# R2GenGPT: Radiology Report Generation with Frozen LLMs

Zhanyu Wang[a], Lingqiao Liu[b], Lei Wang[c] and Luping Zhou[a,*]

[a]*University of Sydney, New South Wales 2006, Australia*
[b]*University of Adelaide, South Australia 5005, Australia*
[c]*University of Wollongong, New South Wales 2522, Australia*

## ARTICLE INFO

## ABSTRACT

Large Language Models (LLMs) have consistently showcased remarkable generalization capabilities when applied to various language tasks. Nonetheless, harnessing the full potential of LLMs for Radiology Report Generation (R2Gen) still presents a challenge, stemming from the inherent disparity in modality between LLMs and the R2Gen task. To bridge this gap effectively, we propose R2GenGPT, which is a novel solution that aligns visual features with the word embedding space of LLMs using an efficient visual alignment module. This innovative approach empowers the previously static LLM to seamlessly integrate and process image information, marking a step forward in optimizing R2Gen performance. R2GenGPT offers the following benefits. First, it attains state-of-the-art (SOTA) performance by training only the lightweight visual alignment module while freezing all the parameters of LLM. Second, it exhibits high training efficiency, as it requires the training of an exceptionally minimal number of parameters while achieving rapid convergence. By employing delta tuning, our model only trains 5M parameters (which constitute just 0.07% of the total parameter count) to achieve performance close to the SOTA levels. Our code is available at https://github.com/wang-zhanyu/R2GenGPT.

## 1. Introduction

The landscape of radiological imaging data is experiencing exponential growth that far surpasses the availability of trained readers, resulting in a significant and unsustainable surge in radiologists' workloads. This surge in both the volume and complexity of cases places big pressure on radiologists to interpret more studies within increasingly tight timeframes. Consequently, radiologists are faced with extended working hours and a heightened risk of reading fatigue, all of which significantly contribute to diagnostic errors. Notably, the situation is particularly precarious during on-call hours for emergency radiology studies. As a result, the demand for automated radiographic report generation has soared, as it promises to alleviate the burden on radiologists, mitigate diagnostic errors, and expedite the clinical workflow. Automated radiographic report generation (R2Gen) is a complex AI task. It aims to produce a coherent paragraph that captures the observations and findings depicted in a given radiology image. There are different R2Gen approaches based on whether the report generation is structured and whether it is template-based. This paper focuses on unstructured multi-sentence report generation.

Given its critical clinical relevance, the field of medical report generation has been garnering increasing attention. Most methodologies are inspired by image/video captioning and adopt the encoder-decoder paradigm [15, 38, 44, 48, 32, 33], with specific improvements tailored to the unique characteristics of the R2Gen task. In summary, recent works in the R2Gen task mainly aim to tackle two major challenges. The first challenge lies in **long text generation**. Unlike the image captioning task which generates a single sentence description, medical report generation requires detailed and coherent paragraph-long descriptions. This requires the model to have a robust capacity for learning long-range dependencies. To address this, many solutions have been proposed [13, 46, 42, 4, 3]. For instance, some research works [13, 46, 42] have employed hierarchically structured LSTM which first produces topic vectors using a sentence LSTM and then creates a description for each generated topic with a word LSTM. Another type of work R2Gen [4] introduced a memory-driven Transformer that can record key information of the generation process, enhancing the model's ability to produce long texts. The second challenge lies in the **bias in visual and textual data**. Due to an over-representation of normal samples in the training data, the model's learning process was biased towards these samples,

---

*Corresponding author

✉ zhanyu.wang@sydney.edu.au (Z. Wang); lingqiao.liu@adelaide.edu.au (L. Liu); leiw@uow.edu.au (L. Wang); luping.zhou@sydney.edu.au (L. Zhou)

🌐 https://wang-zhanyu.github.io/ (Z. Wang); https://lingqiao-adelaide.github.io/lingqiaoliu.github.io/ (L. Liu); https://sites.google.com/view/lei-hs-wang (L. Wang); https://sites.google.com/view/lupingzhou (L. Zhou)

limiting its ability to effectively detect abnormalities and anomalies within the dataset. Some works [42, 47, 39] have addressed this issue by aligning image and text/report features, such as work Self-boost [42] incorporating an image-text matching branch to enhance the model's capability to capture the anomalous features in the image. Other research works mitigate the effects of data bias by incorporating external knowledge, such as medical tags [13, 41], and knowledge graphs [17, 50, 21, 45]. For instance, PPKED [21] utilizes a knowledge graph and introduces the Posterior-and-Prior Knowledge Exploring-and-Distilling framework. Despite many efforts and solutions putting forth, the aforementioned two challenges remain significant issues in this field.

