

Rad-Flamingo: A Multimodal Prompt driven Radiology Report Generation Framework with Patient-Centric Explanations

Anonymous ACL submission

Abstract

In modern healthcare, radiology plays a pivotal role in diagnosing and managing diseases. However, the complexity of medical imaging data and the variability in interpretation can lead to inconsistencies and a lack of patient-centered insights in radiology reports. To address these challenges, we propose a novel multimodal prompt-driven report generation framework that integrates diverse data modalities—such as medical images, and clinical notes—to produce comprehensive and context-aware radiology reports. Our framework leverages innovative prompt engineering techniques to guide vision-language models in synthesizing relevant information, ensuring the generated reports are not only accurate but also tailored to individual patient profiles. A key feature of our framework is its ability to provide patient-centric explanations, offering clear and personalized insights into diagnostic findings and their implications. Experimental results demonstrate this framework’s effectiveness in enhancing report quality, improving understandability, and could foster better patient-doctor communication. This approach represents a significant step toward more intelligent, transparent, and human-centered medical AI systems.

1 Introduction

Radiology reports form the basis for clinical diagnostics and guide medical experts in treating patients. Despite their significance, creating radiology reports is a labor-intensive and expert-intensive process frequently plagued with human errors and differing details based on the radiologist’s level of experience. Given the very low ratio of radiologists to patients, the laborious process of creating full text radiology reports ends up being one of the workflow’s largest obstacles (US, China, and India is 1:10,000, 1:14,772, and 1:100,000, respectively) (Arora, 2014). Given the huge number of cases and

the shortage of radiology experts, time-efficiently generating reports is a major hurdle worldwide. Towards this goal, there has been a huge attempt from both industry and academia, with the landscape of AI-based report generation having seen exponential growth in recent times (Messina et al., 2022). This growth is owed to the evolving capabilities of large language models and vision language models (VLMs) in particular. VLMs have revolutionized the fields of computer vision and natural language processing by integrating visual perception and language understanding. These models have showcased exceptional abilities on a variety of tasks, such as image captioning (Hossain et al., 2019), visual question answering (Lu et al., 2023), and visual common sense reasoning (Zellers et al., 2018). By ingestion of vast amounts of image and text data, these models can learn rich visual representations and align them with the textual token space (Zhu et al., 2023; Alayrac et al., 2022; Radford et al., 2021; Wang et al., 2022c) to generate texts which are consistent with the image. Fine-tuning them with task-specific data improves their alignment with specialized tasks and user needs. VLMs such as (Thawakar et al., 2024; Moor et al., 2023) show promising efficacy in aligning image with text for medical use cases.

1.1 Motivation

Radiology reports are crucial for clinical decision-making, offering critical insights into a patient’s health. VLMs find an excellent application in generation of radiology reports. However, all generative pre-trained models are opaque by design. As report generation is a crucial stage of the medical decision making process, transparency becomes a necessity for such systems. Report generation systems which are able to generate reports with explanations are better placed to build trust and acceptability. This in turn, would help towards the large scale integration of such systems into

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

medical workflows. Zhao et al. (2023) defines various types of explainability namely, gradient-based, perturbation-based and attention-based to name a few. We observe that explainability in case of radiology report generation gets bifurcated into two types: patient-centric and expert-centric. Patient centric explanations are lucid generated texts, that paraphrases medical keywords in the report. Towards this we leverage the language generation capabilities of VLMs which possess an unique feature of generating text in coherence to the prompted image. We employ this unique feature to generate coherent reports and patient centric explanations. Further, recent research has demonstrated that large language models can also rationalize their own generation (Wiegrefe et al., 2021) giving the model an ability to give natural language explanations for its own generated responses. Importantly, generating radiology reports using prompting strategies let alone, multimodal prompting is an under-explored domain. Driven by this motivation, we developed a two step procedure to come up with a system to generate radiology reports along with patient-centric explanations. In the first stage we design few-shot prompts following the standard in-context learning template. For this stage we take a fine-tuned open source VLM model Mini-GPT4 (Zhu et al., 2023) fine-tuned on MIMIC-CXR dataset (Johnson et al., 2019). This stage acts as the synthetic data generator, which annotates each of the image-report instance with a patient-centric explanation. For verifying the explanations we rely only on medical expert evaluations. Following this we go to stage two, where we propose our few-shot multimodal prompting strategy which generates a radiology report along with patient-centric explanations. We apply these few-shot learning capabilities to Med-Flamingo (a fine-tuned Flamingo model) (Moor et al., 2023) and provide our evaluation by utilizing both classical NLG metrics (BLEU, ROUGE, METEOR) and medical expert evaluation score. However, given the nature of medical texts, semantic similarity has paramount importance and therefore, we focus more on semantic similarity scores rather than lexical similarity.

Our contributions are:

1. An Augmented IUX dataset Demner-Fushman et al. (2015) with each of 3995 image-report instances annotated with a patient-centric explanation. We achieve this via a synthetic data generation pipeline which

are then evaluated by medical experts.

2. A multimodal prompt based VLM framework, **Rad-Flamingo**, for automated radiology report generation and patient-centric explanation. Our method improved quantitative and qualitative scores.
3. A first-of-its-kind multimodal in-context learning technique for self-rationalization by adding explicit medical knowledge to the prompt. To the best of our knowledge, this method incorporates explainability for prompt based radiology report generation resulting in a 2.4% increment in performance over existing few-shot prompting techniques.

2 Background and Definitions

Patient-Centric Explanations: Pathophysiology (McCance et al., 2019) is the study of the functional changes that occur in the body as a result of a disease or injury. It focuses on understanding the mechanisms by which diseases disrupt normal physiological processes. In heart failure, for instance, a reduction in cardiac output leads to compensatory mechanisms like fluid retention, which can cause symptoms like edema and shortness of breath. Therefore, such informations serve as a form of medical explanation with the generated report. We extend this idea to patient-centric explanations, where the pathophysiological explanations are provided along-with the medical reports for ease of understanding from the patients' perspective.

Self-Rationalization: Self-rationalization in large language models (LLMs) (Marasovic et al., 2022; Wiegrefe et al., 2021; Camburu et al., 2018) refers to their ability to generate explanations or justifications for their own outputs. This involves creating reasoning pathways that appear coherent, logical, and aligned with the responses they produce, even though these models do not possess true understanding or awareness. LLMs achieve this by leveraging their vast training data to mimic human reasoning patterns, constructing plausible rationales based on context, prior responses, and linguistic structures. However, these explanations do not serve as a pointer to the internal working of the model, they merely act as a justification to the output. In sensitive domains such as healthcare, an explanation, at the very least plays an important role towards building trust.

