

LLaVA-RadZ: Can Multimodal Large Language Models Effectively Tackle Zero-shot Radiology Recognition?

Anonymous ACL submission

Abstract

Recently, Multimodal Large Language Models (MLLMs) have demonstrated exceptional capabilities in visual understanding and reasoning across various vision-language tasks. However, we found that MLLMs cannot process effectively from fine-grained medical image data in the traditional Visual Question Answering (VQA) pipeline, as they do not exploit the captured features and available medical knowledge fully, results in MLLMs usually performing poorly in zero-shot medical disease recognition. Fortunately, this limitation does not indicate that MLLMs are fundamentally incapable of addressing fine-grained recognition tasks. From a feature representation perspective, MLLMs demonstrate considerable potential for tackling such challenging problems. Thus, to address this challenge, we propose *LLaVA-RadZ*, a simple yet effective framework for zero-shot medical disease recognition via utilizing the existing MLLM features. Specifically, we design an end-to-end training strategy, termed *Decoding-Side Feature Alignment Training (DFAT)* to take advantage of the characteristics of the MLLM decoder architecture and incorporate modality-specific tokens tailored for different modalities. Additionally, we introduce a *Domain Knowledge Anchoring Module (DKAM)* to exploit the intrinsic medical knowledge of large models, which mitigates the *category semantic gap* in image-text alignment. Extensive experiments demonstrate that our LLaVA-RadZ significantly outperforms traditional MLLMs in zero-shot disease recognition, achieving the comparable performance to the well-established and highly-optimized CLIP-based approaches.

1 Introduction

With the rapid advancement of deep learning technologies, an increasing number of studies have focused on their applications in medical disease diagnosis, yielding remarkable results (Chan et al.,

2020; Jamshidi et al., 2020; Lee et al., 2022; Tran et al., 2021). However, these approaches typically rely on high-quality annotations provided by clinical experts. Unlike natural image datasets, annotating medical images is both costly and time-consuming. To address this challenge, recent research has explored methods based on paired medical images and textual reports, leveraging contrastive learning techniques. By minimizing the distance between paired samples while maximizing the distance between unpaired ones, these CLIP-based approaches enable zero-shot disease recognition, thereby reducing reliance on extensive medical data annotation to a certain extent. In our in-depth investigation of advanced zero-shot disease recognition methods in the medical domain, several representative CLIP-based models (Lai et al., 2024; Wu et al., 2023; Zhang et al., 2023b; Phan et al., 2024) have achieved significant performance improvements leveraging the capabilities of Large Language Models or incorporate expert domain knowledge to some extent, rather than leveraging the models' intrinsic understanding capabilities.

Recently, Multimodal Large Language Models (MLLMs) (Achiam et al., 2023; Team et al., 2023; Liu et al., 2023; Huang et al., 2024, 2025; You et al., 2025) have demonstrated remarkable capabilities across various user-oriented vision-language tasks, such as image comprehension and reasoning, offering new possibilities for zero-shot disease recognition in medical applications. Among these, LLaVA-Med (Li et al., 2024a) has exhibited exceptional domain-specific medical knowledge in dialogue-based tasks, indicating that it possesses a certain degree of medical expertise. However, a recent study (Zhang et al., 2024) found that MLLMs, *i.e.*, LLaVA (Liu et al., 2023)) performed significantly worse than CLIP (Radford et al., 2021) on standard image classification tasks.

To further validate this observation, we conducted zero-shot classification experiments using

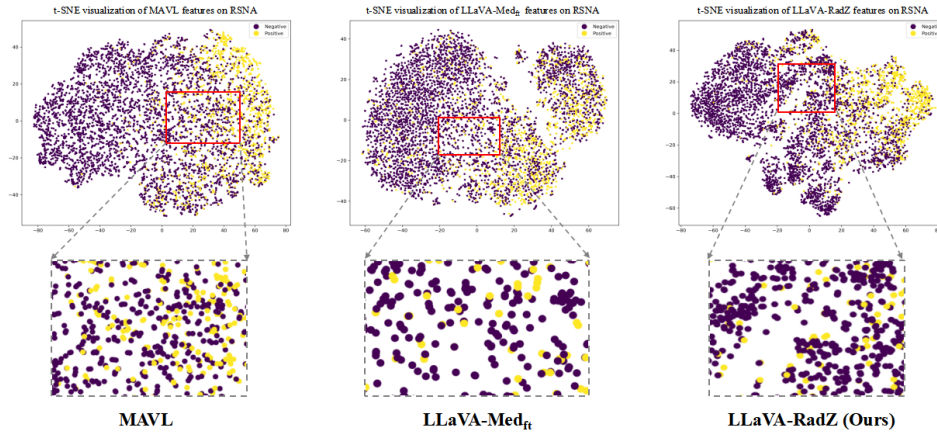


Figure 1: Comparison of Feature Distributions among MAVL, LLaVA-Med_{ft}, and LLaVA-RadZ on the RSNA Dataset.

multiple MLLMs on five medical imaging datasets (see Tab. 1). The experimental results are consistent with previous findings, confirming that MLLMs exhibit suboptimal performance in image classification, particularly when dealing with complex medical images. To enhance the generalization capability of MLLMs in radiology disease recognition tasks, we employed a fine-tuning strategy and performed supervised fine-tuning on the MIMIC-CXR dataset (Johnson et al., 2019). Additionally, inspired by the work of (Zhang et al., 2024), we incorporated a series of optimizations. While these improvements yielded performance gains, the results remained inferior compared to CLIP-based models. This phenomenon raises a critical question: *Can MLLMs effectively perform zero-shot disease recognition?*

As shown in Fig. 1, we visualize the feature distributions of MAVL (Phan et al., 2024), LLaVA-Med_{ft} (fine-tuned by the same dataset of our LLaVA-RadZ) and LLaVA-RadZ on the RSNA (Shih et al., 2019) dataset. The results indicate that MLLM exhibits strong feature extraction capabilities, comparable to the well-established MAVL in the domain. However, in the disease recognition task, MAVL significantly outperforms fine-tuned LLaVA-Med. We hypothesize that this performance gap arises because MLLMs fail to fully utilize the extracted features for effective disease identification via traditional VQA pipeline.

Inspired by this, we propose a simple yet effective LLaVA-RadZ framework for zero-shot disease recognition using the MLLM features. Our proposed framework has the fundamental difference compared with previous CLIP-base methods and traditional MLLM VQA pipeline. As

shown in Fig. 4 in the Appendix, we design a dedicated MLLM feature-based framework to address zero-shot medical disease recognition. Our proposed framework effectively leverages pre-trained MLLM representations to overcome the inherent limitations of the traditional VQA pipeline on this task. Specifically, firstly, we introduce a new training strategy, Decoding-Side Feature Alignment Training (DFAT). Specifically, we introduce special tokens for both image and text modalities and leverage the autoregressive generation capability of the decoder architecture to extract global representations of images and texts. Additionally, we incorporate a cross-modal contrastive loss to optimize the model’s ability to learn discriminative features. Furthermore, to mitigate the semantic category gap encountered during fine-grained alignment between medical images and textual reports, we design a Domain Knowledge Anchoring Module (DKAM). DKAM utilizes the model’s intrinsic medical knowledge to extract the semantic information underlying disease categories, constructing disease description vectors that serve as an intermediary bridge to facilitate the alignment between medical images and textual reports, thereby establishing a stable relationship. To further enhance the correlation among medical images, textual reports, and disease categories, a category knowledge-guided loss strengthens the association between similar images and corresponding textual reports.

Our main contributions can be summarized as follows.

- We analyze the limitations of current MLLMs in addressing complex fine-grained medical disease recognition tasks, investigate the underlying causes of these constraints, and pro-

pose a novel end-to-end feature-based MLLM framework to mitigate these challenges. To the best of our knowledge, we are the *first* work in the field of medical disease recognition to explore how to use MLLM features directly to solve recognition problems.

- We propose the tailored training strategy DFAT, and incorporate a cross-modal contrastive loss to optimize the model’s ability to achieve effective alignment between visual and textual features. Furthermore, we design a DKAM to leverage MLLM’s intrinsic medical knowledge and effectively mitigate semantic gap in image-text alignment, thereby enhancing category-level alignment.
- We conduct extensive experiments on multiple large-scale radiology diagnosis datasets, validating the potential of LLaVA-RadZ in zero-shot disease recognition tasks.

2 Approach

2.1 Can Med-LLMs Be Good Medical Classifiers?

Previous studies have explored the classification capabilities of multimodal large language models (MLLMs), revealing that their performance on image classification tasks is often limited. For example, (Zhang et al., 2024) investigates the performance differences in classification between MLLMs and CLIP, focusing on factors such as inference strategies, training approaches, and datasets. Inspired by this work, we extend the exploration to zero-shot tasks in the medical domain. Unlike natural images and text, the relationship between medical images and reports is more complex. We seek to investigate whether large medical models, leveraging domain-specific knowledge, can achieve superior performance on medical zero-shot tasks.

We first evaluated two open-source MLLMs, i.e., LLaVA-1.5 (Liu et al., 2023) and LLaVA-Med (Li et al., 2024a), on five medical datasets in a zero-shot classification setting. The evaluation followed a general large-model classification approach, where the model selects the correct category from a set of candidate options. As shown in Tab. 1, these models demonstrated limited performance in disease classification tasks and failed to accurately identify various medical conditions. Given the potential knowledge limitations of these

models, we further assessed the performance of more powerful proprietary MLLMs (i.e., Qwen2.5-Max (Yang et al., 2024), Gemini-Pro (Team et al., 2023), and GPT-4o (Achiam et al., 2023)) on zero-shot medical disease recognition tasks. As shown in table 1, these models exhibited superior classification capabilities. However, they still lagged behind the state-of-the-art domain-specific methods in medical classification.