Recently, large language models (LLMs) (e.g., [5, 35]) have demonstrated excellent capabilities to perform tasks with zero in-domain data, conduct logical reasoning, and apply commonsense knowledge in NLP tasks [16, 43]. This leads us to ponder whether we can apply large language models to medical report generation tasks, as pre-trained large language models seem to inherently possess the ability to address the two challenges mentioned above. As for long text generation, LLMs are equipped with an inherent understanding of grammar, syntax, and semantic coherence, making them well-suited for tasks requiring extended text generation, such as medical reporting. Furthermore, their proficiency in context modeling allows them to maintain consistency and relevance throughout a lengthy report. As for the bias stemming from an over-representation of normal samples in medical datasets, LLMs can serve as potential correctives due to their extensive knowledge base. Having been exposed to vast amounts of data, LLMs demonstrate robustness and are less susceptible to the effects of imbalanced datasets. They are even capable of handling numerous zero-shot tasks. Moreover, current methods mitigating bias entail the incorporation of external knowledge, whereas pre-trained LLMs inherently possess a wealth of informative knowledge.

However, applying LLMs to R2Gen tasks poses challenges due to the fundamental disparity between visual and textual modalities. The crucial step in applying LLMs to R2Gen is to bridge the gap between visual information and textual generation. In this paper, we present R2GenGPT and explore three methods for aligning visual features with large language models. We first process chest x-ray images using a Visual Encoder to obtain visual embeddings. These embeddings are then mapped to the LLM's feature space via a Visual Mapper, ensuring uniform dimensions. To identify the most efficient method of aligning visual features with the LLM, we've crafted three alignment modules: 1) shallow alignment, where only the Visual Mapper is trained and other parameters remain fixed; 2) deep alignment, where both the visual encoder and the Visual Mapper are trained simultaneously; and 3) Delta alignment, where the Visual Mapper and a limited set of incremental parameters from the visual encoder are trained, ensuring both effectiveness and efficiency.

Our main contributions are summarized as follows.

• We propose a novel LLMs-based Radiology report generation (R2Gen) framework, dubbed R2GenGPT. This marks the first instance of harnessing pre-trained large language models (LLMs) for the R2Gen task with comprehensive comparisons conducted on two frequently employed benchmark datasets.

• We explored three methods with varying levels of trainable parameters to connect image modalities to large language models, namely: shallow alignment, delta alignment, and deep alignment, enabling the LLM to effectively process visual information.

• Our approach exhibits promising and robust performance on two widely recognized benchmark datasets—IU-Xray and MIMIC-CXR. In comparison to multiple state-of-the-art methods, our framework consistently demonstrates its efficacy, affirming its potential in the field of R2Gen.

## 2. Relate Works

**Radiology report generation** Radiology report generation (R2Gen) has gained significant attention in recent years, with many models being developed based on the encoder-decoder architecture initially used for image captioning tasks [38, 44, 26]. However, R2Gen poses additional challenges compared to image captioning, as medical reports are typically longer and clinical abnormalities in medical images are harder to detect than natural objects due to the data bias existed in the training set. To address these challenges, researchers have proposed various methods.

In [42], Wang et al. 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. [4] and Wang et al. [41] 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 [17, 50, 21, 45, 18, 9]. Zhang et al. [50] and Liu et al. [21] combined pre-constructed graphs representing relationships between diseases and organs using graph neural networks, enabling more effective feature learning for abnormalities. Li et al. [18] developed a dynamic approach that updates the graph with new knowledge in real-time. Huang et al. [9] incorporated knowledge from a symptom graph into the decoding stage using an injected knowledge distiller.

Apart from knowledge graphs, another method for integrating external knowledge involves incorporating semantic information to assist report generation through multi-task learning, such as multi-label classification [13, 41, 39, 12, 34]. Wang et.al. [41] extracted 768 high-frequency medical terms from RadGraph [11] and trained a multi-label classification network. The prediction results from classification were then incorporated as semantic information input to the decoder, assisting the generation of the report. Jin et al. [12] make use of the diagnostic results from the classification via prompts to explicitly guide the generation process. All the above methods employing an encoder-decoder architecture in a traditional way, typically assign equal importance to both the encoder and the decoder, with a comparable number of trainable parameters. In these methods, the output from the encoder serves as the key and value for cross-attention computation in the decoder. In contrast, our approach based on LLM deviates significantly from this traditional encoder-decoder framework. Firstly, the number of parameters in the decoder significantly exceeds that in the encoder. Secondly, the encoder functions more like a "visual tokenizer", converting images into visual tokens that are fed into LLM. The attention mechanism employed within this framework remains self-attention rather than cross-attention. With this innovative approach, our paper pioneers the use of a "decoder-centric" architecture for the task of medical report generation.

