# BIDIRECTIONAL LEARNING FOR THE VISUAL REP RESENTATION IN RADIOLOGY REPORT GENERATION WITH FROZEN LLMS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Radiology report generation (R2Gen) has recently leveraged large language models (LLMs), achieving improved results. However, the generated reports still fall short in both language accuracy and clinical relevance. A key challenge is learning a visual representation of radiology images that an LLM can effectively interpret. To address this, we propose that for a visual representation to be interpretable by an LLM, it shall also be generatable by the LLM. Building on this idea, we introduce a novel bidirectional learning framework for R2Gen, integrating both visionto-text and text-to-vision information to enhance visual representation learning. First, we require that the visual representation aid the LLM in generating reports that closely match the ground truth. Second, we require that the visual representation be maximally generated by the LLM when provided with the ground truth report. To enable the frozen LLM to perform text-to-vision generation, we jointly train a new text encoder for reports. Additionally, through an image reconstruction task, we encourage the visual representation to capture the core features of input radiology images. This bidirectional learning framework is realized using a frozen LLM and incurs no extra computational cost at the inference stage. Experimental results demonstrate better alignment between the learned visual representation and the LLM's word embedding space, along with state-of-the-art performance in both language accuracy and clinical efficacy. Our code will be publicly released.

031 032

033

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

### 1 INTRODUCTION

034 Automated radiology report generation (R2Gen) has emerged as a promising solution to alleviate 035 the heavy workload of radiologists while ensuring the consistency and accuracy of medical reports. Significant progress has been made in this area, with numerous methods proposed in the literature 037 Chen et al. (2020; 2022); Liu et al. (2021a); Yang et al. (2022). The recent development of large 038 language models (LLMs) has further advanced this area, enabling the generation of more accurate reports thanks to their exceptional natural language generation (NLG) capabilities Wang et al. (2023b); Yang et al. (2023); Lee et al. (2023). However, despite their state-of-the-art performance, 040 current LLM-based R2Gen methods continue to face key challenges in achieving the level of lin-041 guistic accuracy and clinical efficacy necessary for real-world medical diagnosis. 042

Among the various challenges, a key issue in LLM-based R2Gen methods is improving the effectiveness of the visual representation that is fed into the LLM to generate reports. Since this representation provides all the information that the LLM can have about the input radiology image (e.g., chest X-ray images), its effectiveness is crucial. Compared with existing approaches that incorporate additional or external information (e.g., lesion-aware augmentations Hou et al. (2024), structured clinical information Dalla Serra et al. (2022), or knowledge bases Huang et al. (2023); Ranjit et al. (2023); Jin et al. (2024a); Bu et al. (2024)) to address this issue, this work takes a different approach that fully exploits the available image-report training pairs to enhance visual representation learning, without introducing any supplementary information.

In this work, we focus on enhancing the compatibility of the visual representation with the LLM, ensuring it is more easily "understood" by the LLM. This challenge, commonly noted in the literature, arises from the well-known modality gap between vision (captured by the visual representation) and

text (the modality on which the LLM is trained). A typical remedy to this case is to use a visual mapper to project visual features onto the space where the LLM's word embedding space resides. However, this remedy alone often falls short in fully bridging the gap between visual and textual modalities, limiting the LLM's ability to fully understand the visual input. To further address this gap, CLIP-based contrastive learning is a natural extension to align visual embeddings with text embeddings more closely. Differently, in this work, we propose a more radical approach.

Our work is inspired by Richard Feynman's famous quote: "*What I cannot create, I do not understand.*" We propose that for a visual representation to be truly "understood" by an LLM, it shall be generatable by the LLM. This requirement is used to achieve a stronger compatibility between the visual representation and the LLM's word embedding space. Furthermore, the generation of the visual representation by the LLM should not be done in a random or arbitrary manner. Instead, considering the correspondence between each image and its ground truth report in the training data, we require that the visual representation be generated when the report is used as input to the LLM.

067 Building on this idea, we propose a novel bidirectional learning framework to learn visual repre-068 sentations for radiology report generation. It is fully built around a frozen LLM, ensuring computa-069 tional efficiency and avoiding performance degradation from improper fine-tuning. This framework 070 involves both vision-to-text and text-to-vision tasks. For the vision-to-text task, we require the vi-071 sual representation to enable the LLM to generate a report that closely matches the ground truth. In an innovative twist, the text-to-vision task requires the visual representation to be maximally 072 generatable by the LLM when provided with the ground truth report. Considering that the LLM, 073 having been trained primarily on textual data, may not effectively handle a text-to-vision task with 074 its built-in word embedding, we jointly learn a new text encoder between the ground truth report 075 and the frozen LLM. Lastly, to enhance the visual representation further, we require it to support 076 the reconstruction of the input radiology image, enabling it to capture the core characteristics of the 077 underlying distribution of these radiology images.

To further understand the core principles of this visual representation learning, we analyze it through the lens of model regularization. From this perspective, the text-to-vision task functions as a regularization constraint, preventing the visual encoder from becoming overly complex and enhancing its generalization capabilities. Once trained, the framework can generate a report when a radiology image is provided. Importantly, at the inference stage, only the vision-to-text branch is needed, operating in the same way as current LLM-based R2Gen methods. This ensures that our method does not introduce any additional computational overhead when deployed.

Extensive experiments on the IU-Xray and MIMIC-CXR datasets demonstrate that the proposed framework consistently outperforms existing methods, achieving the state-of-the-art results in both language accuracy and clinical efficacy. Additionally, ablation studies validate the contribution of the proposed bidirectional learning process, confirming the improved compatibility between the learned visual representation and the LLM's word embedding space.

090 091

# 2 RELATED WORK

092 093 094

Radiology Report Generation. Radiology report generation (R2Gen) has traditionally focused on vision-to-text architectures, where a vision encoder extracts features from CXR scans and a text decoder produces the corresponding report. Early work in this domain used CNNs and LSTMs Jing et al. (2017); Wang et al. (2018); Xue et al. (2018), followed by transformer-based models that significantly improved report quality Chen et al. (2020; 2022); Nicolson et al. (2023); Wang et al. (2023b); Lee et al. (2023). Despite these advancements, achieving language accuracy and clinical efficacy remains a persistent challenge.

