

Large Language Model with Region-guided Referring and Grounding for CT Report Generation

Zhixuan Chen, Yequan Bie, Haibo Jin, and Hao Chen, *Senior Member, IEEE*

Abstract—Computed tomography (CT) report generation is crucial to assist radiologists in interpreting CT volumes, which can be time-consuming and labor-intensive. Existing methods primarily only consider the global features of the entire volume, making it struggle to focus on specific regions and potentially missing abnormalities. To address this issue, we propose Reg2RG, the first region-guided referring and grounding framework for CT report generation, which enhances diagnostic performance by focusing on anatomical regions within the volume. Specifically, we utilize masks from a universal segmentation module to capture local features for each referring region. A local feature decoupling (LFD) strategy is proposed to preserve the local high-resolution details with little computational overhead. Then the local features are integrated with global features to capture inter-regional relationships within a cohesive context. Moreover, we propose a novel region-report alignment (RRA) training strategy. It leverages the recognition of referring regions to guide the generation of region-specific reports, enhancing the model’s referring and grounding capabilities while also improving the report’s interpretability. A large language model (LLM) is further employed as the language decoder to generate reports from integrated visual features, facilitating region-level comprehension. Extensive experiments on two large-scale chest CT-report datasets demonstrate the superiority of our method, which outperforms several state-of-the-art methods in terms of both natural language generation and clinical efficacy metrics while preserving promising interpretability. The code is available at <https://github.com/zhi-xuan-chen/Reg2RG>.

Index Terms—Region-level Understanding, Referring and Grounding, CT Report Generation, Large Language Model.

I. INTRODUCTION

Computed tomography (CT) is widely used in clinical practice and crucial for diagnoses [1]. However, this process is labor-intensive [2] as radiologists need to analyze entire CT

This work was supported by the Hong Kong Innovation and Technology Fund (Project No. MHP/002/22), HKUST (Project No. FS111), and the Research Grants Council of the Hong Kong Special Administrative Region, China (Project Reference Number: T45-401/22-N).

Zhixuan Chen, Yequan Bie, and Haibo Jin are with the Department of Computer Science and Engineering at the Hong Kong University of Science and Technology University, Hong Kong SAR, China. (e-mail: {zchenhi, ybie}@connect.ust.hk, hjinag@cse.ust.hk).

Hao Chen is with the Department of Computer Science and Engineering, Department of Chemical and Biological Engineering and Division of Life Science, Hong Kong University of Science and Technology, Hong Kong, China (e-mail: jhc@cse.ust.hk).

The corresponding author is Hao Chen.

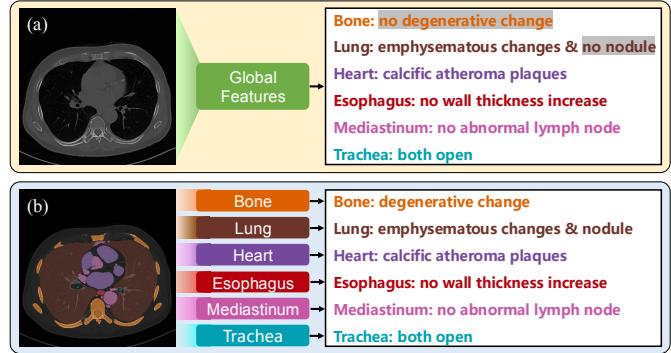


Fig. 1. Our method vs. the vanilla method. The gray background highlights instances of incorrect diagnosis. For ease of comparison, we divide the report into region-level sections. (a) The vanilla method based on global features is prone to neglecting some abnormalities since it fails to explore local details. (b) In contrast, our method can correctly detect all abnormalities with the region-guided local features.

volumes and produce detailed diagnostic reports. To alleviate the workload, automatic CT report generation has attracted increasing attention [3]–[5].

CT report generation involves creating detailed and accurate diagnostic reports from CT volumes, where each anatomical region must be thoroughly analyzed and described to ensure reliable clinical insights. Existing methods [3]–[5] primarily rely on global features extracted from the entire volume to generate reports. While effective to some extent, this approach overlooks the inherent complexity of CT data as a three-dimensional (3D) imaging modality. CT volumes encode rich anatomical details across various regions, presenting significant challenges for decoders to accurately capture region-specific abnormalities using only volume-level embeddings. As depicted in Fig. 1 (a), the vanilla method often fails to identify certain abnormalities when limited to global features alone. Therefore, enhancing the model’s ability to process and integrate region-specific information is crucial for generating comprehensive and clinically relevant CT reports.