**Large language Models** Recently, there has been a surge of interest in Large Language Models (LLMs) due to their superior efficacy in a wide array of Natural Language Processing (NLP) tasks. This began with transformer models like BERT [31], GPT [28], and T5 [29], each designed with distinct pre-training objectives. The introduction of GPT-3 [2] marked a significant shift, demonstrating the model's impressive zero-shot generalization capabilities, owed to a scaled-up parameter and data volume, which allowed it to excel in tasks it had not previously encountered. This catalyzed the development of several LLMs such as OPT [49], BLOOM [30], PaLM [5], and LLaMA [35], heralding the triumph of LLMs. In a parallel endeavor, Ouyang et al. [25] introduced InstructGPT, which brought human instruction and feedback into alignment with GPT-3. These advancements have been leveraged by applications like ChatGPT, which enables human-like dialogue interactions by responding to a vast spectrum of complex and nuanced questions and instructions.
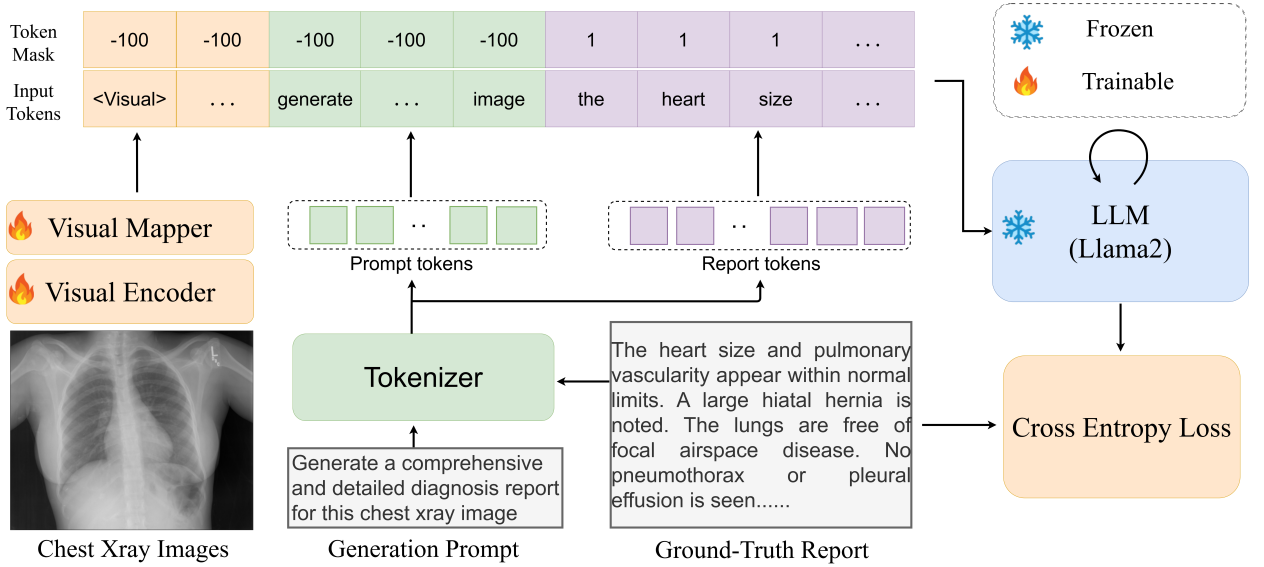
## 3. Methodology

**Overview** As illustrated in Figure 1, R2GenGPT comprises a Visual Encoder, a Visual Mapper, and an LLM (Large Language Model) component. The visual encoder is employed to extract information from chest x-ray images, while the visual mapper serves to project low-dimensional image features into the high-dimensional feature space of the LLM. Utilizing the visual features derived from the chest x-ray images, the LLM generates corresponding diagnostic reports.

**Feature Alignment** For an input chest xray image $\mathbf{X}_v$, we consider the pre-trained Swin Transformer [23] as visual encoder, which provides the visual feature $\mathbf{Z}_v = g(\mathbf{X}_v; \theta_v)$, where $\theta_v$ is the parameters of the Swin Transformer. The grid features of the last transformer layer is utilized in our experiments. We consider a simple linear layer as the Visual Mapper to connect image features into the LLM's word embedding space. Specifically, we apply a trainable projection matrix $\mathbf{W}_m$ to convert $\mathbf{Z}_v$ into language embedding tokens $\mathbf{H}_v$, which have the same dimensionality of the word embedding space in the large language model.

$$\mathbf{H}_v = \mathbf{W}_m \mathbf{Z}_v, \quad \text{with } \mathbf{Z}_v = g(\mathbf{X}_v) \tag{1}$$

Thus we have a sequence of visual tokens $\mathbf{H}_v$. Following the extraction of visual tokens $\mathbf{H}_v$, we propose the following three distinct training strategies to identify the most efficient aligning method by varying the level of trainable parameters.

a) Shallow Alignment: In this mode, we fix the parameters of the pre-trained Swin Transformer and train only the linear Visual Mapper, represented by $\mathbf{W}_m$.