To improve report generation quality, several approaches have been explored. Some methods focus on detecting and describing key anatomical regions Tanida et al. (2023); Dalla Serra et al. (2023), while others emphasize feature alignment between visual and textual modalities Wang et al. (2022; 2023a); Li et al. (2023). Additionally, external knowledge sources, such as lesion-aware augmentations Hou et al. (2024), structured clinical information Dalla Serra et al. (2022), symptom graphs, or knowledge distiller Liu et al. (2021b); Huang et al. (2023); Ranjit et al. (2023); Jin et al. (2024a); Bu et al. (2024), have shown promise. However, constructing specialized knowledge bases for such integration remains resource-intensive and demands domain expertise. Most of aforementioned methods follow a unidirectional vision-to-text approach, focusing on visual feature extraction while not adequately utilizing the semantic richness embedded in radiology reports. In contrast, our approach fully leverages report semantics without introducing external knowledge or supplementary data.

112

The recent work MedM2G Zhan et al. (2024) aims to unify multiple cross-modality medical generation tasks (e.g., among text, X-ray, CT, and MRI) into a single framework Rombach et al. (2022). Our work differs from it in several ways. First, MedM2G focuses on unifying cross-modality generation tasks across different medical imaging modalities, while we concentrate primarily on improving radiology report generation. Second, MedM2G uses the integration of multiple diffusion models, whereas our framework is based on the LLM-based R2Gen setting. Third, MedM2G does not explicitly explore bidirectional learning for aligning images and reports. At last, as will be shown in the experiment, our method achieves better performance on radiology report generation benchmarks.

120 LLM-based R2Gen Models. Recent large language models (LLMs) have demonstrated impressive 121 abilities in text generation across various fields, offering new potential for improving R2Gen sys-122 tems. For instance, RAG Ranjit et al. (2023) treats R2Gen as a retrieval task, utilizing the GPT-3.5-123 turbo and GPT-4 to generate reports based on a retrieval corpus. R2GenGPT Wang et al. (2023b) 124 uses a frozen LLM to generate reports and map visual features to the LLM built-in embedding space. It improves language fluency, yet clinical accuracy remains a key challenge. Other works, 125 like MedXChat Yang et al. (2023) and LLM-CXR Lee et al. (2023), fine-tune LLMs to consolidate 126 multiple tasks, such as image-to-report generation and visual question answering (VQA) into one 127 unified framework. While MedXChat incorporates instruction tuning and fine-tunes Stable Diffu-128 sion Rombach et al. (2022) for CXR report generation and image synthesis, LLM-CXR focuses on 129 image understanding and generation by tokenizing chest X-ray images through VQ-GAN Esser et al. 130 (2021) and using instruction-finetuning to perform tasks such as CXR generation and CXR-related 131 VQA. Both methods leverage the instruction-following capabilities of LLMs to connect text and 132 image modalities in medical applications. 133

Our approach differs from MedXChat and LLM-CXR in several ways: (1) Instead of fine-tuning LLMs, we keep the LLM frozen and focus on improving visual representations without modifying the LLM; (2) As will be shown, our method achieves better performance than MedXChat in radiology report generation. As for LLM-CXR, it concentrates on image synthesis and VQA and does not provide strictly comparable results for radiology report generation; (3) Different from their methods where vision-to-text and text-to-vision tasks operate independently, our bidirectional learning approach ensures that text-to-vision generation enhances the visual encoder's effectiveness, leading to improved report quality.



Figure 1: Overview of the proposed bidirectional learning framework for radiology report generation. It consists of three components: 1) the vision-to-text branch; 2) the text-to-vision branch; 3) the image reconstruction component. To keep the diagram concise, certain detailed components, such as the visual and textual mappers, have been omitted. Details are provided in Section 3.

#### 3 Methodology

162 163 164

166

167

168

**Framework Overview.** As illustrated in Figure 1, the proposed framework consists of three components, corresponding to the vision-to-text task (at the bottom part of the figure), the text-to-vision task (at the top-right part), and the image reconstruction task (at the top-left corner). At the algorithm level, the central issue is to optimise the network parameters of the vision encoder. It is conducted by minimising a combination of three loss terms. They are presented in order as follows.



Figure 2: The format of the input to the LLM for the vision-to-text (a) and text-to-vision (b) tasks.
In both cases, the embedding layer and tokenizer of the frozen LLM (i.e., the parts in dark green) are shared between the branches, while trainable encoders and mappers are specific to each modality. The vision encoder and visual mapper, along with the extracted vision embeddings, are shared between the two branches (indicated in purple). Note that "V-Mapper" and "T-Mapper" are short for visual and textual mappers, respectively.

183 **Vision-to-Text Branch.** As shown in Figure 2(a), given an input image *I*, the vision encoder produces its visual features  $\mathbf{Z}_v = g_v(I; \boldsymbol{\theta}_v)$ , where  $\mathbf{Z} \in \mathbb{R}^{N_p \times d_v}$ , with  $N_p$  being the number of patches 185 and  $d_v$  the dimensionality of the visual features. Then a visual mapper is used to project the visual 186 features onto the space where the LLM's word embedding space resides. That is,  $\mathbf{E}_v = m_v(\mathbf{Z}_v)$ , 187 where  $\mathbf{E}_v \in \mathbb{R}^{N_p \times d_L}$ ,  $d_L$  denotes the dimensionality of the LLM's word embedding space, and 188  $m_v(\cdot)$  represents the trainable visual mapper. The visual embeddings obtained in  $\mathbf{E}_v$  are fed into the 189 frozen LLM to generate a report, instructed by a prompt denoted by  $S_{n2t}$ , as illustrated in Figure 190 3(a). The generated report is compared with the ground-truth report through a vision-text consistency (VTC) loss, defined as the auto-regressive negative log-likelihood of generating the correct 191 report tokens. It is expressed as 192

193

194

196

$$\mathcal{L}_{\text{VTC}} = -\sum_{k=1}^{l} \log P(t_k \mid \mathbf{E}_v, \mathbf{S}_{v2t}, T_{< k}), \tag{1}$$

where l is the length of a report and  $T_{<k}$  represents the ground truth tokens preceding the token  $t_k$  to be generated. Through error back-propagation, this loss encourages the vision encoder to produce visual features that best support the LLM to generate high-quality reports.

**Text-to-Vision Branch.** The text-to-vision branch works with the vision-to-text branch to form the proposed bidirectional learning process. It further tightens the connection between the extracted visual embeddings  $\mathbf{E}_v$  and the ground truth reports. This process is powered by the LLM, which plays a central role in extracting the rich semantic information in the reports and converting it into visual representations. Nevertheless, two subtle issues needs to be addressed beforehand.