To enhance the generalization ability of MLLMs in radiology disease identification, we conducted Supervised Fine-Tuning (SFT) on LLaVA-1.5 (Liu et al., 2023) and LLaVA-Med (Li et al., 2024a) using the publicly available MIMIC-CXR dataset (Johnson et al., 2019). Surprisingly, the fine-tuned models did not achieve consistent performance improvements across the five datasets. In some cases, their classification performance even deteriorated. Further analysis of the model outputs revealed that MLLMs did not always focus on disease-specific information in radiology reports. Instead, they tended to overlearn the textual structures and linguistic patterns of the reports, which limited their classification capability. To mitigate this issue, we incorporated the Chain-of-Thought (CoT) prompting strategy and adjusted the model’s reasoning approach, inspired by the methodology of (Zhang et al., 2024), to optimize the model’s decision-making process. This approach led to moderate improvements in classification performance on medical datasets. Although the models have not yet reached optimal performance, the results suggest that MLLMs still hold significant potential for zero-shot medical disease recognition.

2.2 Motivation

As previously discussed, despite possessing a certain level of domain knowledge, medical MLLMs have not yet demonstrated remarkable performance in zero-shot medical tasks. Even with further instruction tuning, their performance remains inferior to that of existing vision-language models (VLMs). However, it is noteworthy that modifying the inference strategy leads to significant performance improvements, suggesting that MLLMs are indeed capable of capturing medical image and text features. Nevertheless, these features have yet to be fully exploited.

To address this limitation, we propose the LLaVA-RadZ framework, introducing a novel end-to-end training strategy, Decoding-Side Feature Alignment Training (DFAT). This approach lever-

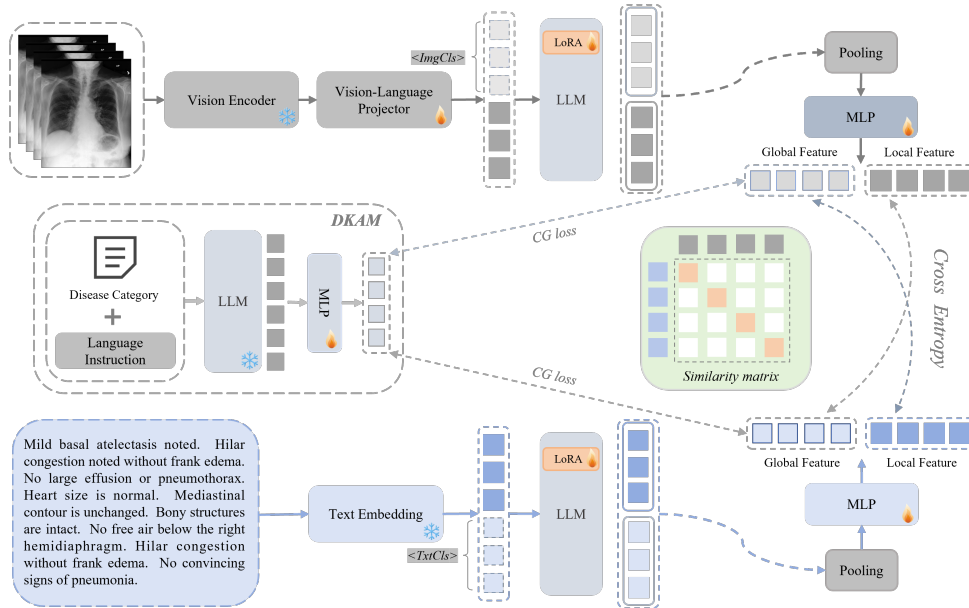


Figure 2: The LLaVA-RadZ framework consists of three components. (A) Construct a category semantic vector repository using the domain knowledge anchoring module (DKAM). (B) Encode medical images and text, appending $\langle \text{ImgCls} \rangle$ and $\langle \text{TxtCls} \rangle$ tokens before feeding them into the LLM. (C) Extract global and local features, optimizing with cross-entropy loss, while leveraging the semantic repository for category-level alignment.

ages the unique properties of the MLLM decoder architecture while incorporating modality-specific special tokens to facilitate effective interaction between medical images and textual features, ultimately achieving more robust cross-modal alignment. As illustrated in Fig. 1, we compare the feature distribution of our model with MAVL (Phan et al., 2024), the current state-of-the-art method, on the RSNA (Shih et al., 2019) dataset. The results clearly demonstrate that our model achieves better clustering of intra-class samples while enhancing inter-class separation, validating the effectiveness of our approach. Furthermore, we introduce the Domain Knowledge Anchor Module (DKAM), which harnesses the intrinsic medical knowledge of LLMs to bridge the semantic gap between images and text, enabling more precise disease classification.

2.3 The Proposed LLaVA-RadZ

We aim to learn generalizable medical image representations from radiology reports to enhance various downstream medical image recognition tasks, particularly when labeled data is scarce. The overall framework is illustrated in Fig. 2. Given a pair of medical images and reports, the image and text are first passed through separate visual and text encoders to obtain their respective encoded features. These encoded features and specially designed tokens are then fed into a language model to obtain the final feature representation. The features are

mapped into a common representational space via an MLP projection layer and optimized with the InfoNCE loss. Furthermore, we propose a Domain Knowledge Anchor Module (DKAM), which leverages domain knowledge inherent in the model to guide the alignment of text and image features at the category level.

2.3.1 End-to-End Training Strategy

Currently, most MLLMs employ generation-based training objectives for instruction fine-tuning. Although this approach effectively captures the features of medical images and textual reports, its performance in zero-shot tasks remains suboptimal, as it fails to fully leverage these features. To address this issue, we propose a novel training strategy, Decoding-Side Feature Alignment Training (DFAT), as illustrated in Fig. 2.

We consider a training dataset consisting of N pairs of medical image-text samples, denoted as $S_{\text{train}} = \{(X_1, Y_1), \dots, (X_N, Y_N)\}$. The medical image $X_i \in \mathbb{R}^{H \times W \times 3}$, with H and W representing the height and width of the image, respectively. Y_i refers to the corresponding medical text report associated with the image.

Specifically, we design special tokens for both image and text modalities, where $\langle \text{ImgCls}_i \rangle (i = 0, \dots, 3)$ denotes image feature tokens and $\langle \text{TxtCls}_i \rangle (i = 0, \dots, 7)$ denotes text feature tokens. These special tokens are attached to the image prompt X_{prompt} and the text

prompt Y_{prompt} , respectively. The image prompt X_{prompt} has a format similar to ‘‘What disease is indicated by the chest X-ray?’’, while the text prompt Y_{prompt} follows a format such as ‘‘What disease is described in this text?’’. By appending special tokens, we obtain the modified prompts $\tilde{X}_{\text{prompt}}$ and $\tilde{Y}_{\text{prompt}}$, which is represented as:

$$\tilde{X}_{\text{prompt}} = [X_{\text{prompt}}; \{\langle \text{ImgCls}_i \rangle\}_{i=0}^3], \quad (1)$$

$$\tilde{Y}_{\text{prompt}} = [Y_{\text{prompt}}; \{\langle \text{TxtCls}_i \rangle\}_{i=0}^7]. \quad (2)$$

When an image and its corresponding prompt $\tilde{X}_{\text{prompt}}$ are input into the MLLM \mathcal{F} to generate a response \hat{R}_{img} . Similarly, when a text sample and its corresponding feature extraction prompt $\tilde{Y}_{\text{prompt}}$ are provided as input, the model produces a response \hat{R}_{txt} . This process can be formally expressed as:

$$\hat{R}_{\text{img}}^i = \mathcal{F}(X_i, \tilde{X}_{\text{prompt}}), \quad \hat{R}_{\text{txt}}^i = \mathcal{F}(Y_i, \tilde{Y}_{\text{prompt}}). \quad (3)$$

Due to the autoregressive nature of the decoder architecture, when the LLM processes visual and textual information to generate responses, its internal representations are stored in the designated special tokens. Specifically, we extract the penultimate layer embedding \tilde{h}_{img} corresponding to the special token $\langle \text{ImgCls}_i \rangle$, which stores the global image features $H_{\text{img}}^{\text{global}} \in \mathbb{R}^{B \times I \times K}$. Here, B denotes the number of image-text pairs in each batch, I represents the number of special image tokens, and K is the dimension of the shared embedding space. After applying a pooling operation followed by an MLP projection layer γ_{img} , we obtain the global image feature representation $X_g \in \mathbb{R}^{B \times K}$. The local image feature $X_l \in \mathbb{R}^{B \times K}$ is obtained by pooling the hidden states of all tokens except those corresponding to special tokens, followed by an MLP projection layer γ_{img} :

$$\begin{aligned} X_g &= \gamma_{\text{img}}(\text{AvgPool}(H_{\text{img}}^{\text{global}})), \\ X_l &= \gamma_{\text{img}}(\text{AvgPool}(H_{\text{img}}^{\text{local}})). \end{aligned} \quad (4)$$

Similarly, we extract the global text representation $Y_g \in \mathbb{R}^{B \times K}$ and the local text representation $Y_l \in \mathbb{R}^{B \times K}$ using the same methodology:

$$\begin{aligned} Y_g &= \gamma_{\text{txt}}(\text{AvgPool}(H_{\text{txt}}^{\text{global}})), \\ Y_l &= \gamma_{\text{txt}}(\text{AvgPool}(H_{\text{txt}}^{\text{local}})). \end{aligned} \quad (5)$$

To further enhance fine-grained alignment across different modalities, we introduce a cross-modal

contrastive loss, L_{CA} . Specifically, for the i -th image-text pair (X_i, Y_i) in a batch, we alternately align the global and local features of images and texts. This procedure yields two symmetric, temperature-normalized InfoNCE objectives: one aligns global image features with local text features, and the other aligns local image features with global text features. These objectives maximize the mutual information between image-text pairs in the latent space.