**Figure 1:** An overview of our proposed R2GenGPT. The input tokens for the Large Language Model (LLM) are sequentially concatenated, consisting of visual tokens, prompt tokens, and report tokens. A token mask of -100 indicates that those particular tokens are excluded from auto-regressive training, while a mask of 1 signifies inclusion in auto-regressive training.

b) Deep Alignment: For this approach, both the Swin Transformer and the Visual Mapper are jointly fine-tuned. Specifically, parameters from both the Visual Encoder (Swin Transformer) and the Visual Mapper, denoted as $\theta_v$ and $\mathbf{W}_m$ respectively, are updated.

c) Delta Alignment: As the Swin Transformer utilized in this paper was originally trained on natural images, the shallow alignment approach hinders the model's ability to capture high-quality radiographic image features. On the other hand, adopting deep alignment substantially impacts the model's training efficiency. Therefore, we propose delta alignment, parameter-efficiently fine-tuning the Swin Transformer model using LoRA [8]. Specifically, for a pre-trained weight matrix $\mathbf{W}_0$ within $\theta_v$, LoRA constrains its update with two smaller matrices using a low-rank decomposition $\mathbf{W}_0 + \Delta\mathbf{W}_0 = \mathbf{W}_0 + \mathbf{BA}$, where $\mathbf{W}_0 \in \mathbb{R}^{d \times k}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times k}$, and the rank $r \ll min(d, k)$. It is noted that in our implementation, we only adjust the query and value projections within the swin transformer to prioritize a simple yet efficient model. The trained parameters are denoted as $\Delta\theta_v$, and both $\Delta\theta_v$ and $\mathbf{W}_m$ are trained in this mode.
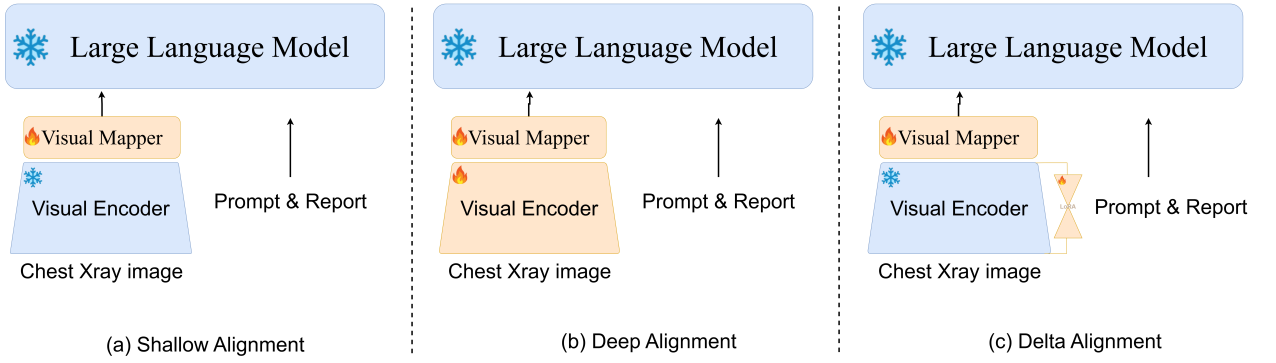
**Large Language Models** We adopt Llama2-7B model for the large language model component. The Llama2-7B stands out for its remarkable capabilities and robustness. Designed with a massive 7-billion-parameter architecture, it encapsulates a rich knowledge base derived from extensive pre-training on diverse datasets. One of its key strengths lies in its extraordinary ability to understand and generate complex language structures, making it particularly well-suited for intricate tasks such as radiology report generation.

Given an chest xray image $\mathbf{X}_v$ and its corresponding report $\mathbf{X}_r$, the detailed prompt inputted into Llama2 is as follows.

*Human: <Img>$\mathbf{X}_v$</Img>, $\mathbf{X}_p$. \n Assistant: $\mathbf{X}_r$ </s>.*

Here $\mathbf{X}_p$ is our designed instruction prompt specific to the R2Gen task. In our current implementation, $\mathbf{X}_p$ = "Generate a comprehensive and detailed diagnosis report for this chest xray image.". For this prompt, before inputting it into LLAMA2 for computation, $\mathbf{X}_v$ will be replaced by visual tokens $\mathbf{H}_v$ processed using Equ. 1 while all other text is tokenized into word tokens using LLAMA's tokenizer.

**Loss Function** We perform instruction-tuning of the LLM only on the report tokens , using its original auto-regressive training objective. Specifically, for a report of length $L$, conditioned on visual information $\mathbf{X}_v$ and instruction prompt

**Figure 2:** Three proposed alignment methods. (a) Shallow Alignment: Training only the Linear Layer. (b) Deep Alignment: Training both the Linear Layer and all parameters of the Visual Encoder. (c) Delta Alignment: Training the Linear Layer and a small subset of incremental parameters of the Visual Encoder.