One issue is that radiology reports often omit whether the CXR image is frontal or lateral, leading to potential mismatches when generating visual embeddings. To avoid this, we train a classifier using only CXR images to identify the view (frontal or lateral) and append this information to the report. This ensures that the LLM generates visual embeddings aligned with the correct view. The classifier is trained strictly on the training set images and serves as a preprocessing step, without influencing the model's training or testing.

The other issue, which is more important, is that we need to jointly learn a new text encoder between the ground truth report and the frozen LLM, instead of directly using the LLM's built-in word embedding layer. This is due to several considerations: 1) the LLM, having been trained primarily on textual data, may not effectively handle a text-to-vision task with its built-in word embedding layer; 2) the text encoder can be optimized to best support the LLM to generate the targeted visual embeddings; and 3) this design ensures the LLM remains frozen, preserving its original capabilities. 216 As shown in Figure 2(b), the ground truth report T, now containing the view information, is fed into 217 a text encoder to extract textual features  $\mathbf{Z}_t = g_t(T; \boldsymbol{\theta}_t)$ , where  $\mathbf{Z}_t \in \mathbb{R}^{l \times d_t}$  with *l* being the length of 218 a report and  $d_t$  the dimensionality of the textual features. As in the vision-to-text branch, a trainable 219 textual mapper  $m_t(\cdot)$  is used to project the textual features onto the space where the LLM's word 220 embedding space resides, and this produces the text embeddings  $\mathbf{E}_t = m_t(\mathbf{Z}_t)$ . The embeddings are then fed into the frozen LLM to generate the visual embeddings denoted by  $\mathbf{E}_g \in \mathbb{R}^{N_p \times d_L}$ . Recall 221 that  $N_p$  is the number of image patches determined by the vision encoder  $q_v(\cdot)$ . Note that the LLM 222 is instructed by a text-to-vision prompt, in which the instruction is  $S_{t2v}$  and the visual embeddings (extracted by the vision encoder and projected by the visual mapper) in  $\mathbf{E}_{v}$  are provided as the 224 target value, as indicated by the purple segment in Figure 2(b). This requires the learning process to 225 optimise the new text encoder to enable the LLM to generate target visual embeddings. Formally, 226 this is conducted via the text-vision consistency (TVC) loss as follows. 227

 $\mathcal{L}_{\text{TVC}} = \frac{1}{N_p} \sum_{i=1}^{N_p} \|\mathbf{E}_{v,i} - \mathbf{E}_{g,i}\|^2.$ 

229

230 231

246 247

253 254

262

264

where the subscript "i" indicates the *i*th vector in the embedding matrix **E**.

232 **Image Reconstruction.** This component requires the visual embeddings,  $E_v$ , to support the re-233 construction of the input radiology image. We firstly use a trainable mapper  $m'_{u}(\cdot)$  to project  $\mathbf{E}_{v}$ 234 from the LLM's word embedding space back to the output space of the vision encoder, that is, 235  $\hat{\mathbf{Z}}_v = m'_v(\mathbf{E}_v)$ , where  $\hat{\mathbf{Z}}_v \in \mathbb{R}^{N_p \times d_v}$ . After that,  $\hat{\mathbf{Z}}_v$  is fed into the vision decoder  $g_{vd}$  to re-236 construct the image as  $\hat{I} = g_{vd}(\hat{\mathbf{Z}}_v; \boldsymbol{\theta}_{vd})$ , where  $\boldsymbol{\theta}_{vd}$  represents the trainable network parameters 237 of  $g_{vd}$ . The reconstructed error is measured by the discrepancies between corresponding pixels as 238  $\mathcal{L}_{\text{REC}} = \frac{1}{CHW} \sum_{c=1}^{C} \sum_{i=1}^{H} \sum_{j=1}^{W} ||I_{c,i,j} - \hat{I}_{c,i,j}||^2$ , where H, W, and C are the height, width, and 239 channel numbers of the input CXR image, and  $I_{c,i,j}$  and  $\hat{I}_{c,i,j}$  represent the pixel values at channel 240 c and position (i, j) of the original and reconstructed images. Minimizing the reconstruction loss 241 encourages the vision encoder to extract the visual representation that can characterise the essential 242 features of the training images. 243

Loss Function and Interpretation. For a given training sample in the form of an image-report pair (I,T), its total loss function is

$$\mathcal{L}_{\text{TOTAL}}(I,T) = \mathcal{L}_{\text{VTC}} + \lambda_1 \mathcal{L}_{\text{TVC}} + \lambda_2 \mathcal{L}_{\text{REC}},\tag{3}$$

(2)

where  $\lambda_1$  and  $\lambda_2$  are hyperparameters that balance the contributions of each loss term.

To interpret this loss, we specially highlight the visual embedding  $\mathbf{E}_v$  and its dependence on the network parameter  $\boldsymbol{\theta}_v$  of the vision encoder, and we temporarily omitting the parameters of the visual and textual mappers. The total loss is then rewritten as

$$\mathcal{L}_{\text{TOTAL}}(I,T) = \underbrace{\mathcal{L}_{\text{VTC}}(\mathbf{E}_{v}(\boldsymbol{\theta}_{v}))}_{\text{Error term}} + \underbrace{\lambda_{1}\mathcal{L}_{\text{TVC}}(\mathbf{E}_{v}(\boldsymbol{\theta}_{v}),\boldsymbol{\theta}_{t}) + \lambda_{2}\mathcal{L}_{\text{REC}}(\mathbf{E}_{v}(\boldsymbol{\theta}_{v}),\boldsymbol{\theta}_{vd})}_{\text{Regularization terms}}.$$
 (4)