In-Context Learning: In-context learning refers to the ability of LLMs to perform tasks by understanding and extrapolating from examples provided within a prompt, without requiring explicit fine-tuning of the model. This technique leverages the model’s parametric knowledge and allows users to define the task through natural language instructions and a few input-output examples (often called few-shot learning). The model infers the pattern from the context and applies it to new instances during the same interaction. In-context learning demonstrates the flexibility of LLMs to adapt to diverse tasks, making them highly versatile for applications like text generation, question answering, and code synthesis (Dong et al., 2024).

3 Related Work

Report Generation: Radiology report generation has been receiving a lot of attention lately, and several models have been developed based on the encoder-decoder architecture that was first used for image captioning tasks (Vinyals et al., 2014; Xu et al., 2015; Pan et al., 2020). However, report generation poses additional challenges compared to image captioning, as medical reports are typically longer and coherent with respect to captions. In an encoder decoder setting it becomes very difficult to generate long-form reports coherent with the medical image. Furthermore, bias in medical datasets makes it difficult to generate comprehensive, long-form reports. To address these challenges, researchers have proposed various methods. Wang et al. (2021), introduced an image-text matching branch to facilitate report generation, utilizing report features to augment image characteristics and consequently minimize the impact of data bias. They also employed a hierarchical LSTM structure for the generation of long-form text. Chen et al. (2020a) and Wang et al. (2022b) introduced additional memory modules to store past information, which can be utilized during the decoding process to improve long-text generation performance. Another type of work aims to mitigate data bias by incorporating external knowledge information, with the most representative approach being the integration of knowledge graphs Li et al. (2019, 2023b); Huang et al. (2023); Liu et al. (2021); Zhang et al. (2020). Zhang et al. (2020) and Liu et al. (2021) combined pre-constructed graphs representing relationships between diseases and organs using graph neural networks, enabling more

effective feature learning for abnormalities. Li et al. (2023b) developed a dynamic approach that updates the graph with new knowledge in real-time. Huang et al. (2023) incorporated knowledge from a symptom graph into the decoding stage using an injected knowledge distiller.

These methods are able to generate reports as caption with very high accuracy. However, they do not have the ability of free-form text generation possessed by pretrained VLMs. Therefore, VLMs become very effective for free-form text generation.

Vision Language Models: A significant area of research in natural language processing (NLP) and computer vision is the exploration of vision language model (VLM) learning techniques. This VLM aims to bridge the gap between visual and textual information, enabling machines to understand and generate content that combines both modalities. Recent studies have demonstrated the potential of VLM models in various tasks, such as image captioning (Zhu et al., 2023), visual question answering (Liu et al., 2023; Maaz et al., 2024), and image generation (Zhang et al., 2023). Developing on these medical VLMs like (Li et al., 2023a) and (Abdin et al., 2024) show impressive performance on medical NLP use cases.

4 Methodology

We propose a two-stage methodology for generating radiology reports that generates patient-centric explanations, aiming to increase the report understanding for a non-expert reader.

In the first stage, as per Figure 1, we use a finetuned MiniGPT4 model to synthetically generate patient-centric explanations for each image report pair. The model is finetuned on MIMIC-CXR Johnson et al. (2019) dataset, a large-scale repository of chest X-ray images and corresponding reports in the form of findings and impressions. Finetuning allows the model to re-parameterize its weights to learn to align a chest X-ray to its corresponding report. Given this finetuned model, we design a three-shot prompt template to generate patient-centric explanations for an X-ray image and its corresponding report. Therefore, this stage appends all the existing dataset samples with a patient-centric explanation. The explanations generated are evaluated by medical-experts which allows us to use the dataset as a gold label for the second stage. In the second stage, we use this newly augmented dataset to

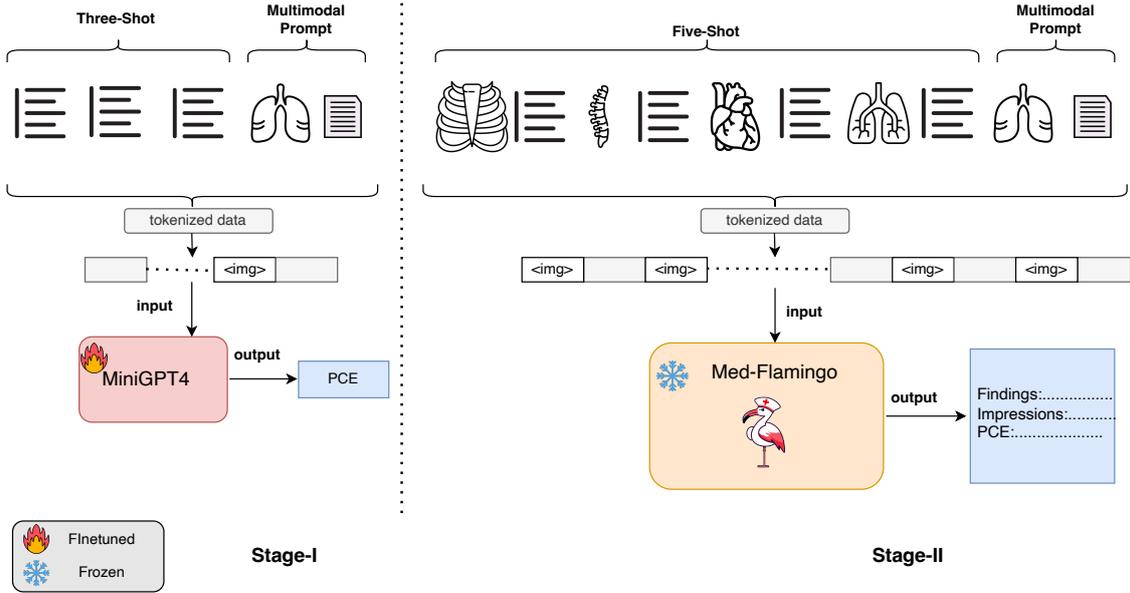


Figure 1: Stage I: Refers to the synthetic data generation stage, which annotate the existing IUX dataset with patient centric explanations. Stage II: Refers to the report generation stage where we design multimodal in-context prompts using the annotated data from stage I. Additionally, the fire symbol represents the finetuned model and ice symbol represent using frozen weights of a model not finetuned by us. PCE refers to the abbreviation of patient-centric explanation.

perform in-context learning with a vision-language model that has been pretrained on a medical question answering dataset. This approach allows the model to incorporate the nuances of patient-centric explanations while maintaining its ability to provide clinically accurate and detailed radiological interpretations. By combining fine-tuning with in-context learning, we aim to achieve a balance between medical precision and accessibility, ensuring the generated reports are useful for both healthcare providers and patients. This methodology showcases a novel application of vision-language models in potentially enhancing medical communication and patient engagement.