$\mathbf{X}_p$, our loss function, captured as the negative log likelihood, is formulated as:

$$\mathcal{L}(\theta; \mathbf{X}_r, \mathbf{X}_v, \mathbf{X}_p) = -\sum_{i=1}^{L} \log p_\theta(x_i | \mathbf{X}_v, \mathbf{X}_p, \mathbf{X}_{r,<i}), \tag{2}$$

where $\theta$ is the trainable parameters, $\mathbf{X}_{r,<i}$ is the report tokens before the current prediction token $x_i$.

## 4. Experiments

### 4.1. Data Collection

We evaluated performance using two datasets: a widely-used benchmark IU-Xray [7] and the currently largest dataset MIMIC-CXR [14] for medical report generation.

**IU-Xray**: Indiana University Chest X-ray Collection (IU-Xray) [7] is the most widely used publicly accessible dataset in medical report generation tasks. It contains 3,955 fully 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 [4] for a fair comparison, with a train/test/val set by 7:1:2 of the entire dataset. All evaluations are done on the test set.

**MIMIC-CXR**: This largest publicly available dataset encompasses both chest radiographs and unstructured textual reports. This comprehensive dataset comprises a total of 377,110 chest X-ray images and 227,835 corresponding reports sourced from 64,588 patients who underwent examination at the Beth Israel Deaconess Medical Center between 2011 and 2016. To ensure equitable comparisons, we adhered to MIMIC-CXR's official partitioning as outlined in [4], resulting in 270790 samples designated for training, while allocating 2130 and 3,858 samples for validation and testing, respectively.

### 4.2. Experimental Settings

**Evaluation Metrics**   Adhering to the established evaluation protocol [1], we employ the prevalent metrics for assessment, namely BLEU scores [27], ROUGE-L [20], METEOR [1] and CIDEr [37], to gauge the quality of the generated textual reports. To measure the accuracy of descriptions for clinical abnormalities, we follow [4, 3, 22] and further report clinical efficacy metrics. Specifically, we employ CheXpert [10] for annotating the generated reports, which are subsequently compared against ground truth annotations across 14 distinct categories related to thoracic

---

[1]https://github.com/tylin/coco-caption

**Table 1**
Comparison on IU-Xray (upper part) and MIMIC-CXR datasets (lower part). † indicates the results are quoted from their respective papers. For the methods without †, their results are obtained by re-running the publicly released codebase [19] on these two datasets using the same training-test partition as our method.
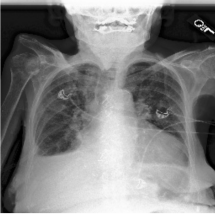
| Dataset | Methods | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 | ROUGE | METEOR | CIDEr |
|---|---|---|---|---|---|---|---|---|
| IU-Xray | Show-Tell | 0.243 | 0.130 | 0.108 | 0.078 | 0.307 | 0.157 | 0.197 |
| | Att2in | 0.248 | 0.134 | 0.116 | 0.091 | 0.309 | 0.162 | 0.215 |
| | AdaAtt | 0.284 | 0.207 | 0.150 | 0.126 | 0.311 | 0.165 | 0.268 |
| | Transformer | 0.372 | 0.251 | 0.147 | 0.136 | 0.317 | 0.168 | 0.310 |
| | M2transformer | 0.402 | 0.284 | 0.168 | 0.143 | 0.328 | 0.170 | 0.332 |
| | R2Gen† | 0.470 | 0.304 | 0.219 | 0.165 | 0.371 | 0.187 | - |
| | R2GenCMN† | 0.475 | 0.309 | 0.222 | 0.170 | 0.375 | 0.191 | - |
| | MSAT† | 0.481 | 0.316 | 0.226 | 0.171 | 0.372 | 0.190 | 0.394 |
| | METransformer† | 0.483 | **0.322** | 0.228 | 0.172 | **0.380** | 0.192 | 0.435 |
| | R2GenGPT (Shallow) | 0.466 | 0.301 | 0.211 | 0.156 | 0.370 | 0.202 | 0.405 |
| | R2GenGPT (Delta) | 0.470 | 0.299 | 0.213 | 0.162 | 0.369 | 0.211 | 0.419 |
| | R2GenGPT (Deep) | **0.488** | 0.316 | **0.228** | **0.173** | 0.377 | **0.211** | **0.438** |
| MIMIC-CXR | Show-Tell | 0.308 | 0.190 | 0.125 | 0.088 | 0.256 | 0.122 | 0.096 |
| | Att2in | 0.314 | 0.198 | 0.133 | 0.095 | 0.264 | 0.122 | 0.106 |
| | AdaAtt | 0.314 | 0.198 | 0.132 | 0.094 | 0.267 | 0.128 | 0.131 |
| | Transformer | 0.316 | 0.199 | 0.140 | 0.092 | 0.267 | 0.129 | 0.134 |
| | M2Transformer | 0.332 | 0.210 | 0.142 | 0.101 | 0.264 | 0.134 | 0.142 |
| | R2Gen† | 0.353 | 0.218 | 0.145 | 0.103 | 0.277 | 0.142 | - |
| | R2GenCMN† | 0.353 | 0.218 | 0.148 | 0.106 | 0.278 | 0.142 | - |
| | PPKED† | 0.36 | 0.224 | 0.149 | 0.106 | 0.284 | 0.149 | 0.237 |
| | GSK† | 0.363 | 0.228 | 0.156 | 0.115 | 0.284 | - | 0.203 |
| | MSAT† | 0.373 | 0.235 | 0.162 | 0.120 | 0.282 | 0.143 | 0.299 |
| | METransformer† | 0.386 | 0.250 | 0.169 | 0.124 | 0.291 | 0.152 | **0.362** |
| | R2GenGPT (Shallow) | 0.365 | 0.237 | 0.163 | 0.117 | 0.277 | 0.136 | 0.145 |
| | R2GenGPT (Delta) | 0.380 | 0.244 | 0.167 | 0.119 | 0.281 | 0.145 | 0.195 |
| | R2GenGPT (Deep) | **0.411** | **0.267** | **0.186** | **0.134** | **0.297** | **0.160** | 0.269 |