As indicated, we interpret the first term on the right hand side as an "error term" because it corresponds to the training error measured by the negative log-likelihood of generating the correct report tokens, which is essentially a cross-entropy classification loss. We interpret the following two terms as "regularization terms" by considering that they impose penalties on the objective function (i.e., the total loss) to constrain the space where  $\mathbf{E}_v$  (or more fundamentally, the vision encoder's parameter  $\theta_v$ ) can reside. Now we focus on the term  $\mathcal{L}_{\text{TVC}}$ , which is the key contribution of this work, to further interpret its regularization effect by rewriting Eq.(2) in a compact form as

$$\mathcal{L}_{\text{TVC}}(\mathbf{E}_{v}(\boldsymbol{\theta}_{v}), \boldsymbol{\theta}_{t}) = \frac{1}{N_{p}} \|\mathbf{E}_{v}(\boldsymbol{\theta}_{v}) - \mathbf{E}_{g}(\boldsymbol{\theta}_{t})\|_{\mathcal{F}}^{2},$$
(5)

where  $\mathcal{F}$  denotes the matrix Frobenius norm. Now it is clear that  $\mathcal{L}_{\text{TVC}}$  essentially imposes a prior on  $\mathbf{E}_v(\boldsymbol{\theta}_v)$ , which is  $\mathbf{E}_g(\boldsymbol{\theta}_t)$ , and measures its deviation from this prior. More interestingly, this prior is not a predefined, static one but optimised to adaptively vary with the ground truth report. Geometrically,  $\mathbf{E}_v(\boldsymbol{\theta}_v)$  is restricted within a high-dimensional sphere centered at  $\mathbf{E}_g(\boldsymbol{\theta}_t)$  while trying to move towards the low-value regions of the error term  $\mathcal{L}_{\text{VTC}}$ . During optimising the total loss,  $\mathcal{L}_{\text{TVC}}$  negotiates with  $\mathcal{L}_{\text{TVC}}$  to find the optimal  $\boldsymbol{\theta}_v$  (and  $\boldsymbol{\theta}_t$ ). 270 In addition to the above regularization perspective, we could also interpret  $\mathcal{L}_{TVC}$  in Eq.(5) as another 271 type of alignment between an image, represented by  $\mathbf{E}_v(\boldsymbol{\theta}_v)$ , and its ground truth report, whose 272 semantic information is conveyed by  $\mathbf{E}_{q}(\boldsymbol{\theta}_{t})$ . Such alignment takes place in the space where the 273 LLM's word embedding space resides and the alignment degree is measured by the distance be-274 tween them. In sum, the introduction of  $\mathcal{L}_{VTC}$  and  $\mathcal{L}_{REC}$ , either interpreted from the perspective of regularization or alignment, further constrains the feasible domain of the vision encoder's parameter 275  $\theta_v$ . According to regularization theory Haykin (2007) (Chapter 7), this helps to reduce the model 276 complexity of the vision encoder and improve its generalization capability. 277

278 279

280

- 4 EXPERIMENTS
- 4.1 DATASETS AND SETTINGS282

IU-Xray. Indiana University Chest X-ray Collection (IU-Xray) Demner-Fushman et al. (2016) contains 3,955 de-identified radiology reports, each of which is associated with frontal and/or lateral chest X-ray images, and 7,470 chest X-ray images in total. Each report is comprised of several sections: Impression, Findings, Indication, etc. In this work, we adopt the same data set partitioning as in the literature Chen et al. (2020) for a fair comparison, with a training/validation/test split set by 7:1:2 of the entire dataset. All evaluations are done on the test set.

MIMIC-CXR. This dataset includes 377,110 chest X-ray images from 65,379 patients, each accompanied by a radiology report. We adhere to the training/validation/test split in the literature Chen et al. (2020) to facilitate comparison with related methods. This split results in 270,790 images for training, 2,130 for validation, and 3,858 for testing, all paired with their respective reports. Keeping the aspect ratio, all images are resized with the shorter edge to be 256 pixels.

Evaluation Metrics. We utilize standard natural language generation (NLG) metrics for assessment, 294 including BLEU scores Papineni et al. (2002), ROUGE-L Chin-Yew (2004), and METEOR Banerjee 295 & Lavie (2005). Specifically, to evaluate the model, we use the BLEU-4 score on the validation set 296 to select the optimal model checkpoint. In addition to the NLG metrics, we employ several clinical 297 efficacy (CE) metrics to assess the model's ability to generate clinically accurate reports. These 298 include the RadCliQ metric Yu et al. (2023), which combines multiple individual metrics to align 299 with radiologist evaluations. We also utilize the Bert Score Zhang et al. (2019), which measures the 300 semantic similarity between generated reports and ground-truth reports by comparing contextualized 301 embeddings, and the RadGraph F1 score Jain et al. (2021), which evaluates the model's ability to 302 correctly extract and describe clinical entities and their relationships. These CE metrics provide a 303 comprehensive assessment of clinical accuracy.

304 **Implementation Details.** The Llama2-7B model<sup>1</sup> is utilised as text decoder, and the text encoder of 305 the base version of the BLIP<sup>2</sup> is used as the new text encoder. The large version of the MAE<sup>3</sup> forms 306 the visual part, i.e., visual encoder and decoder. Between the LLM and the visual text encoders, there 307 is a fully connected layer that can be used to project visual text features to match the dimensions 308 of the LLM's embedding space. Training is conducted on four NVIDIA A100 80GB GPUs for 5 309 epochs on MIMIC-CXR and 2 epochs on IU-Xray, with a mini-batch size of 4, a learning rate of 310 1e-4, AdamW optimizer Loshchilov & Hutter (2018), and a cosine annealing scheduler. Testing 311 involves beam search with a size of 3. Images are randomly cropped to  $224 \times 224$  during training and inference. Both  $\lambda_1$  and  $\lambda_2$  in equation 3 are set as 1 without fine-tuning. 312

313 314

4.2 **RESULTS AND DISCUSSION** 

Table 1 compares the proposed method and relevant state-of-the-art ones on the IU-Xray and MIMIC-CXR datasets. This comparison includes image captioning methods such as Show-Tell Vinyals et al. (2015), AdaAtt Xu et al. (2015), and M2Transformer Cornia et al. (2020), medical report generation methods like R2Gen Chen et al. (2020), R2GenCMN Chen et al. (2022), PP-KED Liu et al. (2021a), GSK Yang et al. (2022), MSAT Wang et al. (2022), METransformer Wang et al. (2023a), CvT2DistilGPT2 Nicolson et al. (2023), KGEER Dalla Serra et al. (2022), D<sup>2</sup>-Net Jin

<sup>2</sup>https://huggingface.co/Salesforce/blip-image-captioning-base <sup>3</sup>https://huggingface.co/facebook/vit-mae-large

<sup>322 &</sup>lt;sup>1</sup>https://huggingface.co/meta-llama/Llama-2-7b-chat-hf