Recently, universal segmentation models [6]–[8] have achieved remarkable progress and possess potent zero-shot capabilities. This inspires us to leverage these models as off-the-shelf tools to extract anatomical masks from CT volumes, providing key regional information for report generation. In this work, we propose **Reg2RG**, a **R**egion-guided **R**eferri-

and **Grounding** framework for **R**eport **G**eneration, leveraging the in-context learning and long-term referencing capabilities of LLMs. Referring and grounding are closely related but distinct concepts in visual understanding tasks. Referring focuses on understanding the semantics of designated regions within an image and providing their descriptions. In contrast, grounding involves locating specific regions based on textual information, effectively linking language to visual elements. With these abilities, the model can accurately refer to specific regions for detailed diagnoses and ground reports to the correct regions for better interpretability. As shown in Fig. 1 (b), our method can generate more accurate reports, which are well-grounded in each referring region of the volume.

To reduce computational costs, previous works [4], [5] often downsample volumes or pool features. However, these methods may lose crucial high-resolution texture details [9]. To address this issue, we propose a local feature decoupling (LFD) strategy to preserve local high-resolution texture details and essential geometry with low computational cost.

Since the global volume-level embedding provides significant inter-regional relationships within a cohesive context, we further integrate global features with local features to enable collaboration between them. However, achieving this effective global-local collaboration is non-trivial. Since global features encompass information from all regions, they may interfere with the diagnosis of a specific region, as validated in Table VI. To address this challenge, we propose a novel region-report alignment (RRA) training strategy that enhances the alignment between local features and their corresponding region reports, reducing the influence of irrelevant global information and ensuring that the generated reports are more robust and accurate.

In medicine, interpretability is vital for radiologists to comprehend the visual basis of the generated reports. Previous works [10], [11] mainly use attention maps to link report keywords with image regions. However, these attention maps are often coarse and ambiguous, lacking precision and reliability. The proposed RRA training strategy can also enhance interpretability by establishing a clear association between the referring region and the report. This explicit alignment ensures that the generated reports are firmly grounded in the identified regions, thereby improving both interpretability and reliability.

In summary, our contributions are as follows:

- We propose **Reg2RG**, a novel **R**egion-guided **R**eferring and **G**rounding framework for **R**eport **G**eneration. It generates accurate reports by focusing on target regions from multiple candidates and aligning them with the correct regions for better interpretability. To our knowledge, this is the first work introducing referring and grounding for CT report generation.
- We design a local feature decoupling (LFD) strategy to decouple texture and geometry, preserving high-resolution texture details and essential geometry with minimal overhead. The global features are integrated with local features to collaboratively capture inter-regional relationships within a cohesive context.
- A region-report alignment (RRA) training strategy is utilized to boost the model's referring and grounding

abilities, making it more interpretable and reliable.

- We highlight the critical role of large-scale LLMs in report generation lies in their exceptional region-level referencing and grounding capabilities. These abilities allow for precise focus on region-specific features and inter-regional relationships, enabling the creation of more accurate and clinically relevant reports.
- Experiments on two large-scale 3D chest CT datasets demonstrate the effectiveness of our method, achieving both superior performance and interpretability.

II. RELATED WORKS

A. Medical Report Generation

To alleviate the heavy workload of pathologists, medical report generation has emerged as an effective solution for the automatic interpretation of medical images. The previous works [10]–[13] mostly focus on the 2D chest X-ray (CXR) report generation. To produce higher-quality reports, recent efforts [10], [13], [14] have focused on incorporating additional information to enhance the accuracy of key abnormality details in the generated reports. For instance, PromptMRG [13] utilizes abnormality classification results as diagnostic prompts for the decoder, enhancing both the clinical relevance and effectiveness of the generated reports. Similarly, ORGAN [10] constructs an observation graph to better aggregate clinically significant information, further improving the quality and coherence of the reports. Building on this trend, RGRG [14] utilizes a detector to introduce regional information, supporting fine-grained report generation. However, this approach has not been thoroughly explored in the more expansive and information-rich 3D space of CT volumes, where capturing regional details is particularly essential.

Leveraging large-scale CT-report datasets [15]–[17], automatic CT report generation has also garnered increasing attention recently. CT2Rep [3] represents the first exploration into 3D CT report generation, leveraging a memory-driven decoder to generate detailed reports directly from global volume features. Inspired by the advancements in multi-modal LLMs [18], [19], recent approaches have attempted to bootstrap LLMs for CT report generation. Dia-LLaMA [20] integrates critical diagnostic prompts as prior medical knowledge, while HILT [9] emphasizes efficient encoding strategies for high-resolution volume features to enhance performance. Additionally, several studies [4], [5], [21], [22] have explored the development of medical generalist models based on LLMs, demonstrating their capacity to perform CT report generation as part of broader diagnostic tasks.