4.1 Stage I (Synthetic Data Generation)

To fine-tune the MiniGPT4 [Zhu et al. \(2023\)](#) model we follow the technique in [Thawakar et al. \(2024\)](#). We combine textual information from a medical large language model (LLM) and visual characteristics from a pre-trained medical vision encoder (VLM) given the X-ray. In particular, our large language model (LLM) is based on the recently developed Vicuna model ([Zheng et al., 2024](#)), and we use MedClip ([Wang et al., 2022c](#)) as a vision encoder.

Given an X-ray $x \in \mathbb{R}^{H \times W \times C}$, the vision en-

coder encodes the image as E_{img} . Then, the raw embeddings are transformed to an output dimension of 512 using a linear projection head.

$$V_p = f_v(E_{img}(x)) \quad (1)$$

where E_{img} is the vision encoder, f_v is the projection head. We use a trainable linear transformation layer to close the gap between the embedding space of the language decoder and image-level features, denoted as t . This layer transforms the image-level features, represented by V_p , into corresponding language-decoder embedding tokens, denoted as L_v :

$$L_v = t(V_p) \quad (2)$$

Following this we employ a few-shot prompting strategy to generate patient-centric explanations for a given image-report pair.

We follow a standard few-shot prompting strategy with three examples in the prompt. In the prompt we write Explanations as a placeholder for patient-centric explanation. The prompt template goes as follows:

Example 1:

```
Findings:.....
Impressions:.....
Explanations:.....
```

330 **Example 2:**

Findings:.....
Impressions:.....
Explanations:.....

331 **Example 3:**

Findings:.....
Impressions:.....
Explanations:.....

332 **Your Turn:**

Findings:.....
Impressions:.....
Looking at the Xray, findings
and impressions generate an explanation

333 For the synthetic data generation we consider the
334 IUX (Demner-Fushman et al., 2015) dataset, the
335 generated explanations are appended to each in-
336 stance of the IUX dataset. For designing the prompt
337 we sample three image-report (findings and im-
338 pressions) pairs from each of the disease classes.
339 We take assistance of medical experts to append
340 each of the samples with patient-centric explana-
341 tions. Subsequently, we pass the prompt as per
342 Stage I in Fig 1 for the fine-tuned model to learn in-
343 context. Fine-tuning the model on a large corpus,
344 such as MIMIC-CXR (Johnson et al., 2019), helps
345 the model to condition on the context provided in
346 the prompt. We provide the full prompt samples
347 in Appendix A. Therefore, the model is able to
348 generate good quality explanations tailoring to our
349 requirement. (the details are in appendix D). An
350 **Augmented Dataset** is now created which consists
351 of Image, report (Findings and Impressions) and
352 patient-centric explanation 2.

353 **4.2 Stage II (Radiology Report Generation)**

354 In this stage we follow the Med-Flamingo model
355 Moor et al. (2023) which is finetuned on a medical
356 visual question answering dataset. Med-Flamingo
357 is developed on the Open-Flamingo Awadalla et al.
358 (2023) architecture which possesses the ability of
359 few-shot learning from multimodal inputs. The
360 language modeling in Med-Flamingo is represented
361 in eq 3

$$p(y_\ell | x_{1:\ell-1}, y_{1:\ell-1}) = \prod_{\ell=1}^L p(y_\ell | y_{1:\ell-1}, x_{1:\ell-1}) \quad (3)$$

where y_ℓ refers to the ℓ_{th} language token, $y_{1:\ell-1}$ to
the set of prior language tokens, and $x_{1:\ell-1}$ to the
set of prior visual tokens. While fine-tuning, the
input is annotated in the form of interleaved image
text data, which makes it effective for multimodal
few-shot learning. We exploit this interleaved tem-
plate to design our proposed prompt as per Stage
II in Fig 1. The interleaved input prompt-design
while fine-tuning enables the model to condition
on the multi-modal context. We choose five exam-
ples for each disease class from the **Augmented
Dataset** compiled in stage I. Pivoting on the idea of
interleaved image text data prompt, we set up our
framework for multimodal in-context learning for
which the prompt template is demonstrated below:

378 **Example 1:**

Findings:.....
Impressions:.....
Explanations:.....

379 **Example 2:**

Findings:.....
Impressions:.....
Explanations:.....

380 **Example 3:**

Findings:.....
Impressions:.....
Explanations:.....

381 **Example 4:**

Findings:.....
Impressions:.....
Explanations:.....

382 **Example 5:**

Findings:.....
Impressions:.....
Explanations:.....

383 **Your Turn:**

Looking at the xray generate
findings and impressions and a explana-
tion

384 Prompt examples are provided in the Appendix B.

Metrics	Models						
	R2GEN (Chen et al., 2020b)	R2GenCMN (Chen et al., 2021)	Joint-TraiNet (Yang et al., 2023)	M2KT (Yang et al., 2022)	Open-Flamingo (Awadalla et al., 2023)	XProNet (Wang et al., 2022a)	Rad-Flamingo
BLEU-1	0.355	0.372	0.359	0.366	0.293	0.353	0.323
BLEU-2	0.223	0.233	0.226	0.213	0.195	0.221	0.232
BLEU-3	0.152	0.153	0.155	0.146	0.155	0.150	0.183
BLEU-4	0.103	0.105	0.102	0.104	0.071	0.105	0.081
METEOR	0.141	0.150	0.142	0.152	0.165	0.141	0.170
ROUGE	0.278	0.282	0.278	0.267	0.223	0.281	0.223

Table 1: Lexical similarity performance of Rad-Flamingo compared to baselines using classical metrics (BLEU, METEOR, ROUGE). The table highlights the limitations of these metrics in evaluating medical text generation, emphasizing the need for domain-specific semantic evaluation.

5 Experiments

5.1 Dataset

In **stage I** we consider the MIMIC-CXR (Johnson et al., 2019) dataset for fine-tuning. MIMIC-CXR dataset comprises 473,057 images and 206,563 reports from 63,478 patients. The official splits, i.e. 368,960 for training, 2,991 for validation, and 5,159 for testing are used for fine-tuning our model. Subsequent to this we follow our prompting technique (Section 4.1) to generate patient-centric explanations and append it to each instance of the IUX dataset (Demner-Fushman et al., 2015).