For the alignment between global image features and local text features, we calculate two similarity matrices, $S_i^{X_g \rightarrow Y_l}$ and $S_i^{Y_l \rightarrow X_g}$, with the following computation:

$$S_i^{X_g \rightarrow Y_l} = \frac{X_{g,i} \cdot Y_{l,i}^T}{\tau}, \quad S_i^{Y_l \rightarrow X_g} = \frac{Y_{l,i} \cdot X_{g,i}^T}{\tau}. \quad (6)$$

where τ is the temperature hyperparameter. Subsequently, we compute the contrastive loss between the global image and the local text, with the following formula:

$$\begin{aligned} L_{CA}^{X_g \rightarrow Y_l, i} &= -\log \frac{\exp(S_i^{X_g \rightarrow Y_l})}{\sum_{k=1}^B \exp(S_k^{X_g \rightarrow Y_l})}, \\ L_{CA}^{Y_l \rightarrow X_g, i} &= -\log \frac{\exp(S_i^{Y_l \rightarrow X_g})}{\sum_{k=1}^B \exp(S_k^{Y_l \rightarrow X_g})}, \\ L_{CA}^{X_g \rightarrow Y_l} &= \frac{1}{2} \sum_{i=1}^B (L_{CA}^{X_g \rightarrow Y_l, i} + L_{CA}^{Y_l \rightarrow X_g, i}). \end{aligned} \quad (7)$$

Similarly, for the alignment between local image features and global text features, we compute the contrastive loss between the local image and global text.

$$\begin{aligned} L_{CA}^{X_l \rightarrow Y_g} &= -\frac{1}{2} \sum_{i=1}^B \left(\log \frac{\exp(S_i^{X_l \rightarrow Y_g})}{\sum_{k=1}^B \exp(S_k^{X_l \rightarrow Y_g})} \right. \\ &\quad \left. + \log \frac{\exp(S_i^{Y_g \rightarrow X_l})}{\sum_{k=1}^B \exp(S_k^{Y_g \rightarrow X_l})} \right). \end{aligned} \quad (8)$$

Finally, we obtain our cross-modal contrastive loss L_{CA} .

$$L_{CA} = \frac{1}{2} (L_{CA}^{X_g \rightarrow Y_l} + L_{CA}^{X_l \rightarrow Y_g}). \quad (9)$$

2.3.2 Domain Knowledge Anchor Module

In aligning medical images with text reports, we observed that the critical entity of the medical disease categories was merely encoded as features by the model, without considering the underlying semantics. To address this limitation and further enhance fine-grained alignment capabilities, we introduce the Domain Knowledge Anchoring Module (DKAM). Initially, we leverage the inherent

Table 1: Comparison of zero-shot disease classification performance of public MLLMs and LLaVA-based methods on five medical benchmarks. “ft” indicates supervised LoRA fine-tuning, “CoT” denotes zero-shot chain-of-thought prompting, and “Inference” refers to CLIP inference. Best results are in bold; second-best are underlined.

Dataset		CheXpert			ChestXray-14			COVIDx CXR-2			RSNA Pneumonia			SIIM-ACR		
Method	Model	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑
MLLM	LLaVA-1.5 (7B) (Liu et al., 2023)	-	7.50	8.28	-	3.33	6.92	-	53.14	50.28	-	40.53	55.34	-	23.66	50.36
	LLaVA-Med (7B) (Li et al., 2024a)	-	6.87	8.94	-	8.02	6.78	-	34.90	50.03	-	18.58	50.00	-	21.91	49.90
	Qwen2.5-Max (Yang et al., 2024)	-	32.23	67.97	-	19.04	76.19	-	<u>75.91</u>	<u>76.81</u>	-	43.58	43.59	-	64.70	<u>72.57</u>
	Gemini-Pro (Team et al., 2023)	-	35.01	76.08	-	14.16	77.78	-	62.84	62.90	-	44.23	51.43	-	61.43	72.03
	GPT-4o (Achiam et al., 2023)	-	<u>45.85</u>	<u>81.14</u>	-	19.85	<u>81.55</u>	-	50.93	77.08	-	54.20	65.33	-	64.57	72.11
Explorative Methods	LLaVA-1.5-7Bft	-	10.61	19.62	-	7.85	19.06	-	27.74	25.18	-	43.60	34.80	-	52.37	50.95
	LLaVA-Med-7Bft	-	14.25	31.46	-	9.00	21.43	-	27.42	24.09	-	46.72	38.88	-	53.11	57.68
	LLaVA-Med-7Bft + CoT (Zhang et al., 2024)	-	8.90	26.23	-	8.33	20.46	-	27.12	26.55	-	49.59	43.80	-	54.06	51.07
	LLaVA-Med-7Bft + Inference (Zhang et al., 2024)	71.00	44.85	75.45	64.30	<u>21.73</u>	70.86	71.07	69.84	60.39	77.51	<u>69.85</u>	<u>72.90</u>	71.25	<u>68.26</u>	71.27
Ours	LLaVA-RadZn	73.36	48.59	82.15	72.61	27.91	84.64	84.36	77.53	74.58	86.98	76.18	83.28	89.92	79.57	84.38

Table 2: Comparison with SOTA methods on four medical datasets for zero-shot classification. AUC, F1, and ACC are reported. Best and second-best results are bold and underlined, respectively. Gray-background methods add supervised classification heads to contrastive learning, while others mainly rely on contrastive representations for zero-shot inference. ProbMED[†] uses official pretrained weights and an evaluation protocol similar to MAVL.

Method	ChestXray-14			COVIDx CXR-2			RSNA Pneumonia			SIIM-ACR		
	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑	AUC ↑	F1 ↑	ACC ↑
ConVIRT (Zhang et al., 2022)	53.15	12.38	57.88	62.78	71.23	63.84	79.21	55.67	75.08	64.25	42.87	53.42
GLoRIA (Huang et al., 2021)	55.92	14.20	59.47	64.52	70.78	60.21	70.37	48.19	70.54	54.71	40.39	47.15
BioViL (Boecking et al., 2022)	57.82	15.64	61.33	61.40	70.92	58.20	84.12	54.59	74.43	70.28	46.45	68.22
CheXzero (Tiu et al., 2022)	66.99	21.99	65.38	73.13	76.13	71.45	85.13	61.49	78.34	84.60	65.97	77.34
MedKLIP (Wu et al., 2023)	72.33	24.18	79.40	76.28	76.54	71.96	86.57	63.28	79.97	89.79	72.73	83.99
MAVL (Phan et al., 2024)	73.50	<u>26.25</u>	<u>82.77</u>	<u>83.86</u>	81.73	78.07	<u>86.91</u>	<u>63.41</u>	<u>82.42</u>	92.04	<u>77.95</u>	87.14
ALTA (Lian et al., 2025)	67.41	24.57	66.65	77.41	77.50	73.69	86.75	62.15	77.97	85.77	67.43	75.36
ProbMED [†] (Gao et al., 2025)	63.35	22.78	66.24	74.67	75.05	67.90	77.22	55.09	69.05	81.08	63.37	76.95
Ours	<u>72.61</u>	27.91	84.64	84.36	<u>77.53</u>	<u>74.58</u>	86.98	76.18	83.28	<u>89.92</u>	79.57	<u>84.38</u>

396 medical domain expertise of an LLM to generate
397 descriptive explanations for each disease category.
398 These generated disease descriptions serve as an in-
399 termediary bridge to guide the alignment between
400 medical images and text reports. Specifically, we
401 input the disease list D_{list} from the training dataset
402 along with a designed prompt template K_{prompt}
403 into the LLM \mathcal{F} . This process is formally ex-
404 pressed as:

$$405 \hat{R}_{\text{dis}} = \mathcal{F}(D_{\text{list}}, K_{\text{prompt}}). \quad (10)$$

406 By fully harnessing the LLM’s exceptional se-
407 mantic understanding, we prompt the model to ex-
408 plore the underlying semantics of the disease cat-
409 egories and discern their distinctions, ultimately
410 producing a refined disease description. The fea-
411 tures extracted from the LLM’s response are then
412 mapped via a multi-layer perceptron (MLP) to yield
413 the disease description vector \hat{D} , which is repre-
414 sented as:

$$415 \hat{D} = \gamma_{\text{dis}} \left(\hat{R}_{\text{dis}} \right). \quad (11)$$