#### 324

Table 1: Comparison on IU-Xray (upper part) and MIMIC-CXR datasets (lower part). † indicates 325 the results are quoted from their respective papers. Those without † are obtained by re-running the 326 publicly released codebase Li et al. (2021) using the same training-test partition as our method. The 327 highest and second highest performance are highlighted by bolding and underlining respectively. 328

320 0		0 1 1		0 0		0		0 1		
329	Dataset	Methods	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	METEOR		
330		Show-Tell	0.243	0.130	0.108	0.078	0.307	0.157		
221		AdaAtt	0.284	0.207	0.150	0.126	0.311	0.165		
331		M2Transformer		- 0.284	0.168	- 0.143	0.328 _	0.170		
332		R2Gen'	0.470	0.304	0.219	0.165	0.371	0.187		
333		MSAT	0.475	0.309	0.222	0.170	0.375	0.191		
334		CvT2DistilGPT2 <sup>†</sup>	0.473	0.304	0.220	0.175	0.372	0.200		
335	IU-Xray	$\overline{R}_{2}\overline{G}en\overline{GP}T^{\dagger}$	<u>0.488</u>	0.316	0.228	- 0.173	0.377	<u>0.211</u>		
336		Ours	0.512	0.341	0.249	0.186	0.392	0.221		
337		Results below are not strictly comparable due to different data partition. For reference only.								
338		<b>PPKED</b> <sup>†</sup>	0.483	0.315	0.224	0.168	0.376	0.187		
339		METransformer <sup>†</sup>	0.483	0.322	0.228	0.172	0.380	0.192		
240		$D^2$ -Net <sup>†</sup>	0.492	0.327	0.231	0.171	0.378	0.204		
340		EKAGen <sup>†</sup>	0.526	0.361	0.267	0.203	0.404	0.214		
341		MedM2G <sup>†</sup>	0.533	0.369	0.278	0.212	0.416	-		
342		Show-Tell	0.308	0.190	0.125	0.088	0.256	0.122		
343		AdaAtt	0.314	0.198	0.132	0.094	0.267	0.128		
344		M2Transformer	$- \frac{0.332}{2.52}$ -	$-\frac{0.210}{0.210}$	$ \frac{0.142}{0.142} -$	$-\frac{0.101}{0.102}$	0.264	0.134		
045		R2Gen'	0.353	0.218	0.145	0.103	0.277	0.142		
345		R2GenUMIN	0.353	0.218	0.148	0.106	0.278	0.142		
346		GSKT	0.360	0.224	0.149	0.100	0.284	0.149		
347		MSAT <sup>†</sup>	0.303	0.220	0.150	0.110	0.282	0 143		
348	MIMIC-CXR	METransformer <sup>†</sup>	0.386	0.250	0.169	0.120	0.291	0.152		
3/10	Minite entr	UniXGen-256 <sup>†</sup>	0.365	0.227	0.147	0.101	0.294	0.156		
343		CvT2DistilGPT2 <sup>†</sup>	0.393	0.248	0.171	0.127	-	0.155		
350		KGEER <sup>†</sup>	0.363	0.245	0.178	0.136	<u>0.313</u>	0.161		
351		$D^2$ -Net	0.365	0.230	0.153	0.107	0.278	0.136		
352		EKAGen <sup>†</sup>	0.419	0.258	0.170	0.119	0.287	0.157		
252		- MedM2G'	0.412 -	- 0.260	0.179	$-\frac{0.142}{0.000}$ -	0.309			
333		LLM-CXR <sup>†</sup>	0.196	0.095	0.054	0.033	0.245	0.081		
354		P2GanGPT <sup>†</sup>	0.307	0.255	0.138	0.111	0.204	0.155		
355		K2GellGF1	0.411	0.207	0.180	0.134	0.297	0.100		
356		Ours	0.427	0.285	0.202	0.144	0.314	0.171		

357

358

359 et al. (2024b), EKAGen Bu et al. (2024), and MedM2G Zhan et al. (2024), as well as LLM-based 360 R2Gen methods LLM-CXR Lee et al. (2023), MedXChat Yang et al. (2023), and R2GenGPT Wang 361 et al. (2023b). Since the IU-Xray dataset lacks an official training-test partition, results for some 362 methods (PPKED, METransformer, D<sup>2</sup>-Net, EKAGen, MedM2G) are not strictly comparable and 363 are provided for reference only. However, all models on MIMIC-CXR use the official partition, ensuring comparability. 364

365 Language Quality Analysis. As demonstrated in Table 1, our framework consistently outperforms 366 prior methods across all key metrics. On IU-Xray, among those strictly comparable methods in 367 Table 1, our method improves the second-best method by BLEU-4 from 0.175 (CvT2DistilGPT2) to 0.186, ROUGE from 0.377 (R2GenGPT) to 0.392, and METEOR from 0.211 (R2GenGPT) to 368 0.221. Among those reference methods (gray part in Table 1), EKAGen and MedM2G show better 369 performance than ours. However, the interpretation of the results requires caution. First, as men-370 tioned, IU-Xray does not provide an official training-test partition, rendering the results not strictly 371 comparable. Second, as IU-Xray is a relatively small dataset, varying training-test partitions can sig-372 nificantly affect performance. Notably, these two methods consistently underperform when strictly 373 compared with our model on the larger MIMIC-CXR dataset, as discussed below. 374

375 More importantly, on the MIMIC-CXR dataset, our method is the best performer, delivering notable improvements. BLEU-1 increases from 0.419 (EKAGen) to 0.427, BLEU-2 from 0.267 376 (R2GenGPT) to 0.285, BLEU-3 from 0.186 (R2GenGPT) to 0.202, and BLEU-4 from 0.142 377 (MedM2G) to 0.144. ROUGE rises from 0.313 (KGEER) to 0.314, and METEOR from 0.161 378 (KGEER) to 0.171, further demonstrating the model's superior contextual accuracy. Notably, on 379 this larger dataset, the two most recent methods EKAGen and MedM2G perform similarly to, if not 380 worse than, R2GenGPT that employs a neat structure to incorporate LLMs for report generation, 381 demonstrating the advantages of LLMs in this task. Our model further advances R2GenGPT by i) 382 integrating the proposed bidirectional learning to enhance the compatibility of visual representation and LLM's comprehension and ii) image reconstruction for more regularization. This synergy leads 383 to the significant improvements in report generation quality observed in both datasets. Focusing 384 on unifying multiple multi-modal analysis tasks, MedXChat and LLM-CXR do not outperform the 385 models specially designed for R2Gen. 386