diseases and support devices. We use precision, recall, and F1 to evaluate model performance for clinical efficacy metrics.

**Implementation Details** In this work, we leveraged the LLAMA2-7B model [2] as the large language model and the base version of the Swin Transformer [3] as the Visual Encoder. Within the parameters of LoRA, we configured the Lora attention dimension to 16, and the alpha parameter for Lora scaling was also set at 16. The training process was conducted on four NVIDIA A100 40GB GPUs using mixed precision for 3 epochs for MIMIC-CXR and 15 epochs for IU-Xray dataset, with a mini-batch size of 6 and a learning rate of 1e-4. During the testing phase, we employed a beam search strategy with a beam size set to 3.

---

[2] https://huggingface.co/meta-llama/Llama-2-7b-chat-hf
[3] https://huggingface.co/microsoft/swin-base-patch4-window7-224

**Figure 3:** Examples of the generated report on MIMIC-CXR dataset. We compared the results generated by the three alignment methods proposed in R2GenGPT. For better illustration, the key medical information in the reports are highlighted using different colors.

**Table 2**
Evaluation of Model Efficiency and Clinical Efficacy on MIMIC-CXR dataset.

| Models | Trainable Components | | | Scale and Efficiency | | Clinical Efficacy | | |
|---|---|---|---|---|---|---|---|---|
| | Mapper | Encoder | LoRA | Trainable Parameter | Time | Precision | Recall | F1 |
| Shallow | ✓ | | | 4.2M | 1.75h/epo | 0.341 | 0.312 | 0.325 |
| Delta | ✓ | | ✓ | 5.0M | 1.83h/epo | 0.366 | 0.350 | 0.358 |
| Deep | ✓ | ✓ | | 90.9M | 2.75h/epo | **0.392** | **0.387** | **0.389** |
| METransformer | - | - | - | 152M | 3.62h/epo | 0.364 | 0.309 | 0.334 |

## 4.3. Results and Discussion

**Comparison with SOTA** Table 1 showcases a performance comparison between the state-of-the-art methods and our R2GenGPT model variants on the IU-Xray and MIMIC-CXR dataset. In terms of standard image captioning methods, the table considers Show-Tell [38], Att2in [44], AdaAtt [24], Transformer [36], and M2Transformer [6]. Furthermore, medical report generation methods such as R2Gen [4], R2GenCMN [3], MSAT [41], METransformer [40], and other methods marked with † in Table 1 are considered.

From Table 1, it is evident that our R2GenGPT model variants, especially R2GenGPT (Deep), outperform the compared methods across nearly all evaluation metrics. In the MIMIC-CXR dataset, apart from CIDEr, we significantly outperform the latest METransformer [40] method across all metrics. For instance, our BLEU_4 score is improved from 0.124 to 0.134, marking an 8.1% increase. However, we achieved a CIDEr score of 0.269, which is lower than METransformer's 0.362. This discrepancy is because METransformer employs an expert voting strategy similar to an ensemble approach to enhance the CIDEr metric. In comparison to methods without this enhancement, such as R2Gen [4] and PPKED [21], we also hold a distinct advantage in terms of the CIDEr metric. It is also noteworthy that our R2GenGPT (Shallow), with only 4.2M trainable parameters, has been able to achieve a performance in par with the well-known R2Gen model [4], if not even better.