In **stage II** we use the **Augmented dataset** from the previous step and design our prompts as per Fig 1. The dataset consists of 7,470 chest X-Ray images and 3,955 radiology reports. The number of patients are equal to the number of reports however, each patient corresponds to two xray images i.e. frontal and lateral. Therefore, number of images are twice the number of reports. We append a patient-centric explanation to each of 3955 radiology reports.

5.2 Experimental Setup

In **stage-1** training, the model is fine-tuned to gain alignment between X-ray image features and corresponding reports by training over a large set of image-report pairs. The result obtained from the injected projection layer is considered as a gentle cue for our medically tuned VLM model, guiding it to generate appropriate report based on the finding and impression that match the given X-ray images. For preprocessing we follow Thawakar et al. (2024) where we utilize high quality interactive report summaries of MIMIC-CXR. The train set contains 213,514 image report pairs for training. During training, the model is trained for 320k total training steps with a batch size of 16 using 3

NVIDIA A100 (80GB) GPUs.

In **stage-II** we utilize predetermined prompts as shown in the previous section (4.2).

For each X-ray image instance we take the corresponding finding, impression and patient centric explanation and put it in the following format:

<image> Findings Impression Explanationlendl of chunkl.

Five of these aforementioned multimodal prompt were followed by the query prompt described below:

<image> + You are a helpful medical assistant. You are provided with images, findings, impressions and explanation. Looking at this image generate Findings, Impressions and Explanations.

6 Result and Analysis

Our results analyse the effectiveness of our multimodal prompt in generating reports with patient-centric explanation. Tables 1 and 2 compare the scores over the generated report and patient-centric explanations.

6.1 Lexical Metrics

In this section, we evaluate the quality of generated reports by **Rad-Flamingo** and compare them against baselines using classical lexical similarity metrics such as BLEU (Papineni et al., 2002), METEOR (Lavie and Agarwal, 2007), and ROUGE (Lin, 2004) as shown in Table 1. These metrics provide a convenient means of measuring word overlap and syntactic similarity between generated and reference texts. **Rad-Flamingo** performs similar to the baselines on lexical similarity metrics. However, these metrics find less application in medical domain. This arises due to their inability to account for the deeper semantic relevance and contextual accuracy required in specialized content, such as medical data. For example, the sentences

Metrics	Rad-Flamingo	Rad-Flamingo w/oI	Open-Flamingo	Open-Flamingo w/oI
BertScore	0.875	0.855	0.863	0.834
BioClinicalBertScore	0.895	0.879	0.885	0.854
RadGraphF1	0.285	0.273	0.279	0.269

Table 2: Performance comparison of Rad-Flamingo and Open-Flamingo models on clinical evaluation metrics using proposed multimodal few-shot prompting framework. The table includes ablation studies highlighting the impact of removing image modalities (w/oI) from the few-shot prompts. We do a metricise significance testing in Appendix D.2

"There is focal consolidation" and "There is no focal consolidation" are lexically very similar yet semantically very dissimilar. Therefore, semantic similarity plays a greater role in evaluating generated medical texts. Our evaluation emphasizes the performance of the Flamingo family of models (Moor et al., 2023) (Awadalla et al., 2023), as these models provide the essential few-shot learning capabilities needed for our prompt-based report generating framework.

Alternative vision-language models, such as Med-Phi (Abdin et al., 2024) and Med-LLaVA (Li et al., 2023a), were deliberately excluded as baselines from our analysis due to their lack of comparable few-shot learning features, making them less suitable for the scenarios we investigate. We stick to few-shot learning abilities as it plays a critical role in data scarce scenario such as medical domain. Our few-shot prompting technique show comparable performance in some of the lexical metrics. While these metrics offer a preliminary measure of performance, they do not fully reflect the real-world utility of generated medical texts. This analysis underscores the need for more domain-specific evaluation frameworks that can assess not only linguistic fluency and coherence but also the contextual alignment of generated texts in medical domain.

6.2 Semantic Metrics

We apply our few-shot prompting framework on all the available Flamingo family models namely, Rad-Flamingo and Open-Flamingo as shown in Table 2. The underlying model in Rad-Flamingo is Med-Flamingo. Med-Flamingo with our proposed multimodal prompt template is referred to as Rad-Flamingo. We choose semantic metrics for clinical evaluation like BioClinicalBERTScore¹

¹BioClinicalBERT is taken from huggingface. Underlying model is BioBERT trained on MIMIC III dataset. https://huggingface.co/emilyalsentzer/Bio_ClinicalBERT

(Lee et al., 2019), BERTScore (Zhang et al., 2019) and RadGraphF1 (Jain et al., 2021). In table 2 column Rad-Flamingo represents the setting where we prompt the Med-Flamingo model with proposed multimodal few-shot prompt. The following column, Rad-Flamingo w/oI represents the setting where we remove only the images from the few-shot prompt examples keeping other components of the prompt similar to Rad-Flamingo. Similar, ablation is carried out for the columns Open-Flamingo and Open-Flamingo w/oI.

Both the BERTScore and ClinicalBERTScore for Rad-Flamingo show a 1.4% increase compared to Open-Flamingo. This shows our proposed multimodal prompt template effectively generates report with better performance than existing models. Similar increase is found in case of RadGraphF1 scores. This result signifies the benefit of our proposed **multimodal prompt** template of Rad-Flamingo, over Open-Flamingo. To show the utility of multimodality in our prompt template, we remove the images from the examples and pass it to the Rad-Flamingo and Open-Flamingo models. Rad-Flamingo w/oI and Open-Flamingo w/oI represents those settings. We see the scores drop significantly by 2.4% percent indicating the utility of the multimodal prompt in integrating different data-modalities and helps the model to generate task-specific outputs. This approach effectively addresses challenges in both unimodal and multimodal data modes. So domain specific metrics are crucial to understand the utility of the multimodal prompt strategy developed by us. Therefore, we observe from the metrics that the semantic similarity scores help us analyze the performance better for task-specific output. Overall, the best performance is given by Rad-Flamingo as the underlying Med-Flamingo model is finetuned on medical data. However, comparing the scores with Open-Flamingo exhibits the effectiveness and utility of

Models	Rad-Flamingo
Cardiomegaly	3.44 ± 0.67
Pulmonary Atelectasis	3.33 ± 1.36
Nodules	3.21 ± 1.05
Opacity	2.06 ± 0.54
Calcified Granuloma	3.03 ± 0.41
Pulmonary Fibrosis	3.0 ± 0.63
Consolidation	3.2 ± 0.39
Pneumothorax	3.6 ± 0.8
Granuloma	3.4 ± 0.95
Bronchiectasis	3.25 ± 0.44

Table 3: The table presents the mean and standard deviation of scores provided by four medical professionals for each of the chosen disease class, highlighting the effectiveness of the proposed prompting method after stage II.

our proposed **multimodal prompt** framework.