416 Subsequently, we introduce the Category of
417 Knowledge-guided Contrastive Loss L_{CG} . Specifi-
418 cally, we calculate the cross-entropy loss between
419 the disease description vector \hat{D} and the global fea-
420 tures of both the images X_g and the text Y_g . This
421 design encourages the model to better capture the

semantic relationships among images, text, and dis-
422 ease categories during training, achieving a more
423 robust category-level alignment.
424

$$425 S_i^{\text{img-disease}} = \frac{X_{g,i} \cdot D^T}{\tau}, \quad S_i^{\text{txt-disease}} = \frac{Y_{g,i} \cdot D^T}{\tau}. \quad (12)$$

$$426 L_{CG,i}^{\text{txt}} = -\log \frac{\exp(S_i^{\text{txt-disease}})}{\sum_{k=1}^N \exp(S_k^{\text{txt-disease}})}, \quad (13)$$

$$427 L_{CG,i}^{\text{img}} = -\log \frac{\exp(S_i^{\text{img-disease}})}{\sum_{k=1}^N \exp(S_k^{\text{img-disease}})}.$$

428 Here, N represents the number of disease cate-
429 gories, B denotes the number of medical image-
430 text pairs in each batch, and τ is the temperature
431 hyperparameter. The final category of knowledge-
432 guided loss is as follows:

$$433 L_{CG} = \frac{1}{2} \sum_{i=1}^B \left(L_{CG,i}^{\text{txt}} + L_{CG,i}^{\text{img}} \right). \quad (14)$$

434 By combining the category knowledge-guided
435 loss and the cross-modal contrastive loss, the final
436 objective function is defined as follows:

$$437 L_{\text{total}} = \lambda L_{CA} + (1 - \lambda) L_{CG}, \quad (15)$$

438 where λ is a balancing factor used to adjust the
439 weights of the two losses, and it is set to 0.5 by
440 default.

Table 3: Comparison of performance with other SOTA methods under different data portions for the fine-tuning classification task. AUC scores are reported. The best and second-best results are highlighted in bold and underlined, respectively. Gray-background methods add supervised classification heads to contrastive learning, while others rely primarily on contrastive representations for inference.

Method	RSNA Pneumonia			Pneumothorax			COVIDx CXR-2		
	1%	10%	100%	1%	10%	100%	1%	10%	100%
Scratch	68.94	83.31	87.12	53.11	76.18	87.48	85.11	93.65	98.86
ConVIRT (Zhang et al., 2022)	78.86	85.42	87.64	72.39	80.41	91.67	90.30	97.74	99.70
GLoRIA (Huang et al., 2021)	79.13	85.59	87.83	75.85	86.20	91.89	92.74	97.18	99.54
BioViL (Boecking et al., 2022)	80.27	86.04	88.29	70.29	79.45	88.05	92.39	98.39	99.68
MedKLIP (Wu et al., 2023)	82.11	87.14	88.58	85.24	89.91	93.02	95.58	98.77	99.77
MAVL (Phan et al., 2024)	86.09	<u>87.90</u>	<u>88.94</u>	91.53	93.00	<u>94.48</u>	97.18	99.15	99.90
ALTA (Lian et al., 2025)	<u>86.61</u>	87.35	87.87	86.00	89.87	93.28	<u>97.33</u>	<u>99.44</u>	<u>99.95</u>
Ours	88.23	88.57	89.49	<u>88.42</u>	<u>89.96</u>	94.50	98.32	99.80	99.96

3 Experiments

In this section, we introduce the datasets used for pre-training and downstream tasks, followed by the implementation details and baseline methods for comparison.

3.1 Dataset

In our experiments, we pre-trained the model using the MIMIC-CXR dataset (Johnson et al., 2019). For downstream tasks, we primarily evaluated the model’s performance in medical disease classification using multiple benchmark datasets, including ChestX-ray14 (Wang et al., 2017), RSNA Pneumonia (Shih et al., 2019), SIIM-ACR Pneumothorax (sui, 2019), CheXpert (Irvin et al., 2019), and COVIDx CXR-2 (Pavlova et al., 2022). Detailed information on these datasets can be found in the supplementary material.

3.2 Evaluation Metrics

For the zero-shot classification task, we employ standard classification evaluation metrics, including Accuracy, AUC score, and F1 score. The macro-average metrics are reported for all diseases present in the target dataset.

3.3 Zero-shot evaluation

As shown in Tab. 2, we compare the performance of established methods in the field on the zero-shot classification task for radiological diseases, evaluated on four officially released test datasets. Our findings demonstrate that, compared to conventional CLIP-style models such as ConVIRT (Zhang et al., 2022), GLoRIA (Huang et al., 2021), BioViL (Boecking et al., 2022), and CheXzero (Tiu et al., 2022), our approach exhibits significant advantages. Even when com-

pared to state-of-the-art models incorporating external models or domain-specific expert knowledge, our method remains highly competitive. Specifically, on the multi-class dataset ChestXray-14, our model surpasses the supervised learning method MAVL (Phan et al., 2024) by 1.87% in accuracy. Moreover, on the RSNA Pneumonia dataset, we achieve a 12.77% improvement in F1 score. These results indicate that multimodal large language models (MLLMs) possess strong feature extraction capabilities, further underscoring their immense potential in medical disease classification tasks.

3.4 Fine-tuning evaluation

Consistent with previous studies (Phan et al., 2024; Wu et al., 2023), we fine-tune the model on downstream datasets using 1%, 10%, and 100% of the available data and further evaluate its performance. Tab. 3 presents the fine-tuning results across three datasets, demonstrating that our model consistently maintains a competitive advantage. Notably, when fine-tuned with only 1% data, our proposed LLaVA-RadZ outperforms the MAVL (Phan et al., 2024) model by 2.14% on the RSNA Pneumonia and by 1.14% on COVIDx. Even when fine-tuned with 100% data, our model continues to deliver performance improvements. This enhancement is likely attributed to our decoder-side alignment training strategy, which effectively captures global modality information and leverages the interaction between global and local features to achieve fine-grained cross-modal alignment, further strengthening the model’s disease recognition capability.

3.5 Ablation Study

Ablation Study of DKAM. To validate the effectiveness of our proposed Domain Knowledge An-

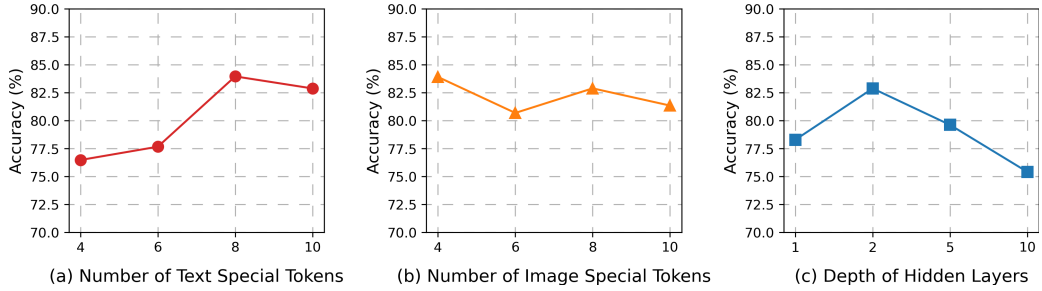


Figure 3: Effect of Special Token Numbers and Hidden Layer Depth on ChestXray-14 Classification.

Table 4: Ablation study of DKAM on ChestXray-14. D_1 represents a semantic vector set of 75 medical entities, and D_2 represents a semantic vector set of 14 disease categories.

#	DKAM	D_1	D_2	AUC \uparrow	F1 \uparrow	ACC \uparrow
a				69.31	27.30	82.32
b	✓	✓		68.67	25.73	81.84
c	✓		✓	72.61	27.91	84.64

Table 5: Ablation study of feature representations on ChestXray-14.

#	Global	Local	Prompt	AUC \uparrow	F1 \uparrow	ACC \uparrow
a		✓		67.14	25.11	77.82
b		✓	✓	68.29	26.42	78.63
c	✓		✓	70.13	26.22	82.50
d	✓	✓	✓	72.61	27.91	84.64

510 chor Module (DKAM), we conduct an ablation
 511 study on the ChestXray-14 dataset. With DKAM
 512 incorporated, we further examine the impact of dif-
 513 ferent category semantic vector repositories on the
 514 model’s fine-grained alignment capability. Follow-
 515 ing the MedKLIP study, we select 75 primary med-
 516 ical entities from the MIMIC-CXR dataset. How-
 517 ever, unlike MedKLIP, we leverage the model’s
 518 intrinsic domain knowledge to construct a category
 519 semantic vector repository, denoted as D_1 . In ad-
 520 dition, we build a disease-specific semantic vector
 521 repository for the 14 medical disease categories in
 522 the MIMIC-CXR training dataset, denoted as D_2 .

523 As shown in Tab. 4 (a vs. c), introducing DKAM
 524 significantly improves model performance. Using
 525 disease category semantics as an intermediary en-
 526 ables more accurate alignment between medical im-
 527 ages and textual descriptions at the category level.
 528 Further comparisons in Tab. 4 (b vs. c) indicate
 529 that, compared to a larger set of medical entities,
 530 a repository focused on primary disease categories
 531 provides stronger guidance for image-text align-
 532 ment. Moreover, additional medical entities in D_1 ,
 533 such as tip, tube, PICC, and device, may introduce
 534 noise and adversely affect disease-level alignment,
 535 as evidenced by the results in Tab. 4 (a vs. b).