387 388

Table 2: Evaluation of Clinic-related Metrics on MIMIC-CXR.

hods	RadGraph F1 ( $\uparrow$ )	Bert Score (†)	$RadCliQ\left(\downarrow\right)$
Jen	0.172	0.406	1.228
JenCMN	0.182	0.418	1.182
2DistilGPT2	0.196	0.374	1.220
0ialog-RG <sup>†</sup>	-	0.40	-
JenGPT	0.187	0.415	1.207
's	0.203	0.427	1.169
'S	0.203	0.427	1.10

Clinical Efficacy Analysis. Clinical efficacy scores, like RadGraph F1, Bert Scores, and RadCliQ, 397 are only computable on MIMIC-CXR, shown in Table 2. Our framework achieves the best results across varied clinical efficacy metrics when compared to state-of-the-art methods. Specifically, 398 it improves the RadGraph F1 score, which measures clinical entity extraction accuracy, from 0.196 399 (CvT2DistilGPT2) to 0.203, showcasing the model's superior ability to capture clinical concepts and 400 their relationships. Additionally, Bert Score, which assesses the semantic similarity between gen-401 erated and ground-truth reports, increases from 0.418 (R2GenCMN) to 0.427, indicating enhanced 402 report fluency and clinical relevance. For RadCliQ, where a lower score reflects better clinical qual-403 ity, our framework decreases the score from 1.182 (R2GenCMN) to 1.169, further demonstrating its 404 ability to produce more clinically accurate and error-free reports. These improvements highlight the 405 strength of bidirectional visual representation learning and the image reconstruction in refining the 406 model's clinical efficacy.

407

418

419

420

421

408 409

Table 3: Ablation study of the components in our model ("-" means not applicable).

Dataset	V2T	T2V	Img Rec	Epoch	BLEU-4	ROUGE	METEOR	RadGraph F1	Bert Score	RadCliQ $(\downarrow)$
IU-Xray	~			8	0.176	0.380	0.213	-	-	-
	√(+CLIP)			6	0.178	0.383	0.211	-	-	-
	$\checkmark$	$\checkmark$		5	0.184	0.388	0.219	-	-	-
	✓	$\checkmark$	$\checkmark$	2	0.186	0.392	0.221	-	-	-
MIMIC-CXR	~			9	0.135	0.302	0.163	0.188	0.416	1.201
	√(+CLIP)			7	0.135	0.307	0.162	0.188	0.418	1.197
	$\checkmark$	$\checkmark$		7	0.141	0.315	0.167	0.196	0.425	1.178
	$\checkmark$	$\checkmark$	$\checkmark$	5	0.144	0.314	0.171	0.203	0.427	1.169

**Ablation Study.** As summarized in Table 3, we conduct ablation studies to assess the contribution of each module in our framework including vision-to-text (V2T), text-to-vision (T2V), and image reconstruction (Img Rec), using both NLG metrics and CE metrics. Also, CLIP-based contrastive learning, a widely adopted approach in the literature, is implemented as a reference.<sup>4</sup> Additionally, the number of epochs required for convergence is reported.

As seen, the baseline model with only V2T achieves moderate performance, capturing basic patterns but lacking clinical detail extraction. Adding CLIP-based contrastive learning (V2T+CLIP) seems to have minimal impact on both IU-Xray and MIMIC-CXR in terms of NLG metrics. It slightly improves CE metrics on MIMIC-CXR, e.g., Bert Score increasing from 0.416 to 0.418, and RadCliQ decreasing from 1.201 to 1.197. These results suggest the complexity of visual representation learning for LLM integration, where a CLIP-based feature alignment is not sufficient. Introducing T2V yields stronger improvements, particularly in clinical efficacy, e.g., on MIMIC-CXR, RadGraph F1

429

 <sup>&</sup>lt;sup>4</sup>Specifically, the visual embeddings extracted by the vision encoder from the input CXR image and the
 LLM-based word embeddings of the corresponding ground truth report are projected into a 512-dimensional space via their respective linear layers. The CLIP loss is then applied to bring these embeddings close.

Table 4: Evaluation of our T2V module using R2GenGPT's backbone on IU-Xray dataset

Methods	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	METEOR
R2GenGPT	0.488	0.316	0.228	0.173	0.377	0.211
R2GenGPT+T2V	0.509	0.336	0.246	0.188	0.403	0.218

436 437 438

457

458

461

462

463

464

465

467

468

469 470

471

472

473

474

475

476

477

478

479

480

432

433 434 435

increasing from 0.188 to 0.196 and RadCliQ decreasing from 1.201 to 1.178, showing the effective-439 ness of refining visual embeddings by the embedding generated by the LLM upon the ground turth 440 report. If cross-referencing Figure. 4, introducing T2V better reduces the distances between the vi-441 sual embedding and the LLM's word embedding space than V2T alone. The full model, combining 442 V2T, T2V, and Img Rec, achieves the best results, converging in just 5 epochs on MIMIC-CXR 443 and 2 epochs on IU-Xray. In Figure 4, this setup achieves the shortest distances between the visual 444 embedding and the LLM's embedding space. In sum, T2V plays a key role in moving visual em-445 beddings toward the LLM's operational range, while Img Rec accelerates optimization and further 446 refines the visual representation.

447 Backbone. In addition, to verify the effectiveness of our proposed bidirectional visual representa-448 tion learning, we embed our T2V module into the backbone of R2GenGPT, which employs Swin-449 Transformer as the visual encoder rather than the MAE used in our model. The performance is 450 shown in Table 4. As seen, integrating our bidirectional learning strategy can also significantly 451 improve R2GenGPT, with consistent improvements across different metrics. Specifically, BLUE-4 452 increases from 0.173 to 0.188, ROUGE from 0.377 to 0.403, and METEOR from 0.211 to 0.218. 453 These results reinforce the advantages of our bidirectional visual feature learning when integrating LLM for R2Gen on different backbone models. 454



Figure 3: Comparison of the reports generated by the model trained under different settings on MIMIC-CXR. The key medical information are highlighted using different colors (See the text).