6.3 Qualitative Evaluation

Owing to the subjective nature and the semantic complexity which medical data possesses, evaluation by medical expert becomes very important to have a rigorous examination of a proposed system. We consulted four expert-medical professionals to evaluate our generated reports and corresponding patient-centric explanations. Our expert evaluation was fixed on the primary criteria of Patient Centric Understandability. Following this we created five levels of grading: 1 (very poor), 2 (poor), 3 (good), 4 (very good), 5 (excellent). We choose the most common disease classes and provide to each of the medical experts. Subsequently, for each disease class we get four scores and the table shows a mean and standard deviation over these four scores. The expert evaluation also shows our proposed prompting method, gives promising performance. An output sample of our system is given in Appendix B.

6.4 Readability measure

To evaluate the patient centricity of the generated explanations we evaluate them using the Lexile Reading Measure (Stenner, 2023). A Lexile measure is a standardized score that assesses both the reading ability of individuals and the complexity of

Models	Rad-Flamingo	
	Generated	Ground Truth
Lexile Measure	69.28	63.6
CharBLEU	0.298	0.283

Table 4: The table highlights the readability and spelling accuracy of the generated texts, demonstrating their alignment with patient comprehension needs and medical domain standards.

written texts, represented on a scale typically ranging from below 200L to above 1600L. This measure helps educators, parents, and students identify reading materials that align with a reader’s current ability level, ensuring an appropriate level of challenge to support comprehension and skill development. We also evaluate on CharBLEU metrics (Denoual and Lepage, 2004) since in medical text spelling plays a crucial role. Table 4 shows a 8.9% increase in the readability of the generated explanations. The score provided is an average over all the ten selected diseases as per table 3. Averaging over all the values gives a rise in the readability measure, however for some disease class we do not find any increment. The overall readability of the explanations increase as per the scores. The explanations generated in stage II demonstrate improved readability compared to those from stage I, highlighting the effectiveness of our proposed prompt design in enhancing explanation clarity. This improvement underscores the potential of carefully crafted prompts in generating lucid explanations for users to comprehend and utilize.

7 Conclusion

Rad-Flamingo proposes a radiology report generation framework by combining multimodal data with prompt-driven methodologies and patient-centric explanations. This framework enhances accuracy, interpretability, and communication in medical imaging, setting new standards for personalized, explainable AI in healthcare. Rad-Flamingo highlights the potential of AI to automate routine reporting tasks, allowing radiologists to focus on complex cases and clinical decision-making. In conclusion, Rad-Flamingo shows a potential option for more efficient and impactful healthcare delivery. Future efforts can focus on better alignment of vision and language components of VLMs. Thereby, generating better reports with explanations.

602 Limitations

603 In this section we discuss the main limitations of
604 our proposed framework. A notable limitation in
605 our study is the absence of a number of VLMs
606 which possess the same few-shot learning capabil-
607 ity as the Flamingo family of models. This restricts
608 us from evaluating the generalizability of our ap-
609 proach. While our method shows promise, validat-
610 ing its performance against a diverse set of few-
611 shot models would provide deeper insights into its
612 strengths and weaknesses. The inclusion of these
613 models would also allow us to better understand
614 how our approach fares in broader scenarios and
615 under varying conditions, such as domain shifts or
616 noisy inputs.

617 Class imbalance in machine learning occurs
618 when certain classes dominate the training data,
619 causing the model to be biased toward these over-
620 represented classes and perform poorly on minor-
621 ity classes. This is particularly problematic in ap-
622 plications like medical diagnosis, where minority
623 classes are crucial, and can be addressed using tech-
624 niques like re-sampling, loss adjustment, or robust
625 algorithms.

626 Another constraint in our evaluation is the lack
627 of a direct comparison with ChatGPT, a widely
628 recognized benchmark in conversational AI. The
629 prompt template we use would be require high com-
630 putational and financial cost to perform a rigorous
631 analysis. These constraints underscore the need
632 for collaborative efforts and accessible research
633 resources to enable comprehensive benchmarking.

634 Ethical Considerations

635 The development of the Rad-Flamingo framework,
636 designed for multimodal prompt-driven radiology
637 report generation with patient-centric explanations,
638 adheres to the highest standards of ethical con-
639 duct. We prioritize patient privacy and data secu-
640 rity by ensuring that all medical information used
641 in the model is anonymized and handled in com-
642 pliance with relevant regulations. We augment a
643 standard dataset where each data sample is already
644 anonymized. Therefore, our proposed data augmen-
645 tation does not create any scare for identity leak-
646 age. The framework is designed to support, rather
647 than replace, radiologists and clinicians, with a
648 focus on improving diagnostic accuracy and foster-
649 ing transparent communication between healthcare
650 providers and patients.

651 We are committed to minimizing bias by ensur-

ing that the training data used is diverse, repre- 652
senting a wide range of demographic groups and 653
medical conditions. Additionally, patient explana- 654
tions generated by the model are designed to be 655
comprehensible and respectful, avoiding harmful 656
or misleading interpretations. 657

In all instances, human oversight is maintained 658
to validate outputs and ensure the model’s align- 659
ment with clinical practice standards and ethical 660
guidelines. 661