536 **Ablation Study of Special Tokens.** As shown
 537 in Tab. 1, we have demonstrated the effectiveness
 538 of the Decoding-Side Feature Alignment Training
 539 (DFAT) strategy. To further investigate the design
 540 of the critical special tokens integral to this ap-
 541 proach, we conducted an in-depth analysis on the
 542 ChestXray-14 dataset. As illustrated in Fig. 3, we
 543 observed that the number of text and image to-
 544 kens significantly influences model performance,

545 with both an excessive and an insufficient count
 546 potentially resulting in a loss of modal informa-
 547 tion. Moreover, our study indicates that the optimal
 548 global features are not stored in the final hidden
 549 layer but rather in the penultimate layer, which may
 550 be attributed to the loss of fine-grained information
 551 due to deeper feature aggregation, thereby affecting
 552 overall performance.

553 **Ablation Study of Features.** During cross-modal
 554 alignment, we analyzed the impact of global and
 555 local features on model performance and exam-
 556 ined the effectiveness of prompts (Tab. 5). Results
 557 show that using only local features yields the worst
 558 performance, while global features perform bet-
 559 ter, likely because modality-specific tokens capture
 560 global information more accurately. Combining
 561 global and local features achieves the best results,
 562 and incorporating prompts further improves the
 563 model’s feature representation.

564 4 Conclusion

565 This paper proposes LLaVA-RadZ, a simple yet
 566 effective framework for zero-shot medical disease
 567 recognition. We introduce an end-to-end decoding-
 568 side feature alignment strategy to leverage MLLM
 569 characteristics and store modality-specific infor-
 570 mation. Cross-modal contrastive learning fur-
 571 ther enhances feature alignment and cross-modal
 572 understanding. Additionally, a Domain Knowl-
 573 edge Anchoring Module facilitates category-level
 574 alignment between medical images and textual de-
 575 scriptions. Experiments show that LLaVA-RadZ
 576 achieves strong performance across multiple bench-
 577 marks, demonstrating the potential of MLLMs for
 578 zero-shot radiological disease recognition.

5 Limitations

First, our experiments primarily focus on chest disease recognition tasks and have not yet been systematically validated on other types of radiological disease datasets. Although the proposed framework demonstrates promising effectiveness in chest image classification, its generalizability across different anatomical regions and disease categories warrants further exploration and validation in future research.

Second, the current study is predominantly centered on classification tasks. While this provides valuable insights into the model’s foundational performance in disease identification, future work could extend this framework to a broader range of clinical applications, such as lesion detection, segmentation, and medical report generation, to more comprehensively assess its clinical utility.

Third, constrained by computational resources, this study has not yet expanded the framework to larger-scale vision-language models. Thus, investigating the potential of more extensive models within this task represents a meaningful direction for future research.

References

2019. Society for imaging informatics in medicine: Siim-acr pneumothorax segmentation. <https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation>.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Olivier Bodenreider. 2004. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270.

Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, and 1 others. 2022. Making the most of text semantics to improve biomedical vision–language processing. In *European conference on computer vision*, pages 1–21. Springer.

Heang-Ping Chan, Lubomir M Hadjiiski, and Ravi K Samala. 2020. Computer-aided diagnosis in the era of deep learning. *Medical physics*, 47(5):e218–e227.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, and 1 others. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Shivang Desai, Ahmad Baghal, Thidathip Wongsurawat, Piroon Jenjaroenpun, Thomas Powell, Shaymaa Al-Shukri, Kim Gates, Phillip Farmer, Michael Rutherford, Geri Blake, and 1 others. 2020. Chest imaging representing a covid-19 positive rural us population. *Scientific data*, 7(1):414.

Yuan Gao, Sangwook Kim, Jianzhong You, and Chris McIntosh. 2025. Probmed: A probabilistic framework for medical multimodal binding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20157–20167.

Tianyu Han, Lisa C Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno K Bresssem. 2023. Medalpaca—an open-source collection of medical conversational ai models and training data. *arXiv preprint arXiv:2304.08247*.

Hulingxiao He, Geng Li, Zijun Geng, Jinglin Xu, and Yuxin Peng. 2025. Analyzing and boosting the power of fine-grained visual recognition for multi-modal large language models. *arXiv preprint arXiv:2501.15140*.

Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena Yeung. 2021. Gloria: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3942–3951.

Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu, and Shaohui Lin. 2025. Vision-r1: Incentivizing reasoning capability in multimodal large language models. *arXiv preprint arXiv:2503.06749*.

Wenxuan Huang, Zijie Zhai, Yunhang Shen, Shaoshen Cao, Fei Zhao, Xiangfeng Xu, Zheyu Ye, and Shaohui Lin. 2024. Dynamic-llava: Efficient multimodal large language models via dynamic vision-language context sparsification. *arXiv preprint arXiv:2412.00876*.

Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, and 1 others. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597.

Mohammad Jamshidi, Ali Lalbakhsh, Jakub Talla, Zdeněk Peroutka, Farimah Hadjilooei, Pedram Lalbakhsh, Morteza Jamshidi, Luigi La Spada, Mirhamed Mirmozafari, Mojgan Dehghani, and 1 others. 2020. Artificial intelligence and covid-19: deep learning approaches for diagnosis and treatment. *Ieee Access*, 8:109581–109595.

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

Haoran Lai, Qingsong Yao, Zihang Jiang, Rongsheng Wang, Zhiyang He, Xiaodong Tao, and S Kevin Zhou. 2024. Carzero: Cross-attention alignment for radiology zero-shot classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11137–11146.

690	Junghwan Lee, Cong Liu, Junyoung Kim, Zhehuan Chen,	Vu Minh Hieu Phan, Yutong Xie, Yuankai Qi, Lingqiao Liu,	748
691	Yingcheng Sun, James R Rogers, Wendy K Chung, and	Liyang Liu, Bowen Zhang, Zhibin Liao, Qi Wu, Minh-Son	749
692	Chunhua Weng. 2022. Deep learning for rare disease:	To, and Johan W Verjans. 2024. Decomposing disease	750
693	A scoping review. <i>Journal of biomedical informatics</i> ,	descriptions for enhanced pathology detection: A multi-	751
694	135:104227.	aspect vision-language pre-training framework. In <i>Pro-</i>	752
695	Chunyu Li, Cliff Wong, Sheng Zhang, Naoto Usuyama,	<i>ceedings of the IEEE/CVF Conference on Computer Vision</i>	753
696	Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung	<i>and Pattern Recognition</i> , pages 11492–11501.	754
697	Poon, and Jianfeng Gao. 2024a. Llava-med: Training a	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh,	755
698	large language-and-vision assistant for biomedicine in one	Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda	756
699	day. <i>Advances in Neural Information Processing Systems</i> ,	Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021.	757
700	36.	Learning transferable visual models from natural language	758
701	Jiachun Li, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo	supervision. In <i>International conference on machine learn-</i>	759
702	Chen, Daojian Zeng, Kang Liu, and Jun Zhao. 2024b. Fo-	<i>ing</i> , pages 8748–8763. PmLR.	760
703	cus on your question! interpreting and mitigating toxic	Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan	761
704	cot problems in commonsense reasoning. <i>arXiv preprint</i>	Zong, Letian Wang, Yu Liu, and Hongsheng Li. 2024.	762
705	<i>arXiv:2402.18344</i> .	Visual cot: Advancing multi-modal language models with a	763
706	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a.	comprehensive dataset and benchmark for chain-of-thought	764
707	Blip-2: Bootstrapping language-image pre-training with	reasoning. <i>Advances in Neural Information Processing</i>	765
708	frozen image encoders and large language models. In <i>In-</i>	<i>Systems</i> , 37:8612–8642.	766
709	<i>ternational conference on machine learning</i> , pages 19730–	George Shih, Carol C Wu, Safwan S Halabi, Marc D Kohli, Lu-	767
710	19742. PMLR.	ciano M Prevedello, Tessa S Cook, Arjun Sharma, Judith K	768
711	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang,	Amorosa, Veronica Arteaga, Maya Galperin-Aizenberg,	769
712	and You Zhang. 2023b. Chatdoctor: A medical chat model	and 1 others. 2019. Augmenting the national institutes	770
713	fine-tuned on a large language model meta-ai (llama) using	of health chest radiograph dataset with expert annotations	771
714	medical domain knowledge. <i>Cureus</i> , 15(6).	of possible pneumonia. <i>Radiology: Artificial Intelligence</i> ,	772
715	Chenyu Lian, Hong-Yu Zhou, Dongyun Liang, Jing Qin, and	1(1):e180041.	773
716	Liansheng Wang. 2025. Efficient medical vision-language	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi,	774
717	alignment through adapting masked vision models. <i>IEEE</i>	Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tan-	775
718	<i>Transactions on Medical Imaging</i> .	wani, Heather Cole-Lewis, Stephen Pfohl, and 1 others.	776
719	Haotian Liu, Chunyu Li, Qingyang Wu, and Yong Jae Lee.	2023. Large language models encode clinical knowledge.	777
720	2023. Visual instruction tuning. <i>Advances in neural infor-</i>	<i>Nature</i> , 620(7972):172–180.	778
721	<i>mation processing systems</i> , 36:34892–34916.	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste	779
722	Yanming Liu, Xinyue Peng, Tianyu Du, Jianwei Yin, Wei-	Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, An-	780
723	hao Liu, and Xuhong Zhang. 2024. Era-cot: improv-	drew M Dai, Anja Hauth, Katie Millican, and 1 others.	781
724	ing chain-of-thought through entity relationship analysis.	2023. Gemini: a family of highly capable multimodal	782
725	<i>arXiv preprint arXiv:2403.06932</i> .	models. <i>arXiv preprint arXiv:2312.11805</i> .	783
726	Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam	Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, An-	784
727	Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xi-	drew Y Ng, and Pranav Rajpurkar. 2022. Expert-level de-	785
728	anzhi Du, Futang Peng, Anton Belyi, and 1 others. 2024.	tection of pathologies from unannotated chest x-ray images	786
729	Mml: methods, analysis and insights from multimodal llm	via self-supervised learning. <i>Nature biomedical engineer-</i>	787
730	pre-training. In <i>European Conference on Computer Vision</i> ,	<i>ing</i> , 6(12):1399–1406.	788
731	pages 304–323. Springer.	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann	789
732	Debjyoti Mondal, Suraj Modi, Subhadarshi Panda, Rituraj	LeCun, and Saining Xie. 2024. Eyes wide shut? exploring	790
733	Singh, and Godawari Sudhakar Rao. 2024. Kam-cot:	the visual shortcomings of multimodal llms. In <i>Proceed-</i>	791
734	Knowledge augmented multimodal chain-of-thoughts rea-	<i>ings of the IEEE/CVF Conference on Computer Vision and</i>	792
735	soning. In <i>Proceedings of the AAAI conference on artificial</i>	<i>Pattern Recognition</i> , pages 9568–9578.	793
736	<i>intelligence</i> , volume 38, pages 18798–18806.	Khoa A Tran, Olga Kondrashova, Andrew Bradley, Eliza-	794
737	Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga,	beth D Williams, John V Pearson, and Nicola Waddell.	795
738	Yash Dalmia, Jure Leskovec, Cyril Zakka, Eduardo Pontes	2021. Deep learning in cancer diagnosis, prognosis and	796
739	Reis, and Pranav Rajpurkar. 2023. Med-flamingo: a multi-	treatment selection. <i>Genome medicine</i> , 13:1–17.	797
740	modal medical few-shot learner. In <i>Machine Learning for</i>	Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammad-	798
741	<i>Health (ML4H)</i> , pages 353–367. PMLR.	hadi Bagheri, and Ronald M Summers. 2017. Chestx-ray8:	799
742	Maya Pavlova, Naomi Terhjan, Audrey G Chung, Andy Zhao,	Hospital-scale chest x-ray database and benchmarks on	800
743	Siddharth Surana, Hossein Aboutalebi, Hayden Gunraj, Ali	weakly-supervised classification and localization of com-	801
744	Sabri, Amer Alaref, and Alexander Wong. 2022. Covid-	mon thorax diseases. In <i>Proceedings of the IEEE confer-</i>	802
745	net cxr-2: An enhanced deep convolutional neural network	<i>ence on computer vision and pattern recognition</i> , pages	803
746	design for detection of covid-19 cases from chest x-ray	2097–2106.	804
747	images. <i>Frontiers in Medicine</i> , 9:861680.	Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and	805
		Weidi Xie. 2023. Medklip: Medical knowledge enhanced	806
		language-image pre-training for x-ray diagnosis. In <i>Pro-</i>	807
		<i>ceedings of the IEEE/CVF International Conference on</i>	808
		<i>Computer Vision</i> , pages 21372–21383.	809