Visualisation. Figure. 3 visualizes the generated reports for the same CXR image using different model configurations. Key medical terms are highlighted in different colors to compare clinical accuracy and relevance. The report produced by V2T (baseline) provides basic observations but lacks clarity on the differential diagnosis for the mild retrocardiac opacity. With the addition of the T2V module (V2T+T2V), the generated report improves by clearly stating the differential diagnoses of pneumonia versus atelectasis and confirming that there is no evidence of new or worsening pleural effusion or pneumothorax, while also indicating that the left lower lobe opacity is stable. Finally, the full model (V2T+T2V+Rec) generates a more comprehensive and precise summary by reiterating the differential diagnoses and emphasizing the stability of the retrocardiac opacity, as well as detailing the unchanged status of the endotracheal and nasogastric tubes, and the PICC line, aligning closely with the ground-truth report.

481 To further investigate the characteristics of the learned visual representation, we calculate the mini-482 mal Euclidean distance of each visual embedding from the LLM's built-in word embedding space, and visualize the histograms of such distances obtained from the whole test set of MIMIC-CXR. 483 Specifically, for each CXR image, the visual embeddings extracted by the vision encoder and pro-484 jected by visual mapper are compared with all token embeddings within the LLM's built-in embed-485 ding space. The Euclidean distance of each visual embedding to the closest token is recorded and

9

493

494

495

496

497

498

499 500

501

504 505

506

507

508 509 510

511

the distribution of these minimum distances across all 3858 test samples of MIMIC-CXR is visualized in the histogram in Figure. 4(a). As seen, by the proposed bidirectional learning, our V2T+T2V
method significantly drags the visual embeddings closer to the LLM built-in word embedding space.
This trend becomes more pronounced by further incorporating the image reconstruction component
(V2T+T2V+REC). In contrast, introducing the CLIP-like contrastive loss (V2T+CLIP) does not
substantially reduce the minimum distances. These results are consistent with the performance observations in Table 3. The similar conclusion can be drawn from the IU-Xray dataset in Figure. 4(b).



Figure 4: Histograms of the minimum Euclidean distances of the learned visual embeddings to the LLM's built-in word embedding space when different visual representation learning schemes are used. The results are based on all test samples in MIMIX-CXR and IU-XRay datasets, respectively.

# 5 CONCLUSION, LIMITATION, AND FUTURE WORK

512 As we explore leveraging the powerful capabilities of LLMs to enhance radiology report generation, 513 a key challenge lies in improving the compatibility between the visual representation produced by 514 the vision encoder and the operational scope of the LLM—essentially making the visual data more 515 easily "understood" by the LLM. This work addresses this challenge by integrating both vision-516 to-text and text-to-vision tasks, resulting in a bidirectional learning framework for radiology report 517 generation. The framework is intentionally built on a frozen LLM, preserving computational effi-518 ciency and performance. As demonstrated through experiments, this bidirectional learning approach 519 improves report quality and strengthens the desired compatibility. In a broader sense, by fostering mutual reinforcement of visual and textual embeddings, our work underscores the potential of LLMs 520 as a foundation for cross-modal generation. This paradigm not only advances radiology report gen-521 eration but also opens new pathways for bridging modality gaps in future research. 522

523 Meanwhile, we also observe in our investigation that while LLMs excel at generating coherent 524 and accurate text, the visual embeddings generated by LLMs are insufficient for synthesizing highfidelity images (e.g., the input chest X-ray scans). This difficulty arises from the modality gap, 525 as LLMs are inherently designed for text generation, complicating the translation of embeddings 526 between textual and visual domains. Directly generating high-quality images from frozen LLMs 527 remains challenging. We anticipate that by fine-tuning LLMs or utilizing truly multimodal large 528 models, we can extend our work to achieve fully bidirectional generation, producing both high-529 quality reports and high-fidelity images. This will be an intriguing direction for our future work. 530

Ethical considerations are essential when developing radiology report generation techniques. Current methods are still not sufficiently reliable for practical medical diagnosis. Even as these techniques mature, issues of fairness, transparency, and explainability must be thoroughly addressed.
Furthermore, the development of radiology report generation techniques requires large amounts of medical data, making it crucial to fully respect and protect patient privacy during the collection, curation, and use of benchmark datasets.

Our work builds on publicly available models, techniques, benchmark datasets, software packages,
and programming languages. The framework, method, and implementation settings are elaborated
in the sections of methodology and experimental study. Upon publication, we will release the source
code and model checkpoints, along with detailed instructions to ensure full reproducibility.

540	REFERENCES
541	KLIEKLIGES

<b>34</b> I	
542	Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved
543	correlation with human judgments. In Proceedings of the acl workshop on intrinsic and extrinsic
544	evaluation measures for machine translation and/or summarization, pp. 65–72, 2005.
545	Shenshen Bu, Taiji Li, Yuedong Yang, and Zhiming Dai. Instance-level expert knowledge and ag-
546	gregate discriminative attention for radiology report generation. In Proceedings of the IEEE/CVF
547	Conference on Computer Vision and Pattern Recognition, pp. 14194–14204, 2024.
548	Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via
549	memory-driven transformer. arXiv preprint arXiv:2010.16056, 2020.
550	Zhihang Chan Valing Shan Van Sang and Vieng Wan. Cross model memory naturalis for redial
551 552	ogy report generation. <i>arXiv preprint arXiv:2204.13258</i> , 2022.
553	Lin Chin-Yew. Rouge: A package for automatic evaluation of summaries. In Proceedings of the
554	Workshop on Text Summarization Branches Out, 2004, 2004.
555	Marcalla Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara, Meshed memory trans
556 557 558	former for image captioning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 10578–10587, 2020.
559	Francesco Dalla Serra, William Clackett, Hamish MacKinnon, Chaoyang Wang, Fani Deligianni
560	Jeff Dalton, and Alison O O'Neil. Multimodal generation of radiology reports using knowledge-
561	grounded extraction of entities and relations. In <i>Proceedings of the 2nd Conference of the Asia-</i>
562	Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint
563	Conference on Natural Language Processing (Volume 1: Long Papers), pp. 615–624, 2022.
564	Francesco Dalla Serra, Chaoyang Wang, Fani Deligianni, Jeffrey Dalton, and Alison O O'Neil.
565	Finding-aware anatomical tokens for chest x-ray automated reporting. In <i>International Workshop</i>
566	on Machine Learning in Medical Imaging, pp. 413-423. Springer, 2023.
567	Dina Demner-Fushman Marc D Kohli Marc B Rosenman Sonya E Shooshan Laritza Rodriguez
568	Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiol-
569	ogy examinations for distribution and retrieval. Journal of the American Medical Informatics
570	Association, 23(2):304–310, 2016.
570	Patrick Esser Robin Rombach and Biorn Ommer Taming transformers for high-resolution image
573	synthesis. In Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-
574	<i>tion</i> , pp. 12873–12883, 2021.
575	Simon Havkin Neural Networks: A Comprehensive Foundation (3rd Edition) Prentice-Hall Inc.
576	USA, 2007. ISBN 0131471392.
577	
578	report consistency of radiology report generation via lesion aware mix up sugmentation arXiv
579	preprint arXiv:2402.12844, 2024.
580	
581	Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang. Kiut: Knowledge-injected u-transformer
582	and Pattern Recognition pp. 19809–19818, 2023
583	and Fallern Recognition, pp. 19809–19818, 2025.
595	Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven QH Truong, Du Nguyen Duong, Tan Bui,
586	Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, et al. Radgraph: Extracting
587	chinear churdes and relations from radiology reports. arXiv preprint arXiv:2100.14405, 2021.
588	Haibo Jin, Haoxuan Che, Yi Lin, and Hao Chen. Promptmrg: Diagnosis-driven prompts for medical
589	report generation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38,
590	pp. 2007–2015, 2024a.
591	Yuda Jin, Weidong Chen, Yuanhe Tian, Yan Song, Chenggang Yan, and Zhendong Mao. Improving
592	radiology report generation with d 2-net: When diffusion meets discriminator. In ICASSP 2024-
593	2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2215–2219. IEEE, 2024b.