References 662

- 663 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed
664 Awadallah, Ammar Ahmad Awan, Nguyen Bach,
665 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat
666 Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck,
667 Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav
668 Chaudhary, Dong Chen, Dongdong Chen, Weizhu
669 Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng,
670 Parul Chopra, Xiyang Dai, Matthew Dixon, Ron-
671 nen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao,
672 Min Gao, Amit Garg, Allie Del Giorno, Abhishek
673 Goswami, Suriya Gunasekar, Emman Haider, Jun-
674 heng Hao, Russell J. Hewett, Wenxiang Hu, Jamie
675 Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi,
676 Xin Jin, Nikos Karampatziakis, Piero Kauffmann,
677 Mahoud Khademi, Dongwoo Kim, Young Jin Kim,
678 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi
679 Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui
680 Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu,
681 Weishung Liu, Xiaodong Liu, Chong Luo, Piyush
682 Madan, Ali Mahmoudzadeh, David Majercak, Matt
683 Mazzola, Caio César Teodoro Mendes, Arindam Mi-
684 tra, Hardik Modi, Anh Nguyen, Brandon Norick,
685 Barun Patra, Daniel Perez-Becker, Thomas Portet,
686 Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang
687 Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy,
688 Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil
689 Salim, Michael Santacroce, Shital Shah, Ning Shang,
690 Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia
691 Song, Masahiro Tanaka, Andrea Tupini, Praneetha
692 Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan
693 Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel
694 Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia
695 Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu,
696 Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang,
697 Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu,
698 Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen
699 Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan
700 Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219. 701
702
- 703 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
704 Antoine Miech, Iain Barr, Yana Hasson, Karel
705 Lenc, Arthur Mensch, Katherine Millican, Malcolm
706 Reynolds, Roman Ring, Eliza Rutherford, Serkan
707 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei,
708 Marianne Monteiro, Jacob L Menick, Sebastian
709 Borgeaud, Andy Brock, Aida Nematzadeh, Sahand

710	Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. Flamingo: a visual language model for few-shot learning . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 23716–23736. Curran Associates, Inc.	766
711		767
712		768
713		769
714		770
715		
716	Richa Arora. 2014. The training and practice of radiology in india: current trends. <i>Quant. Imaging Med. Surg.</i> , 4(6):449–450.	771
717		772
718		773
719	Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models . <i>Preprint</i> , arXiv:2308.01390.	774
720		775
721		776
722		777
723		778
724		779
725		780
726		781
727	Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations . In <i>Advances in Neural Information Processing Systems</i> , volume 31. Curran Associates, Inc.	782
728		783
729		784
730		785
731		
732	Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. 2021. Cross-modal memory networks for radiology report generation . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5904–5914, Online. Association for Computational Linguistics.	786
733		787
734		788
735		789
736		790
737		791
738		792
739		793
740	Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. 2020a. Generating radiology reports via memory-driven transformer . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1439–1449, Online. Association for Computational Linguistics.	794
741		795
742		796
743		797
744		798
745		799
746	Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. 2020b. Generating radiology reports via memory-driven transformer . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1439–1449, Online. Association for Computational Linguistics.	800
747		801
748		802
749		803
750		804
751		805
752	Dina Demner-Fushman, Marc D. Kohli, Marc B. Rosenman, Sonya E. Shooshan, Laritza M. Rodriguez, Sameer Kiran Antani, George R. Thoma, and Clement J. McDonald. 2015. Preparing a collection of radiology examinations for distribution and retrieval . <i>Journal of the American Medical Informatics Association : JAMIA</i> , 23 2:304–10.	806
753		807
754		808
755		809
756		810
757		811
758		
759	Etienne Denoual and Y. Lepage. 2004. Bleu in characters: Towards automatic mt evaluation in languages without word delimiters . In <i>International Joint Conference on Natural Language Processing</i> .	812
760		813
761		814
762		815
763	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A survey on in-context learning . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.	816
764		817
765		
	MD. Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2019. A comprehensive survey of deep learning for image captioning . <i>ACM Comput. Surv.</i> , 51(6).	818
		819
	Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. 2023. Kiut: Knowledge-injected u-transformer for radiology report generation . <i>2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 19809–19818.	820
		821
	Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Truong, D. Duong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew P. Lungren, Andrew Y. Ng, Curt P. Langlotz, and Pranav Rajpurkar. 2021. Radgraph: Extracting clinical entities and relations from radiology reports . <i>ArXiv</i> , abs/2106.14463.	822
		823
	Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih ying Deng, Yifan Peng, Zhiyong Lu, Roger G. Mark, Seth J. Berkowitz, and Steven Horng. 2019. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs . <i>Preprint</i> , arXiv:1901.07042.	
	Alon Lavie and Abhaya Agarwal. 2007. Meteor: an automatic metric for mt evaluation with high levels of correlation with human judgments. In <i>Proceedings of the Second Workshop on Statistical Machine Translation, StatMT '07</i> , page 228–231, USA. Association for Computational Linguistics.	
	Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining . <i>Bioinformatics</i> , 36(4):1234–1240.	
	Christy Y. Li, Xiaodan Liang, Zhiting Hu, and Eric P. Xing. 2019. Knowledge-driven encode, retrieve, paraphrase for medical image report generation . In <i>Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI'19/IAAI'19/EAAI'19</i> . AAAI Press.	
	Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day . <i>Preprint</i> , arXiv:2306.00890.	
	Mingjie Li, Bingqian Lin, Zicong Chen, Haokun Lin, Xiaodan Liang, and Xiaojun Chang. 2023b. Dynamic graph enhanced contrastive learning for chest x-ray report generation . <i>2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3334–3343.	

935 [free-text rationales](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10266–10284, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

940 Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho,
941 Aaron Courville, Ruslan Salakhudinov, Rich Zemel,
942 and Yoshua Bengio. 2015. [Show, attend and tell: Neural image caption generation with visual attention](#).
943 In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2048–2057, Lille, France. PMLR.

948 Shuxin Yang, Xian Wu, Shen Ge, S. Kevin Zhou, and
949 Li Xiao. 2022. [Radiology report generation with a learned knowledge base and multi-modal alignment](#).
950 *Preprint*, arXiv:2112.15011.

952 Yan Yang, Jun Yu, Jian Zhang, Weidong Han, Hanliang
953 Jiang, and Qing-An Huang. 2023. [Joint embedding of deep visual and semantic features for medical image report generation](#). *IEEE Transactions on Multimedia*, 25:167–178.

957 Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin
958 Choi. 2018. [From recognition to cognition: Visual commonsense reasoning](#). *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6713–6724.

962 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. [Adding conditional control to text-to-image diffusion models](#). *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3813–3824.

966 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.
967 Weinberger, and Yoav Artzi. 2019. [Bertscore: Evaluating text generation with bert](#). *ArXiv*, abs/1904.09675.

970 Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu,
971 Alan Loddon Yuille, and Daguang Xu. 2020. [When radiology report generation meets knowledge graph](#).
972 In *AAAI Conference on Artificial Intelligence*.

974 Haiyan Zhao, Hanjie Chen, F. Yang, Ninghao Liu, Huiqi
975 Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin,
976 and Mengnan Du. 2023. [Explainability for large language models: A survey](#). *ACM Transactions on Intelligent Systems and Technology*, 15:1 – 38.