810 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,
811 Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei
812 Huang, Haoran Wei, and 1 others. 2024. Qwen2. 5 techni-
813 cal report. *arXiv preprint arXiv:2412.15115*.

814 Ling You, Wenxuan Huang, Xinni Xie, Xiangyi Wei, Bangyan
815 Li, Shaohui Lin, Yang Li, and Changbo Wang. 2025.
816 Timesoccer: An end-to-end multimodal large language
817 model for soccer commentary generation. *arXiv preprint*
818 *arXiv:2504.17365*.

819 Sheng Zhang, Yanbo Xu, Naoto Usuyama, Jaspreet Bagga,
820 Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen
821 Valluri, Cliff Wong, and 1 others. 2023a. Large-scale
822 domain-specific pretraining for biomedical vision-language
823 processing. *arXiv preprint arXiv:2303.00915*, 2(3):6.

824 Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yan-
825 feng Wang. 2023b. Knowledge-enhanced visual-language
826 pre-training on chest radiology images. *Nature Communi-*
827 *cations*, 14(1):4542.

828 Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D
829 Manning, and Curtis P Langlotz. 2022. Contrastive learn-
830 ing of medical visual representations from paired images
831 and text. In *Machine Learning for Healthcare Conference*,
832 pages 2–25. PMLR.

833 Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruva Ghosh,
834 Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy.
835 2024. Why are visually-grounded language models bad
836 at image classification? *The Thirty-eighth Annual Confer-*
837 *ence on Neural Information Processing Systems*.

838 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mo-
839 hamed Elhoseiny. 2023. Minigt-4: Enhancing vision-
840 language understanding with advanced large language mod-
841 els. *arXiv preprint arXiv:2304.10592*.

A Related Work

Multi-modal Large Language Models. Inspired by the exceptional reasoning capabilities of large language models (LLMs), researchers are actively exploring ways to extend these abilities to the visual domain, driving advancements in multimodal LLMs. With the release of GPT-4 (Vision) (Achiam et al., 2023) and Gemini (Team et al., 2023), these models have demonstrated remarkable multimodal understanding and generation capabilities, further fueling research in this field.

To bridge the gap between vision encoders and LLMs, BLIP-2 (Li et al., 2023a) introduces a Q-Former that transforms image features into a format compatible with LLMs, enabling seamless integration with text embeddings. LLaVA (Liu et al., 2023) and MiniGPT-4 (Zhu et al., 2023) further enhance generalization and task performance by leveraging large-scale multimodal pretraining, followed by instruction tuning for specific applications. In the medical domain, LLMs have shown immense potential for advancing research and practical applications. Med-Flamingo (Moor et al., 2023) extends Flamingo to the medical field by pretraining on multimodal knowledge sources spanning multiple medical disciplines. LLaVA-Med (Li et al., 2024a) refines image-text pairs from PMC-15M (Zhang et al., 2023a) and trains a biomedical-specialized MLLM using a limited dataset, building upon the pre-trained parameters of LLaVA. Similarly, Med-PaLM (Singhal et al., 2023) fine-tunes PaLM (Chowdhery et al., 2023) using domain-specific medical instructions, demonstrating strong performance under human evaluation frameworks. Other notable models, such as Chat-Doctor (Li et al., 2023b) and Med-Alpaca (Han et al., 2023), have been tailored for medical question-answering and dialogue applications.

Despite the significant progress of MLLMs, several challenges remain (McKinzie et al., 2024; Tong et al., 2024; Zhang et al., 2024; He et al., 2025). Recent studies (Zhang et al., 2024; He et al., 2025) highlight the suboptimal performance of MLLMs in image classification, particularly in fine-grained category recognition. We find that this issue is especially pronounced in the medical domain, where precise classification is crucial for medical applications. To address these shortcomings, we are refining traditional MLLM training paradigms to enhance classification performance and improve fine-grained category comprehension.

Prompt Engineering. Prompting enhances the ability of pre-trained large language models (LLMs) to understand tasks by incorporating language instructions into the input text (Mondal et al., 2024; Shao et al., 2024; Liu et al., 2024; Li et al., 2024b). Recently, prompt-based techniques have also been applied to vision-language models to improve performance. In medical vision-language models (VLMs), GloRIA (Huang et al., 2021) generates a set of textual prompts to describe potential subtypes, severity levels, and anatomical locations for each disease category. Med-KLIP (Wu et al., 2023) enhances model performance by retrieving descriptions of medical entities from the UMLS knowledge base (Bodenreider, 2004). CARZero (Lai et al., 2024) introduces a prompt-alignment strategy based on LLMs, integrating prompt templates into the training dataset to ensure alignment during both training and inference. MAVL (Phan et al., 2024) uses visual descriptions of pathological features to guide the model in effectively detecting diseases in medical images.

Although these approaches have successfully improved model performance through prompt-based strategies, they all rely on external models or expert knowledge, without fully leveraging the model’s intrinsic understanding capabilities. Fortunately, recent research on LLaVA-Med (Li et al., 2024a) has demonstrated remarkable domain-specific conversational abilities, proving that it possesses a certain level of medical knowledge. Building upon LLaVA-Med (Li et al., 2024a), we further explore the feasibility of utilizing the model’s inherent comprehension to enhance zero-shot medical classification performance.

B Dataset Details

MIMIC-CXR v2 (Johnson et al., 2019). In our experiments, we pre-trained the model using the MIMIC-CXR, a publicly available collection of chest radiographs paired with corresponding radiology text reports. The MIMIC-CXR dataset comprises 377,110 images corresponding to 227,835 radiographic studies from 65,379 patients. Since all downstream tasks utilize frontal-view images, we exclude all lateral-view images from the dataset. Moreover, we selectively retain only the findings and impressions sections from these reports.

