuned while processing retrieved
986 image–caption pairs as contextual input.

987 **Perplexity Distillation Loss (PDist)** (Izacard
988 et al., 2023). For each query and candidate docu-
989 ment, we compute the reduction in perplexity (or
990 equivalently, the increase in likelihood) of the cor-
991 rect output when the language model is conditioned
992 on that document. From these likelihoods, we de-
993 rive a “posterior” over documents (proportional
994 to their contributions), and train the retriever to
995 mimic that posterior via a KL divergence objective.
996 In doing so, the retriever is guided to prioritize docu-
997 ments that most effectively support the language
998 model’s generation, thereby aligning retrieval with
999 generation quality.

1000 **C.2 Medical Pre - trained baseline models**

1001 We benchmark against state-of-the-art medical
1002 large vision–language models (LVLMs) that un-
1003 derwent large-scale medical pre-training:

1004 **BiomedGPT.** An open multimodal generative
1005 pre-trained transformer for biomedicine that aligns
1006 diverse biological/biomedical modalities with nat-
1007 ural language; we include it as a strong domain
1008 baseline (Luo et al., 2023).

1009 **LLaVA-Med.** A biomedical adaptation of
1010 LLaVA trained via curriculum (biomedical
1011 figure–caption alignment followed by instruction
1012 tuning); we evaluate three commonly used releases
1013 as separate baselines (Li et al., 2023).

1014 **MedVInT-TE and MedVInT-TD.** Two imple-
1015 mentations of Medical Visual Instruction Tuning
1016 from PMC-VQA (Zhang et al., 2023b). Both adopt
1017 a PMC-CLIP vision encoder (Lin et al., 2023). TE
1018 uses an encoder-style text pathway feeding a multi-
1019 modal decoder, while TD concatenates text tokens
1020 with visual features to a decoder-only pathway; we
1021 report both as distinct baselines.

Model	Breast		Derma		Retina		Vindr		BREST		VQARAD	SLAKE	PathVQA	Mean	
	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	ACC	ACC	ACC	F1
FT RAG _{Qwen2-vl}	.84	.81	.73	.38	.63	.45	.54	.09	.49	.27	.77	.89	.92	.72	.40
CLARE _{Qwen2-vl}	.87	.82	.75	.40	.65	.50	.58	.15	.50	.39	.79	.91	.92	.74	.45
FT RAG _{Pixtral}	.88	.84	.79	.57	.57	.46	.55	.09	.44	.31	.74	.87	.88	.72	.45
CLARE _{Pixtral}	.91	.87	.80	.59	.59	.49	.56	.15	.50	.39	.77	.89	.89	.74	.50

Table 10: Performance comparison of CLARE against FT RAG baseline across medical imaging classification and VQA tasks. Both Qwen2-VL and Pixtral backbones show consistent improvements with CLARE, particularly in F1 scores. Results demonstrate that BiomedCLIP-based retrieval enhances performance similarly to the general-purpose retriever. VQA results are reported for the closed-question subset using exact match accuracy.

Model	Breast		Derma		Retina		VQARAD
	ACC	F1	ACC	F1	ACC	F1	ACC
Reader only	.82	.77	.67	.54	.57	.40	.65
CLARE	.81	.77	.68	.54	.60	.51	.69

Table 11: Med-Flamingo comparison with and without CLARE retrieval augmentation across medical imaging tasks. CLARE achieves best performance on 4 out of 7 metrics, with notable improvements over Reader Only in Retina classification (ACC: 0.57→0.60, F1: 0.40→0.51) and VQA-RAD (0.65→0.69). However, the overall performance gains are more modest than those observed with other backbone architectures, indicating that Med-Flamingo may have limited capacity to leverage additional retrieval context.

InternVL-based baselines: MedDr, MedDr+RAD, and GSCo. MedDr is a generalist medical VLLM trained on large-scale instruction-style data curated from diagnosis-guided bootstrapping and medical image descriptions, covering multiple modalities and tasks (He et al., 2024a). MedDr+RAD augments MedDr at inference with Retrieval-Augmented Diagnosis (RAD): for a test image, similar cases are retrieved and summarized as contextual guidance in the prompt (He et al., 2024b). GSCo (Generalist-Specialist Collaboration) is a two-stage framework that (i) builds a generalist GFM (MedDr) and lightweight specialist models, and (ii) performs collaborative inference via Mixture-of-Expert Diagnosis (MoED; using specialist predictions as references) and RAD (using specialists to retrieve similar cases) (He et al., 2024b). GSCo evaluates across a large multi-dataset benchmark; we include GSCo as

a strong InternVL-based baseline along with its MedDr and MedDr+RAD components (He et al., 2024b; Chen et al., 2024).