979 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
980 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
981 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
982 Joseph E. Gonzalez, and Ion Stoica. 2024. Judging
983 llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA. Curran Associates Inc.

987 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and
988 Mohamed Elhoseiny. 2023. [Minigt-4: Enhancing vision-language understanding with advanced large language models](#). *Preprint*, arXiv:2304.10592.

A Appendix

991

A.1 Prompt sample from Stage I

992

Example 1: Findings: Borderline cardiomegaly. Midline sternotomy. Enlarged pulmonary arteries. Clear lungs. Inferior.

Impression: No acute pulmonary findings.

Explanation: The borderline cardiomegaly suggests slight enlargement of the heart, which could be due to chronic pressure overload (possibly from hypertension) or cardiomyopathy. The midline sternotomy indicates a previous cardiac surgery, which could relate to coronary artery disease or valve replacements, considering the context of other findings. Enlarged pulmonary arteries might suggest pulmonary hypertension, which could be secondary to left heart disease or chronic lung conditions, though the lungs are reported as clear. The absence of acute pulmonary findings despite these chronic alterations indicates that there are no new or worsening issues in the pulmonary system at this time.

993

Example 2: Findings: The cardiac contours are normal. basilar atelectasis. The lungs are clear. Thoracic spondylosis. Lower cervical arthritis.

Impression: Basilar atelectasis. No confluent lobar consolidation or pleural effusion.

Explanation: The chest X-ray shows basilar atelectasis, a partial collapse or incomplete expansion of the lower lung segments, likely due to shallow breathing or reduced lung volume. The lungs are otherwise clear, without any signs of confluent lobar consolidation or pleural effusion, suggesting no infection or fluid accumulation. The cardiac contours are normal, ruling out cardiomegaly or abnormal heart positioning. Degenerative changes, including thoracic spondylosis and lower cervical arthritis, are noted, indicating age-related or chronic wear-and-tear changes in the spine. The overall impression is of minor pulmonary findings with no acute cardiopulmonary disease.

994

Example 3: Findings: The cardiomeastinal silhouette is within normal limits for size and contour. The lungs are normally inflated without evidence of focal airspace disease, pleural effusion, or pneumothorax. Stable calcified granuloma within the right upper lung. No acute bone abnormality.

Impression: No acute cardiopulmonary process.

Explanation: The chest X-ray reveals a normal cardiomeastinal silhouette, indicating no enlargement or abnormality in heart or mediastinal contours. The lungs are well inflated, with no signs of focal airspace disease such as infection, fluid, or collapse. A stable, calcified granuloma in the right upper lung suggests a past, healed infection (e.g., tuberculosis or histoplasmosis) with no current active disease. The absence of pleural effusion, pneumothorax, or acute bone abnormalities supports the lack of any acute thoracic issues. Overall, there is no evidence of an acute cardiopulmonary process, and the findings are consistent with a stable, chronic condition.

995

Your Turn: Findings: Heart size within normal limits. No focal alveolar consolidation, no definite pleural effusion seen. No typical findings of pulmonary edema. No pneumothorax.

Impression: No acute cardiopulmonary findings.

Looking at the Xray, findings and impressions generate a patient- centric explanation

996

B Appendix

997

B.1 Prompt sample from Stage II

998

The images in figure 1 are placeholders to exhibit the multimodal prompt design. While experiments we took examples per disease class mentioned in the results section.

999

1000

Example 1: Findings: Heart size mildly enlarged, stable mediastinal and hilar contours. No focal alveolar consolidation, no definite pleural effusion seen. No typical findings of pulmonary edema. No pneumothorax.

Impressions: Mild cardiomegaly, no acute pulmonary findings.

Explanation: It is clear the patient is in an enlarged heart condition (mild cardiomegaly) and has been for some time. The pulmonary system is clear and no evidence of acute lung changes. The overall pathophysiology suggests that the patient is in a chronic, stable state with ongoing changes. The findings are normal and do not suggest any acute or severe events.

Example 2: Findings: The heart is top normal in size. The mediastinum is stable. The lungs are clear.

Impressions: Borderline cardiomegaly without acute disease.

Explanation: It is a case of borderline cardiomegaly, but it is difficult to correlate the finding of enlarged heart with pathophysiology. Borderline cardiomegaly is a radiological finding. It could mean that the heart size is in between normal and big (normal is 55 mm) but in this case it comes to 57 mm.

Example 3: Findings: Stable borderline cardiomegaly, stable mediastinal and hilar contours. No alveolar consolidation, no findings of pleural effusion or pulmonary edema. No pneumothorax.

Impressions: No acute cardiopulmonary findings.

Explanation: 1. No acute cardiopulmonary findings. 2. No pulmonary edema. 3. Bilateral patchy pulmonary opacities noted. 4. No focal alveolar consolidation, no definite pleural effusion seen. 5. No large pleural effusions. 6. Multifocal scattered bibasilar patchy and pulmonary opacities again noted. 7. No pneumothorax. 8. No left apical basal consolidation. 9. No left apical, basal patchy consolidation. 10. No pulmonary vascular congestion. 11. No pulmonary infarction.

Example 4: Findings: Persistent cardiomegaly. No abnormal airspace consolidation. Resolved interstitial edema. No pneumothorax or pleural effusion.

Impressions: Stable cardiomegaly without acute abnormality.

Explanation: No pneumothorax, no effusion, no infiltrate, no pulmonary congestion, no pleural erythema all point towards non-infectious etiology. Mild cardiomegaly without acute abnormality is also non-specific and without clinical significance. The unremarkable spine suggests degenerative changes and nothing else.

Example 5: Findings: The outside is normal except for slight cardiomegaly.

Impressions: Heart size upper limits normal. Lungs are clear. No evidence of active tuberculosis. No change from prior exam.

Explanation: Slight cardiomegaly. Clear lungs indicate no pulmonary congestion or active disease.