594 Baoyu Jing, Pengtao Xie, and Eric Xing. On the automatic generation of medical imaging reports. arXiv preprint arXiv:1711.08195, 2017. 596 Suhyeon Lee, Won Jun Kim, Jinho Chang, and Jong Chul Ye. Llm-cxr: Instruction-finetuned llm 597 for cxr image understanding and generation. arXiv preprint arXiv:2305.11490, 2023. 598 Yaowei Li, Bang Yang, Xuxin Cheng, Zhihong Zhu, Hongxiang Li, and Yuexian Zou. Unify, align 600 and refine: Multi-level semantic alignment for radiology report generation. In Proceedings of the 601 IEEE/CVF International Conference on Computer Vision, pp. 2863–2874, 2023. 602 603 Yehao Li, Yingwei Pan, Jingwen Chen, Ting Yao, and Tao Mei. X-modaler: A versatile and highperformance codebase for cross-modal analytics. In Proceedings of the 29th ACM International 604 Conference on Multimedia, pp. 3799-3802, 2021. 605 606 Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. Exploring and distilling posterior and 607 prior knowledge for radiology report generation. In Proceedings of the IEEE/CVF conference on 608 computer vision and pattern recognition, pp. 13753–13762, 2021a. 609 610 Fenglin Liu, Chenyu You, Xian Wu, Shen Ge, Xu Sun, et al. Auto-encoding knowledge graph for 611 unsupervised medical report generation. Advances in Neural Information Processing Systems, 34: 612 16266–16279, 2021b. 613 Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018. 614 615 Aaron Nicolson, Jason Dowling, and Bevan Koopman. Improving chest x-ray report generation by 616 leveraging warm starting. Artificial intelligence in medicine, 144:102633, 2023. 617 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic 618 evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association 619 for Computational Linguistics, pp. 311–318, 2002. 620 621 Mercy Ranjit, Gopinath Ganapathy, Ranjit Manuel, and Tanuja Ganu. Retrieval augmented chest x-622 ray report generation using openai gpt models. In Machine Learning for Healthcare Conference, 623 pp. 650-666. PMLR, 2023. 624 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-625 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-626 ence on computer vision and pattern recognition, pp. 10684–10695, 2022. 627 628 Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. Interactive and explainable 629 region-guided radiology report generation. In Proceedings of the IEEE/CVF Conference on Com-630 puter Vision and Pattern Recognition, pp. 7433–7442, 2023. 631 Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural 632 image caption generator. In Proceedings of the IEEE conference on computer vision and pattern 633 recognition, pp. 3156-3164, 2015. 634 635 Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, and Ronald M Summers. Tienet: Text-image 636 embedding network for common thorax disease classification and reporting in chest x-rays. In 637 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9049–9058, 638 2018. 639 Zhanyu Wang, Mingkang Tang, Lei Wang, Xiu Li, and Luping Zhou. A medical semantic-assisted 640 transformer for radiographic report generation. In International Conference on Medical Image 641 Computing and Computer-Assisted Intervention, pp. 655–664. Springer, 2022. 642 643 Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. Metransformer: Radiology report gen-644 eration by transformer with multiple learnable expert tokens. In Proceedings of the IEEE/CVF 645 Conference on Computer Vision and Pattern Recognition, pp. 11558–11567, 2023a. 646 Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. R2gengpt: Radiology report generation 647

with frozen llms. *Meta-Radiology*, 1(3):100033, 2023b.

- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich
   Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual
   attention. In *International conference on machine learning*, pp. 2048–2057. PMLR, 2015.
- Yuan Xue, Tao Xu, L Rodney Long, Zhiyun Xue, Sameer Antani, George R Thoma, and Xiaolei Huang. Multimodal recurrent model with attention for automated radiology report generation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st Interna-tional Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part I*, pp. 457–466. Springer, 2018.
- Ling Yang, Zhanyu Wang, and Luping Zhou. Medxchat: Bridging cxr modalities with a unified
   multimodal large model. *arXiv preprint arXiv:2312.02233*, 2023.
- Shuxin Yang, Xian Wu, Shen Ge, S Kevin Zhou, and Li Xiao. Knowledge matters: Chest radiology report generation with general and specific knowledge. *Medical image analysis*, 80:102510, 2022.
- Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser
  Ururahy Nunes Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein Abad, Andrew Y Ng,
  et al. Evaluating progress in automatic chest x-ray radiology report generation. *Patterns*, 4(9),
  2023.
  - Chenlu Zhan, Yu Lin, Gaoang Wang, Hongwei Wang, and Jian Wu. Medm2g: Unifying medical multi-modal generation via cross-guided diffusion with visual invariant. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11502–11512, 2024.
  - Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

# A APPENDIX

You may include other additional sections here.