**Model Efficiency and Clinical Efficacy Analysis** In Table 2, we have presented both model efficiency and clinical efficacy metrics. It's evident that our model exhibits higher training efficiency compared to METransformer. For instance, R2GenGPT (Deep) requires training with only 90.9 million parameters, which is significantly less than METransformer's 152 million. Furthermore, R2GenGPT (Delta) achieves comparable performance with just 5 million parameters. To assess the model's training efficiency, we conducted evaluations on four A100 40G GPUs and recorded

the time required for one training epoch. Notably, R2GenGPT Shallow, Delta, and Deep each completed this epoch in just 1.75, 1.83, and 2.75 hours, respectively, compared to METransformer's 3.62 hours, highlighting our model's superior training efficiency. In terms of clinical efficacy metrics, it can be observed that our R2GenGPT(deep) and R2GenGPT(delta) achieved F1 scores of 0.389 and 0.358, respectively, surpassing the current SOTA method, METransformer, with a score of 0.334. This demonstrates the ability of R2GenGPT to generate crucial clinical information.

**Qualitative results** In Figure 3, we compare the reports generated by the three alignment methods of R2GenGPT. To provide a better visualization, we highlight key medical information in both the ground truth and generated reports using different colors. From the figure, it can be observed that reports generated using the Shallow Alignment method are notably inferior to those generated using the Delta Alignment and Deep Alignment methods. For instance, in the first example (top), the Shallow Alignment method erroneously identifies a sample with mild pulmonary edema as a normal sample, whereas Delta Alignment and Deep Alignment methods can accurately identify it.

## 5. Conclusions

In this paper, we present R2GenGPT, an innovative framework at the forefront of Radiology Report Generation (R2Gen) that capitalizes on the capabilities of Large Language Models (LLMs). Through a comprehensive exploration of three alignment methods—shallow, delta, and deep— this research highlights the game-changing potential of LLMs in elevating the R2Gen landscape. R2GenGPT not only attains competitive SOTA performance but also achieves a remarkable reduction in computational complexity. This dual achievement positions R2GenGPT as a promising solution to automate and improve radiology reporting.

## References

[1] Banerjee, S., Lavie, A., 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments, in: ACL.
[2] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners. Advances in neural information processing systems 33, 1877–1901.
[3] Chen, Z., Shen, Y., Song, Y., Wan, X., 2022. Cross-modal memory networks for radiology report generation.
[4] Chen, Z., Song, Y., Chang, T.H., Wan, X., 2020. Generating radiology reports via memory-driven transformer, in: EMNLP.
[5] Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., et al., 2022. Palm: Scaling language modeling with pathways. arXiv preprint arXiv:2204.02311 .
[6] Cornia, M., Stefanini, M., Baraldi, L., Cucchiara, R., 2020. Meshed-memory transformer for image captioning, in: CVPR.
[7] Dina, D.F., Kohli, M.D., Rosenman, M.B., Shooshan, S.E., Laritza, R., Sameer, A., Thoma, G.R., Mcdonald, C.J., 2015. Preparing a collection of radiology examinations for distribution and retrieval. Journal of the American Medical Informatics Association Jamia .
[8] Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., 2022. LoRA: Low-rank adaptation of large language models, in: International Conference on Learning Representations. URL: https://openreview.net/forum?id=nZeVKeeFYf9.
[9] Huang, Z., Zhang, X., Zhang, S., 2023. Kiut: Knowledge-injected u-transformer for radiology report generation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19809–19818.
[10] Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., Seekins, J., Mong, D.A., Halabi, S.S., Sandberg, J.K., Jones, R., Larson, D.B., Langlotz, C.P., Patel, B.N., Lungren, M.P., Ng, A.Y., 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison.
[11] Jain, S., Agrawal, A., Saporta, A., Truong, S.Q., Duong, D.N., Bui, T., Chambon, P., Zhang, Y., Lungren, M.P., Ng, A.Y., et al., 2021. Radgraph: Extracting clinical entities and relations from radiology reports.
[12] Jin, H., Che, H., Lin, Y., Chen, H., 2023. Promptmrg: Diagnosis-driven prompts for medical report generation. arXiv preprint arXiv:2308.12604 .
[13] Jing, B., Xie, P., Xing, E.P., 2018. On the automatic generation of medical imaging reports, in: ACL.
[14] Johnson, A.E.W., Pollard, T.J., Greenbaum, N.R., Lungren, M.P., Deng, C.Y., Peng, Y., Lu, Z., Mark, R.G., Berkowitz, S.J., Horng, S., 2019. Mimic-cxr: A large publicly available database of labeled chest radiographs. CoRR .
[15] Karpathy, A., Li, F., 2015. Deep visual-semantic alignments for generating image descriptions, in: CVPR.
[16] Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y., 2022. Large language models are zero-shot reasoners. Advances in neural information processing systems 35, 22199–22213.
[17] Li, C.Y., Liang, X., Hu, Z., Xing, E.P., 2019. Knowledge-driven encode, retrieve, paraphrase for medical image report generation.
[18] Li, M., Lin, B., Chen, Z., Lin, H., Liang, X., Chang, X., 2023. Dynamic graph enhanced contrastive learning for chest x-ray report generation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3334–3343.
[19] Li, Y., Pan, Y., Chen, J., Yao, T., Mei, T., 2021. X-modaler: A versatile and high-performance codebase for cross-modal analytics, in: Proceedings of the 29th ACM International Conference on Multimedia, pp. 3799–3802.
[20] Lin, C.Y., 2004. ROUGE: A package for automatic evaluation of summaries, in: ACL.