Your Turn: You are a helpful medical assistant. You are provided with images, findings, impressions and explanation. Looking at this image generate Findings, Impressions and Explanations

C Appendix

1007

C.1 Augmented IUX dataset instance

1008

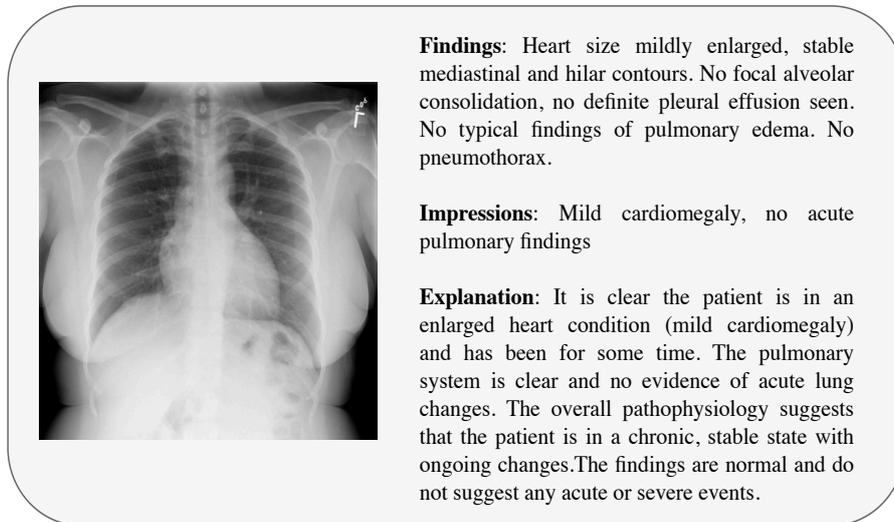


Figure 2: Augmented dataset instance showcasing input modalities (e.g., medical images, clinical text) and corresponding annotated outputs, illustrating the report (findings and impression) and patient-centric explanation

C.2 Radiology Report with patient-centric explanation generated by Rad-Flamingo

1009

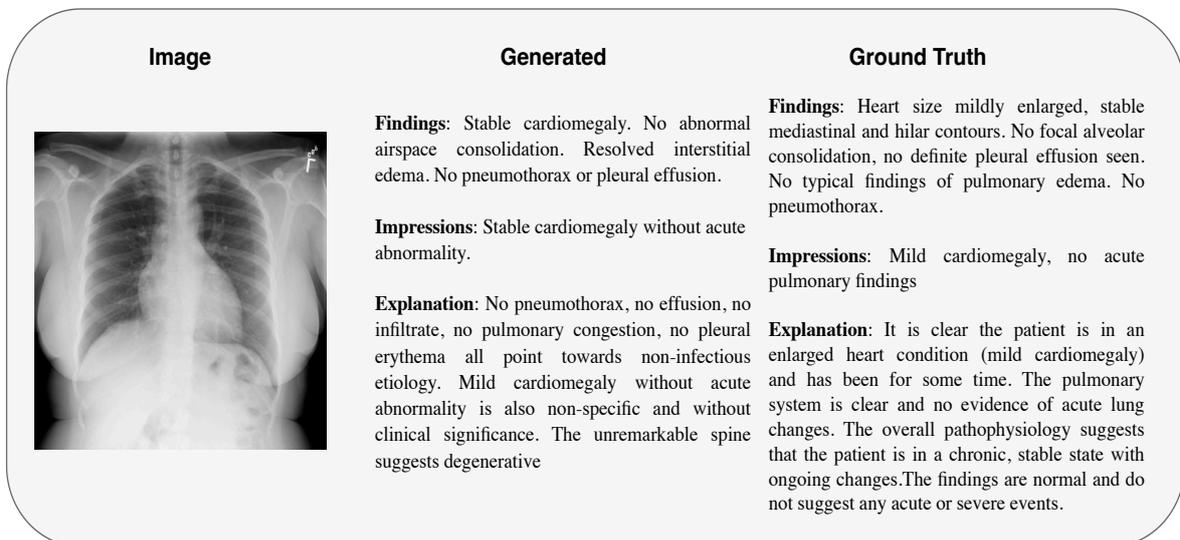


Figure 3: Example of output given by Rad-Flamingo. Image and ground truth are from the proposed augmented dataset.

D Appendix

1010

D.1 Medical Expert Evaluation for Stage I outputs

1011

D.2 Significance testing for Semantic Metrics

1012

Extending our analysis in the results section, we further provide significance testing for the BERTScore, BioClinicalBERTScore, and RadGraphF1 scores of Rad-Flamingo, Rad-Flamingo w/oI, Open-Flamingo,

1013

1014

Models	Rad-Flamingo
Cardiomegaly	3.72 ± 0.89
Pulmonary Atelectasis	3.15 ± 1.22
Nodules	3.38 ± 1.12
Opacity	2.11 ± 0.68
Calcified Granuloma	3.11 ± 0.58
Pulmonary Fibrosis	3.27 ± 0.73
Consolidation	3.05 ± 0.49
Pneumothorax	3.51 ± 0.92
Granuloma	3.52 ± 1.01
Bronchiectasis	3.18 ± 0.65

Table 5: The table presents the mean and standard deviation of scores provided by four medical professionals for each of the chosen disease class, highlighting the effectiveness of the proposed prompting method after stage I.

Metrics	F-statistic	p-value
BioClinicalBertScore	30.00	0.0001
BertScore	30.01	0.0001
RadGraphF1	30.00	0.0001

Table 6: Statistical significance analysis using one-way ANOVA for BERTScore, BioClinicalBERTScore, and RadGraphF1 scores across four evaluation settings: Rad-Flamingo, Rad-Flamingo w/oI, Open-Flamingo, and Open-Flamingo w/oI. The results indicate significant differences in scores, as determined by F -statistics and p -values ($p < 0.05$).

and Open-Flamingo w/oI.

Null Hypothesis (H_0): There is no significant difference between the <score-name>. **Alternative Hypothesis (H_1):** There is significant difference between the <score-name>. As each of the output from the models are mean of generated reports over the chosen disease classes, we take them as the group mean for the one-way ANOVA test (Ross and Willson, 2017). Therefore, we consider the four evaluation setting as four groups of data, We get F -statistic = 30.00 and p -value ≈ 0.0001 respectively. Consequently, F -statistic $> F_{critical}$ and p -value < 0.05 , satisfying these conditions we can reject the Null Hypothesis thereby establishing the values are significantly different. Similarly, we get F -statistic = 30.01 and p -value ≈ 0.0001 respectively. As the BioClinicalBERTScores are similar to the BERTScore we get similar F -statistic and p -value. Consequently, F -statistic $> F_{critical}$ and p -value < 0.05 , satisfying these conditions we can reject the Null Hypothesis thereby establishing the values are significantly different. Lastly, we get F -statistic = 30.00 and p -value ≈ 0.0001 respectively. Consequently, F -statistic $> F_{critical}$ and p -value < 0.05 , satisfying these conditions we can reject the Null Hypothesis thereby establishing the values are significantly different.