ChestX-ray14 (Wang et al., 2017). ChestX-ray14 consists of 112,120 frontal-view chest X-ray im-

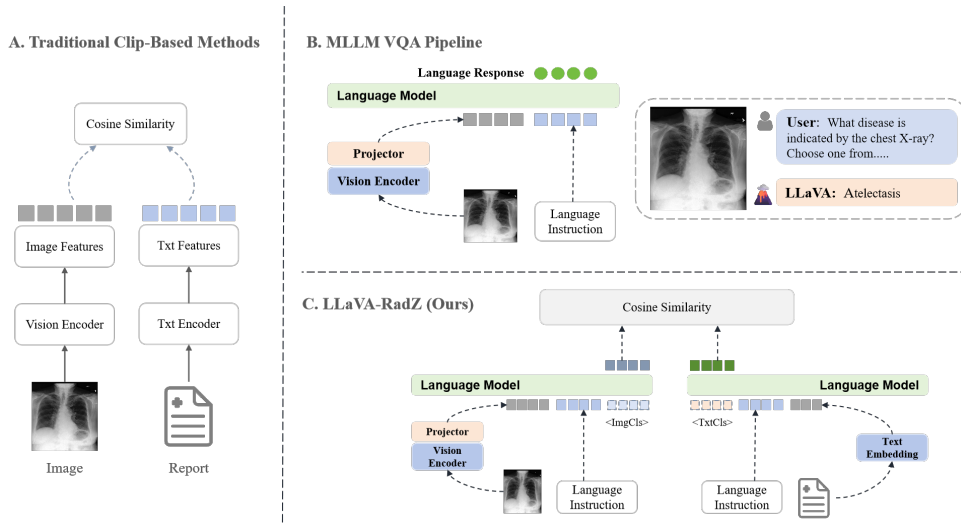


Figure 4: Framework comparison of traditional CLIP-based methods, MLLM VQA pipeline, and the proposed LLaVA-RadZ.

ages from 30,805 unique patients, collected between 1992 and 2015. The official test set, comprising 22,433 images, has been meticulously annotated by board-certified radiologists. For evaluation purposes, we restrict our testing to the official test set.

RSNA Pneumonia (Shih et al., 2019). RSNA Pneumonia includes over 260,000 frontal-view chest X-rays with annotated pneumonia masks, collected by the Radiological Society of North America (RSNA). This dataset supports both pneumonia segmentation and classification tasks (Wu et al., 2023; Phan et al., 2024). We partition the dataset into training, validation, and test sets with a ratio of 0.6/0.2/0.2, respectively.

SIIM-ACR Pneumothorax (sui, 2019). SIIM-ACR Pneumothorax contains 12,954 chest X-ray images, along with image-level pneumothorax annotations and pixel-level segmentation masks where pneumothorax is present. Like the RSNA Pneumonia dataset, it can be used for both classification and segmentation tasks. We divide the dataset into training, validation, and test sets with a ratio of 0.6/0.2/0.2.

CheXpert (Irvin et al., 2019). CheXpert contains 224,316 chest X-ray images from 65,240 patients, collected by Stanford Hospital. The official test set includes images from 500 patients, annotated through consensus by five board-certified radiologists. We evaluated all disease categories in this test dataset.

COVIDx CXR-2 (Pavlova et al., 2022) and COVID Rural (Desai et al., 2020). The COVIDx CXR-2 and COVID Rural are designed for eval-

uating COVID-19 diagnosis. COVIDx CXR-2 (Pavlova et al., 2022) consists of 29,986 images from 16,648 COVID-19 patients, each labeled with a classification tag. The dataset is split into training, validation, and test sets with a ratio of 0.7/0.2/0.1, used for evaluating classification performance. The COVID Rural dataset contains over 200 chest X-ray images with annotated segmentation masks, used for the COVID-19 segmentation task. This dataset is partitioned into training, validation, and test sets with a ratio of 0.6/0.2/0.2.

C Medical Category Semantic Vector Library

We draw inspiration from the work of Med-CLIP (Wu et al., 2023) and incorporate 75 frequently occurring medical entities from clinical reports as input to our model. By designing prompts, we stimulate the model’s intrinsic medical knowledge, enabling it to infer the semantic representations of various entity categories. The resulting semantic descriptions of these 75 medical entities are presented in table 8.

Furthermore, to achieve a more precise representation of major disease categories, we construct a dedicated disease semantic vector library, which facilitates a more nuanced understanding of disease-related semantics. The generated disease descriptions are detailed in table 7.

D Implementation Details

Unless otherwise specified, we use LLaVA-Med (Li et al., 2024a) as the foundational MLLM \mathcal{F} . We employ the LoRA strategy for parameter-

Table 6: Training cost comparison of LLaVA-Med and LLaVA-RadZ under different configurations.

Model	LLaVA-Med	LLaVA-RadZ (wo/DKAM)	LLaVA-RadZ
Training Cost (GPU-hours)	32.80	33.89	34.26

efficient fine-tuning, with training managed via the DeepSpeed engine.

For optimization, we utilize the AdamW optimizer with a learning rate of $2e-5$ and no weight decay. A cosine learning rate decay schedule is applied, with 3% of the total training steps allocated for warm-up. The number of special tokens for images $\langle \text{ImgCls} \rangle$ is set to 4, while the number of special tokens for text $\langle \text{TxtCls} \rangle$ is set to 8. The temperature hyperparameter τ is configured as 0.05, and the loss weight coefficient λ is set to 0.5. Furthermore, the batch size per GPU is set to 64.

Building on these training configurations, we further analyzed the computational complexity of the LLaVA-RadZ framework, focusing on the impact of the proposed training strategies on training and inference efficiency. The training durations under different configurations are summarized in table 6.

As shown in table 6, under the same hardware setup (2x A100 80GB GPUs) and identical parameter configurations, LLaVA-RadZ introduces only about a 4.5% increase in training cost compared with the baseline model LLaVA-Med. This minimal overhead yields substantial gains in zero-shot generalization, demonstrating an excellent balance between efficiency and effectiveness.

In terms of inference efficiency, LLaVA-RadZ shows even more pronounced advantages. Its per-sample computational cost is approximately 4329.26 GFLOPs, noticeably lower than the baseline’s 5122.74 GFLOPs. Moreover, the average inference latency is reduced to 95.4 ms, compared with 539.3 ms for LLaVA-Med, corresponding to a **5.6x speed-up**.

Taken together, these results show that the training strategies introduced in LLaVA-RadZ not only enhance semantic alignment and domain knowledge integration but also significantly improve overall computational efficiency, providing a strong foundation for practical deployment.

E Use of LLMs

This paper employed large language models (i.e., ChatGPT, Claude) solely for language editing and polishing purposes, including but not limited to grammar checking, expression optimization, and text refinement. All core research content, includ-

ing experimental design, data analysis, and conclusion derivation, was carried out independently by the authors. The authors take full responsibility for the entire content of this paper and have thoroughly verified and validated all AI-assisted modifications.

F Ethical Considerations and Potential Risks

This work focuses on the research exploration of multimodal representation learning and zero-shot medical image disease recognition, and is not intended for direct clinical diagnosis or automated medical decision-making. Although the proposed model demonstrates promising results on public benchmark datasets, its predictions may still be inaccurate or incomplete. Direct deployment in real-world clinical settings without appropriate medical expert supervision may therefore pose potential risks.

Accordingly, this work should be used for research purposes only. Any potential practical application must be conducted under rigorous clinical validation, with qualified human experts involved and comprehensive risk assessment in place. Future research may further investigate more efficient, robust, and safer multimodal medical artificial intelligence methods.

Table 7: Semantic Descriptions of 14 Medical Disease Categories

Disease	Description
Fibrosis	Fibrosis refers to excessive deposition of collagen and extracellular matrix during abnormal tissue repair after inflammation or injury, leading to the replacement of normal lung tissue with reticular or band-like high-density shadows, commonly seen in the lower and peripheral lungs. Imaging may show honeycombing and traction bronchiectasis. Clinically, patients often present with progressive dyspnea, dry cough, and reduced exercise tolerance.
Edema	Pulmonary edema refers to the abnormal accumulation of fluid in the pulmonary interstitium and alveoli, usually caused by cardiogenic or non-cardiogenic factors. Imaging shows patchy or 'bat-wing' distributed heterogeneous high-density shadows in the middle or entire lung, often accompanied by Kerley lines and cardiac enlargement. Clinically, patients typically experience acute dyspnea, cough, cyanosis, and bilateral lung crackles.
Pneumothorax	Pneumothorax refers to the presence of air in the pleural cavity, leading to partial or complete lung collapse. Imaging typically shows a low-density black air space along the pleura, with a clear demarcation from the normal lung tissue, along with lung collapse. In tension pneumothorax, mediastinal shift may occur. Clinically, patients often present with sudden unilateral chest pain, dyspnea, and decreased breath sounds, sometimes accompanied by subcutaneous emphysema.
Cardiomegaly	Cardiomegaly refers to the enlargement of the heart due to hypertension, cardiomyopathy, or valvular disease, causing chamber dilation or wall thickening. Imaging shows significant cardiac enlargement with an expanded and smooth contour, often marked by an increased cardiothoracic ratio, potentially accompanied by pulmonary congestion and bronchial congestion. Clinically, patients may experience reduced exercise tolerance, dyspnea, lower limb edema, and arrhythmias.
Atelectasis	Atelectasis refers to the collapse of part or all of the lung tissue due to airway obstruction, external thoracic pressure, or intrapulmonary pathology. Imaging shows increased local lung density, volume reduction, bronchial displacement, and visceral pleural traction, commonly affecting the lower lobes. Clinically, patients may exhibit rapid shallow breathing, localized decreased or absent breath sounds, and a history of recent surgery or inadequate airway clearance.
Nodule	A lung nodule is a localized lesion with a diameter of less than 3 cm. Imaging typically shows a round or oval localized density, with either well-defined or spiculated edges. Some nodules may contain calcifications or low-density necrotic areas. Clinically, most patients are asymptomatic, but growing or malignant nodules may present with cough and hemoptysis.
Emphysema	Emphysema is a chronic obstructive pulmonary disease caused by the permanent destruction of alveolar walls and airspace enlargement. Imaging shows scattered or diffuse low-density areas in both lungs, reduced lung markings, often with bullae or cystic lesions, a flattened diaphragm, and hyperinflated lungs. Clinically, patients typically have a history of chronic cough, sputum production, and progressive dyspnea, often associated with smoking or long-term occupational exposure.
No Finding	No finding refers to the absence of radiographic abnormalities detected in the chest X-ray.