[21] Liu, F., Wu, X., Ge, S., Fan, W., Zou, Y., 2021a. Exploring and distilling posterior and prior knowledge for radiology report generation, in: CVPR.

[22] Liu, F., You, C., Wu, X., Ge, S., Sun, X., et al., 2021b. Auto-encoding knowledge graph for unsupervised medical report generation. Advances in Neural Information Processing Systems .

[23] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B., 2021c. Swin transformer: Hierarchical vision transformer using shifted windows, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 10012–10022.

[24] Lu, J., Xiong, C., Parikh, D., Socher, R., 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning, in: CVPR.

[25] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al., 2022. Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems 35, 27730–27744.

[26] Pan, Y., Yao, T., Li, Y., Mei, T., 2020. X-linear attention networks for image captioning, in: CVPR.

[27] Papineni, K., Roukos, S., Ward, T., Zhu, W., 2002. Bleu: a method for automatic evaluation of machine translation, in: ACL.

[28] Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al., 2018. Improving language understanding by generative pre-training .

[29] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. The Journal of Machine Learning Research .

[30] Scao, T.L., Fan, A., Akiki, C., Pavlick, E., Ilić, S., Hesslow, D., Castagné, R., Luccioni, A.S., Yvon, F., Gallé, M., et al., 2022. Bloom: A 176b-parameter open-access multilingual language model. arXiv preprint arXiv:2211.05100 .

[31] Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F., Dai, J., 2019. Vl-bert: Pre-training of generic visual-linguistic representations. arXiv preprint arXiv:1908.08530 .

[32] Tang, M., Wang, Z., Liu, Z., Rao, F., Li, D., Li, X., 2021. Clip4caption: Clip for video caption, in: Proceedings of the 29th ACM International Conference on Multimedia, pp. 4858–4862.

[33] Tang, M., Wang, Z., Zeng, Z., Li, X., Zhou, L., 2022. Stay in grid: Improving video captioning via fully grid-level representation. IEEE Transactions on Circuits and Systems for Video Technology .

[34] Tanida, T., Müller, P., Kaissis, G., Rueckert, D., 2023. Interactive and explainable region-guided radiology report generation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7433–7442.

[35] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al., 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 .

[36] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: NIPS.

[37] Vedantam, R., Zitnick, C.L., Parikh, D., 2015. Cider: Consensus-based image description evaluation, in: CVPR.

[38] Vinyals, O., Toshev, A., Bengio, S., Erhan, D., 2015. Show and tell: A neural image caption generator, in: CVPR.

[39] Wang, Z., Han, H., Wang, L., Li, X., Zhou, L., 2022a. Automated radiographic report generation purely on transformer: A multi-criteria supervised approach. IEEE Transactions on Medical Imaging .

[40] Wang, Z., Liu, L., Wang, L., Zhou, L., 2023. Metransformer: Radiology report generation by transformer with multiple learnable expert tokens, in: CVPR, pp. 11558–11567.

[41] Wang, Z., Tang, M., Wang, L., Li, X., Zhou, L., 2022b. A medical semantic-assisted transformer for radiographic report generation, in: MICCAI.

[42] Wang, Z., Zhou, L., Wang, L., Li, X., 2021. A self-boosting framework for automated radiographic report generation, in: CVPR.

[43] Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al., 2022. Emergent abilities of large language models. arXiv preprint arXiv:2206.07682 .

[44] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A.C., Salakhutdinov, R., Zemel, R.S., Bengio, Y., . Show, attend and tell: Neural image caption generation with visual attention, in: ICML, 2015.

[45] Yang, S., Wu, X., Ge, S., Zhou, S.K., Xiao, L., 2021. Knowledge matters: Radiology report generation with general and specific knowledge. Medical Image Analysis .

[46] Yin, C., Qian, B., Wei, J., Li, X., Zhang, X., Li, Y., Zheng, Q., 2020. Automatic generation of medical imaging diagnostic report with hierarchical recurrent neural network, in: ICDM.

[47] You, D., Liu, F., Ge, S., Xie, X., Zhang, J., Wu, X., 2021. Aligntransformer: Hierarchical alignment of visual regions and disease tags for medical report generation, in: MICCAI.

[48] You, Q., Jin, H., Wang, Z., Fang, C., Luo, J., 2016. Image captioning with semantic attention, in: CVPR.

[49] Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X.V., et al., 2022. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 .

[50] Zhang, Y., Wang, X., Xu, Z., Yu, Q., Xu, D., 2020. When radiology report generation meets knowledge graph. Proceedings of the AAAI Conference on Artificial Intelligence .