D Additional Analysis

D.1 Retrieval Relevancy Analysis

We explore whether LVLM-aware multimodal retrieval fine-tuning improves the quality of retrieval candidates compared to those produced before fine-tuning. Importantly, our setup does not include ground-truth retrieval labels; instead, the retriever learns which candidates better steer the frozen LVLM toward the correct downstream prediction. To measure candidate improvement, we compare the retriever’s image and text heads before vs. after fine-tuning in a head-to-head ranking task. We then prompt GPT-5.2 to judge which candidate more closely resembles the ground-truth prediction, assuming that higher similarity to the ground-truth prediction indicates better evidence steering. When both candidates are identical, we record a tie. To avoid any location biases, we conduct shuffling between entering the model and reducing cost; we use a subset of 10% of the samples. Prompt 2 details the prompt for the task. Table 12 shows that LVLM-aware multimodal retrieval fine-tuning improves candidate relevance, with a mean gain of 0.17 for both Qwen2-VL and Pixtral, demonstrating the practical utility of our method.

D.2 Analysis of Model Components

We analyze the contribution of each part in the model in Table 7a and in Table 7b. The result of the first stage of training - the reader retrieval augmentation fine-tuning is the FT RAG then we check the contribution of training only the text retriever head and then finally is the contribution of training

Backbone	Winner	Breast	Derma	Retina	VinDr	BRSET	VQA-RAD	SLAKE	PathVQA	Mean
Qwen2-vl	Tie	.20	.59	.23	.25	.50	.24	.50	.20	.34
	FT RAG Wins	.30	.09	.32	.36	.24	.25	.16	.35	.26
	CLARE Wins	.50	.30	.44	.46	.25	.50	.34	.45	.41
Pixtral	Tie	.24	.51	.51	.47	.35	.12	.15	.29	.33
	FT RAG Wins	.26	.12	.15	.19	.40	.32	.37	.30	.26
	CLARE Wins	.50	.35	.34	.34	.52	.54	.47	.41	.43

Table 12: Analysis of relevance of retrieval candidates before and after LVLM-aware multimodal retrieval fine-tuning. Training the multimodal retrieval significantly improves performance, with a mean improvement of 0.17 across Qwen2-vl and Pixtral.

both the text and image retriever head as CLARE. We also explore the utility of converting the open question into a closed one in our retrieval loss. If we don’t apply o3 model to convert the open question, we only trained on the closed one, which is referred to in Table 7b as Closed Question Only.

D.3 Performance with a State-of-the-Art LVLM Reranker

We evaluate whether a state-of-the-art large vision–language model (LVLM), used as an optional reranker, can improve performance on inconsistent-retrieval predictions. Specifically, we test the o3 reasoning model (OpenAI, 2025) as a reranker: given a query image to classify and four retrieved image–caption pairs, o3 is asked to select the single pair that is most informative for predicting the correct label. We compare this reranking to four alternatives: (i) using the top-1 retrieved candidate (no reranking), (ii) selecting the model’s prediction with the highest confidence (maximum logit), (iii) using o3 with captions only (no images), and (iv) the CLARE aggregation strategy.

Our results show that o3 generally outperforms the maximum-logit baseline but remains inferior to CLARE’s logit-fusion aggregation. We evaluated performance on the Breast, Derma, and Retina datasets. For the remaining datasets—VinDr-PCXR, BRSET, and the visual question answering (VQA) datasets—we did not run this evaluation because the images and/or retrieved candidates come from restricted-access sources (e.g., MIMIC-CXR, VinDr-PCXR, and BRSET). See results in Table 9.

D.4 CLARE Variants with Medical Pre-Trained Backbones

Although our work focuses on general-purpose LVLMs, we also evaluated a medical-specific baseline. We conducted the evaluation on four benchmarks—Breast, Derma, Retina, and VQA-

RAD (closed-question subset). For the medical pre-trained baseline, we used Med-Flamingo and paired it with a medically pre-trained retriever (BiomedCLIP), using only one retrieval head (image retrieval without text retrieval). The performance gap between Med-Flamingo and our reader approach was small—0.01 in accuracy and 0.03 in F1. We hypothesize that this narrow gap may be due to Med-Flamingo’s relatively older LLaMA backbone (Touvron et al., 2023). Results are presented in Table 11. We also tested our general-purpose backbone with a medically pre-trained retriever, which showed improvements over FT RAG, with a similar pattern to the general-purpose retriever. Full results are presented in Table 10.