Continued on next page

Disease	Description
Mass	A mass refers to an abnormal localized tissue overgrowth. Imaging shows a focal high-density lesion, which may have regular or irregular shapes with spiculated margins, often accompanied by internal necrosis, calcification, or hemorrhage. Surrounding features may include bronchial distortion or lymphadenopathy. Clinically, patients may present with cough, weight loss, or hemoptysis, requiring further pathological examination.
Pleural Thickening	Pleural thickening refers to fibrotic or calcified pleural changes due to chronic inflammation, infection, or asbestos exposure. Imaging shows localized or diffuse thickening along the pleural surface, appearing as streaky or patchy high-density shadows, sometimes with nodular changes. Clinically, patients may be asymptomatic, but a history of pleuritis or exposure to harmful substances is often present.
Effusion	Pleural effusion refers to the abnormal accumulation of fluid in the pleural cavity, which may be caused by infection, heart failure, malignancy, or other inflammatory diseases. Typically seen in the lower lung fields and posterior chest cavity, imaging shows a homogeneous or layered fluid density with a clear meniscus sign, with CT revealing low-density regions. Severe effusion may cause lung compression or bronchial displacement. Clinically, patients may present with dyspnea, chest pain, and cough, with physical signs of reduced breath sounds, dull percussion, and abnormal auscultation.
Infiltration	Infiltration refers to localized or diffuse high-density changes in lung tissue due to inflammation, infection, or malignancy. Imaging typically shows patchy or ill-defined high-density areas, sometimes with a ground-glass appearance or consolidation, occasionally accompanied by air bronchograms or bronchial wall thickening. Clinically, patients may present with cough, fever, dyspnea, and fatigue, often with elevated inflammatory markers.
Pneumonia	Pneumonia refers to lung parenchyma inflammation caused by bacteria, viruses, fungi, or other microorganisms, leading to alveolar filling with inflammatory exudates. Imaging shows localized or patchy consolidation with irregular margins, often accompanied by air bronchograms, pleural reaction, and mild pleural effusion. Clinically, patients present with fever, cough, sputum production, chest pain, and fatigue, with elevated white blood cell counts and inflammatory markers.
Consolidation	Consolidation refers to the complete filling of alveolar spaces with liquid, pus, blood, or cellular material, replacing the normal air content. Imaging shows homogeneous, dense, well-defined opacities, often with air bronchograms and pleural reactions, sometimes with minimal pleural effusion. Clinically, patients often have fever, cough, sputum production, chest pain, and dyspnea, with significantly elevated inflammatory markers.

Table 8: Semantic Descriptions of 75 Medical Categories

Disease	Description
normal	Indicates that the structure appears within standard parameters without signs of pathology.
clear	The imaging reveals no obscuring abnormalities, ensuring clear visualization of the structure.
sharp	Boundaries are precisely defined, accentuating the distinct separation between tissues.
sharply	The structure is rendered with exceptional clarity, facilitating detailed evaluation.
unremarkable	No significant deviations or abnormalities are observed in the examined area.
intact	The structure remains whole and undamaged, with no disruption detected.
stable	The tissue exhibits consistent appearance over time without progressive changes.
free	Presence of extraluminal air in unexpected locations, possibly indicating a perforation.
effusion	Accumulation of fluid between the pleural layers, often reflecting an underlying pathology.
opacity	An area of increased radiodensity that obscures normal lung markings, suggesting fluid or tissue replacement.
pneumothorax	Air present in the pleural space that may lead to partial or complete lung collapse.
edema	Diffuse fluid accumulation within lung tissue, frequently associated with cardiac or inflammatory issues.
atelectasis	Collapse of lung segments resulting in volume loss and increased local density.
tube	A medical tube visible on imaging, such as for drainage or airway management.
consolidation	Region where alveolar air is replaced by fluid or cells, producing homogeneous density.
process	Denotes an active pathological condition altering the tissue's normal appearance.
abnormality	A generic term for any deviation from normal structure suggestive of disease.
enlarge	Indicates that a structure appears larger than typical normal values.
tip	The distal or pointed end of a structure or medical device.
low	Underinflation of the lungs, often implying a restrictive process.
pneumonia	Inflammatory infection of lung parenchyma, typically showing consolidation and air bronchograms.
line	A linear structure that may represent a fissure, pleural interface, or artifact.
congestion	Increased blood or fluid accumulation in tissues, often indicating impaired circulation.
catheter	A slender, flexible tube inserted for drainage or medication delivery, visible in imaging.
cardiomegaly	An enlarged cardiac silhouette, frequently associated with chronic heart conditions.
fracture	A break or discontinuity in bone structure evident on radiographs.
air	Regions of radiolucency indicating the presence of gaseous content.
tortuous	Describes a vessel or structure exhibiting excessive curvature or winding.
lead	The foremost or guiding portion of a device or anatomical feature.
disease	A general term for any pathological process affecting normal tissue function.
calcification	Deposition of calcium salts within tissue, appearing as bright foci on imaging.

Continued on next page

Disease	Description
prominence	An area that appears more pronounced than surrounding tissues, suggesting an increase in size or density.
device	Any implanted or externally applied apparatus used for diagnostic or therapeutic purposes.
engorgement	Excessive filling of vessels or tissues with blood, leading to a swollen appearance.
picc	A long, thin catheter introduced via a peripheral vein and advanced into the central circulation for long-term therapy.
clip	A small metallic or plastic fastener used during surgery to secure tissues or vessels.
elevation	An upward displacement or raised position of an anatomical structure relative to its usual location.
expand	Describes a structure that appears dilated or increased in volume.
nodule	A small, rounded lesion typically less than 3 cm in diameter that can be benign or malignant.
wire	A thin, flexible metallic strand often used in surgical fixation or as part of medical devices.
fluid	The presence of liquid within tissues or cavities, altering the normal radiographic appearance.
degenerative	Changes in tissue structure resulting from chronic wear, aging, or repeated stress.
pacemaker	An implanted device that regulates heart rhythm, visible through its leads and generator.
thicken	Describes a structure that appears denser or more layered, possibly due to fibrotic changes.
marking	Visible patterns or lines that may represent vascular or connective tissue features.
scar	Fibrotic tissue that replaces normal parenchyma following injury, typically seen as an irregular opacity.
hyperinflate	Denotes lungs that are over-expanded, often with increased radiolucency and flattened diaphragms.
blunt	Loss of sharp definition in anatomical borders, leading to a less distinct appearance.
loss	Indicates a reduction or absence of normal tissue volume or density.
widen	Suggests that a structure or space is broader than the standard measurement.
collapse	A significant reduction or complete loss of volume in lung tissue due to obstruction or injury.
density	Reflects the compactness of a tissue, with higher density appearing whiter on radiographs.
emphysema	A chronic condition marked by alveolar wall destruction and abnormal enlargement of air spaces.
aerate	Indicates that the lung tissue is adequately filled with air, supporting effective gas exchange.
mass	A malignant tumor arising from lung tissue, typically presenting as an irregular mass with possible cavitation.
crowd	Compaction of airways and vessels, often due to volume loss or infiltrative processes.
infiltrate	Diffuse or patchy opacities in the lung that suggest inflammation, infection, or neoplastic involvement.

Continued on next page

Disease	Description
obscure	Describes anatomical structures that are not clearly visualized, often due to overlapping tissues or technical factors.
deformity	An abnormal shape or structure resulting from congenital anomalies, trauma, or disease progression.
hernia	The protrusion of an organ or tissue through an abnormal opening in the surrounding structure.
drainage	The process or presence of fluid removal from a body cavity, often via an inserted tube.
distention	Abnormal expansion or swelling of a structure due to accumulation of fluid or gas.
shift	Displacement of anatomical structures from their usual positions, indicating mass effect or volume change.
stent	A small mesh tube used to maintain the patency of a vessel or duct.
pressure	The force exerted per unit area by fluids or tissues, which can influence organ function.
lesion	Any abnormal area of tissue that deviates from the standard architecture, potentially indicative of pathology.
finding	A generic term for an observed abnormality or noteworthy feature on imaging.
borderline	The heart appears at the upper limit of normal size, without clear evidence of enlargement.
hardware	Any implanted or externally attached device used for diagnostic, therapeutic, or supportive purposes.
dilation	The widening or expansion of a hollow structure, often reflecting increased internal pressure.
chf	A clinical syndrome characterized by the heart's reduced pumping ability, leading to systemic fluid accumulation.
redistribution	A shift in the normal pattern of blood or air distribution within the lungs, often due to altered hemodynamics.
aspiration	Inhalation of foreign material into the airways, potentially leading to inflammatory or infectious complications.
rare diseases	Conditions that occur infrequently in the population and often require specialized diagnostic and management approaches.
Covid-19	An infectious disease caused by the SARS-CoV-2 virus, with a broad spectrum of respiratory and systemic manifestations.