D.5 Retrieval Robustness Analyses.

We strive for our model to be robust in cases where the retrieval is noisy. Meaning, that the intrinsic ability in the model parameters is not overly dependent on the retrieval content. To test this property without ground truth retrieval labels, we evaluate the model performance with random candidates instead of searching according to cosine similarity. Overall as shown in Table 8 in Appendix D, we observed that for all datasets the model performance is above the reader-only, showing the model succeeded more than a model trained without any retrieval.

We also test the ability of the model to still function without any retrieved context. We do this to check whether some intrinsic capability is still maintained in the model’s reader without any retrieved information. This could also be helpful in cases where model users would like it to be versatile and support both retrieval-augmented prediction and non-retrieval-augmented prediction (e.g., in cases where no retrieved contexts were discovered). As shown in Table 8 in Appendix D, removing retrieval from CLARE resulted in performance

Backbone	Model	Breast		Derma		Retina		VinDr		BRSET	
		ACC	F1	ACC	F1	ACC	F1	ACC	F1	ACC	F1
Qwen2-vl	CLARE N=2	.87	.83	.74	.50	.63	.50	.56	.14	.47	.35
	CLARE N=4	.87	.84	.76	.50	.65	.50	.57	.14	.49	.37
	CLARE N=6	.87	.84	.75	.50	.63	.52	.58	.14	.48	.37
Pixtral	CLARE N=2	.89	.86	.80	.63	.58	.48	.54	.13	.49	.35
	CLARE N=4	.90	.87	.80	.62	.60	.51	.56	.14	.51	.37
	CLARE N=6	.90	.87	.80	.63	.61	.51	.55	.13	.49	.36

Table 13: Ablation on the number of retrieved candidates (N) for classification tasks.

Backbone	Model	VQA-RAD		SLAKE		PathVQA	
		Closed	Open	Closed	Open	Closed	Open
Qwen2-vl	CLARE N=2	.78	.46	.88	.82	.92	.37
	CLARE N=4	.79	.48	.90	.84	.93	.38
	CLARE N=6	.78	.47	.89	.83	.93	.38
Pixtral	CLARE N=2	.78	.45	.89	.83	.93	.38
	CLARE N=4	.78	.47	.90	.84	.93	.37
	CLARE N=6	.78	.47	.89	.84	.93	.38

Table 14: Ablation on the number of retrieved candidates (N) for visual question answering tasks.

comparable or slightly lower than the reader-only model.

Finally, we tested the model without the input image, relying only on the (image, report / caption) retrievals, to further assess the information held in retrieved data for predictions. In this scenario, we expected a substantial performance drop because the model no longer had direct knowledge of the patient’s condition, aside from what was provided by retrieved candidates. As a baseline, we used a simple random guess (coin flip), representing a scenario in which retrieval was not used effectively. As shown in Table 8 in Appendix D, our model’s performance exceeded this baseline by a large margin, with improvements of 0.32 in Accuracy, 0.28 in F1. These results indicate that the model learns from the retrieved data, yet it does not become overly dependent on it as we previously showed.

D.6 Different number of retrieved candidates

Our model demonstrates robustness across varying numbers of retrieved candidates. We explore performance for N = 2, 4, and 6, showing consistent results across the different configurations. Results are presented in Table 13 for classification tasks and in Table 14 for visual question answering tasks.

E Prompt Design

We detail the prompts used in our experiments and analyses. To convert open questions into closed ones with the o3 model, we use the prompt in Prompt 1. To assess the utility of o3 as a re-ranker of the most predictive image–caption pair (Section D.3, Appendix D), we provide both LVLM backbones (Pixtral and Qwen2-VL) with the task-specific prompt in Prompt 3.

We also include the training prompts for CLARE across datasets—Prompts 4, 6, 5, 8, 7—and the VQA prompt (Prompt 9).

```
Convert the following open-ended question into a closed yes/no question based on the given answer
.
The new question should be answerable with "{expected_answer}".

Original Question: {question}
Original Answer: {answer}
Expected New Answer: {expected_answer}

Please provide only the new closed question without any additional text or explanation.
```

Listing 1: Prompt used for o3 to convert open question into closed one

```
You are an expert medical AI assistant specializing in evaluating medical image and text
relevance for clinical decision-making.

Compare two medical (image, text) pairs and determine which is MORE RELEVANT to the ground truth
label.

Ground Truth Label: "{gt_text}"

First Image: {image_1}
First Text: {text_1}

Second Image: {image_2}
Second Text: {text_2}

Consider:
- Medical image findings
- Medical terminology accuracy in text
- Clinical relevance to the ground truth
- Diagnostic information alignment
- How useful each (image, text) pair would be for the given condition

Return JSON format:
{{
  "choice": <1 or 2>,
  "confidence": <"high", "medium", or "low">,
  "explanation": "<brief 1-sentence explanation>"
}}

Be concise and accurate. Choose the option that is MORE medically relevant to the ground truth.
```

Listing 2: Prompt used for retrieval candidate relevance evaluation

Task:

You are an expert assistant selecting the most informative reference for medical image classification.

Description:

A patient's [MRI/CT/X-ray] image needs to be classified into one of the following categories: [List of possible labels].

You are given four candidate reference pairs, each consisting of a caption and its associated image, drawn from PubMed literature.

Your goal is to determine which single candidate provides the most useful information to help an AI model correctly classify the provided image.

Input Information:

- Medical image: The patient's MRI/CT/X-ray image is provided here.
- Classification task: Classify the image into one of the following categories: [List of possible labels].
- Candidates:
 - Candidate 0:
 - Caption: [Caption 0]
 - Image: [Image 0]
 - Candidate 1:
 - Caption: [Caption 1]
 - Image: [Image 1]
 - Candidate 2:
 - Caption: [Caption 2]
 - Image: [Image 2]
 - Candidate 3:
 - Caption: [Caption 3]
 - Image: [Image 3]

Instructions:

1. For each candidate, carefully assess how informative its caption and image are for the current classification task.
2. Specifically, evaluate whether the candidate describes imaging features, findings, or clinical context relevant to distinguishing among the listed categories.
3. Compare all candidates and select the one that gives the clearest, most discriminative information to aid the classification.
4. Output only the number of the selected candidate in the following format:
Caption Number: [your selected number]

Listing 3: Prompt used for o3

```
<retrieved image>background:<retrieved
text>\n\n<query image>Does this breast
ultrasound image show signs of cancer?
```

Listing 4: Prompt used by JOMED for Breast classification

```
<retrieved image>background:<retrieved
text>\n<query image>what is the severity of
diabetic retinopath?
```

Listing 5: Prompt used by JOMED for Retina classification

```

<retrieved image>background:<retrieved text>\n\nProvide answer according to the labels:\n
"0 - 'Actinic keratoses and intraepithelial carcinoma'\n"
"1 - 'Basal cell carcinoma'\n"
"2 - 'Benign keratosis-like lesions'\n"
"3 - 'Dermatofibroma'\n"
"4 - 'Melanoma'\n"
"5 - 'Melanocytic nevi'\n"
"6 - 'Vascular lesions'\n\n"<query image>What type of skin lesion does the patient have?

```

Listing 6: Prompt used by JOMED for Derma classification

```

<retrieved image>background:<retrieved text> Provide answer according to the labels:\n
0 - 'No finding'\n
1 - 'Bronchitis'\n
2 - 'Brocho-pneumonia'\n
3 - 'Other disease'\n
4 - 'Bronchiolitis'\n
5 - 'Situs inversus'\n
6 - 'Pneumonia'\n
7 - 'Pleuro-pneumonia'\n
8 - 'Diagphramatic hernia'\n
9 - 'Tuberculosis'\n
10 - 'Congenital emphysema'\n
11 - 'CPAM'\n
12 - 'Hyaline membrane disease'\n
13 - 'Mediastinal tumor'\n
14 - 'Lung tumor'\n\n
\n\n<image>Question: Look at this X-ray scan and select all abnormalities you see from the
given labels Answer:

```

Listing 7: Prompt used by JOMED for Vindr-PCXR classification

```

<retrieved image>background:<retrieved text> Provide answer according to the labels:\n
0 - 'no findings'\n
1 - 'diabetic_retinopathy'\n
2 - 'macular_edema'\n
3 - 'scar'\n
4 - 'nevus'\n
5 - 'amd'\n
6 - 'vascular_occlusion'\n
7 - 'hypertensive_retinopathy'\n
8 - 'drusens'\n
9 - 'hemorrhage'\n
10 - 'retinal_detachment'\n
11 - 'myopic_fundus'\n
12 - 'increased_cup_disc'\n
13 - 'other'\n\n
\n\n<image>Question: Look at retinal fundus image and select all abnormalities you see from the given
labels

```

Listing 8: Prompt used by JOMED for BRSET classification

```
<retrieved image>background:<retrieved  
text>\n<query image><query question>?
```

Listing 9: Prompt used by JOMED for visual question answering tasks