However, these approaches generate reports solely based on global features, overlooking the critical role of local anatomical information. Furthermore, LLMs' referencing and context-learning capabilities are underutilized. To address these issues, we propose a region-guided referring and grounding framework that leverages local features to enhance regional comprehension, fully harnessing the capabilities of LLM to optimize CT report generation.

B. Region-level Referring and Grounding

To promote region-level understanding, several works in the general domain [23]–[27] have explored integrating LLMs with region-level features, showcasing their potential to enhance fine-grained reasoning and context-aware interpretations. However, the regional interpretation of medical images remains relatively underexplored. RGRG [14] extracts region features for CXR report generation but generates reports independently for each region, neglecting inter-regional relationships and becoming inefficient with more regions. Some methods [28]–[30] attempt to ground diagnostic text to targeted regions in medical images, yet they fail to utilize region-level features for generating fine-grained, context-aware descriptions. MedRegA [31] refers to specific regions using bounding box coordinates in prompts but may be less accurate compared to approaches that leverage detailed regional features. The ideal integration and utilization of regional information in medical imaging tasks remains an open question.

Given the importance of geometric information in assessing lesion size and position in medicine, we propose a novel local feature decoupling strategy that preserves complete geometric information while retaining detailed texture features. Inspired by Groma [25], we achieve grounded report generation by referencing regions from a pool of candidates, enabling focused and relevant analysis. Additionally, we integrate global features with local features to capture inter-regional relationships and provide a cohesive contextual understanding of the entire image. To enhance regional analysis, we propose a region-report alignment training strategy to explicitly link visual regions with their reports, enabling accurate and coherent multi-region report generation in a single inference.

III. METHODS

In this section, we first overview our framework in Sec. III-A. Next, we detail the local feature decoupling (LFD) strategy in Sec. III-B, and explain its integration with global features in Sec. III-C. Finally, we describe the region-report alignment (RRA) training strategy in Sec. III-D.

A. Overview

The overview of our method **Reg2RG** is shown in Fig. 2. CT report generation involves creating report **R** based on CT volumes $\mathbf{V} \in \mathbb{R}^{H \times W \times D}$. Unlike the previous works [3], [5] that use only global features \mathcal{G} of \mathbf{V} , we additionally extract a set of local features $\mathcal{L} = \{\mathcal{L}_1, \dots, \mathcal{L}_n\}$ with the universal segmentation module. To preserve the local high-resolution texture details and significant geometry information with minimal computational overhead, local features are decoupled into texture and geometry. The decoupled local features work alongside global features to generate the report **R** using the LLM. Furthermore, we propose a training strategy that strengthens the alignment between the local feature \mathcal{L}_i and the region-specific report R_i , thereby improving diagnostic accuracy and enhancing the reliability of the generated report. The overall generating process can be described as follows:

$$\mathbf{R} = \text{LLM}(\mathcal{G}, \mathcal{L}) = \text{LLM}(\mathcal{G}, \mathcal{L}_1, \dots, \mathcal{L}_n), \quad (1)$$

where n is the number of referring regions in \mathbf{V} .

B. Local Feature Decoupling Strategy

As shown in Fig. 2, we first utilize the existing universal segmentation module f_S to extract mask M_{A_j} for the anatomical area A_j given the CT volumes \mathbf{V} . The segmentation process can be formulated as follows:

$$\{M_{A_1}, \dots, M_{A_n}\} = f_S(\mathbf{V}). \quad (2)$$

The region mask M_{A_j} is then utilized to construct the corresponding local features \mathcal{L}_{A_j} . To preserve higher-resolution details without increasing much computational burden, we design a local feature decoupling (LFD) strategy that separates texture and geometry information. Texture information refers to patterns within regions of interest that capture surface characteristics, which are essential for disease diagnosis. For example, in lung CT scans, the texture feature “ground-glass opacity” indicates increased lung tissue density. Geometry information encompasses the size, shape, and spatial location of regions of interest, playing a crucial role in medical imaging. For example, the geometry feature “enlarged heart” may indicate conditions such as heart failure or pericardial effusion.

For the texture information, we first use region mask M_{A_j} to extract the region volume V_{A_j} from \mathbf{V} by element-wise multiplication. This results in a large redundant area outside the region of interest, which is not informative. Therefore, we crop the region volume V_{A_j} to exclude these irrelevant parts. Since V_{A_j} is much smaller than the entire volume \mathbf{V} , we can retain higher-resolution details without increasing the input size. Next, we utilize a 3D volume encoder f_V to extract local texture features $\mathcal{L}_{A_j}^t$. The following adapter module f_A is used to compress and align these local texture features with the embedding space of the LLM. This process can be formulated as follows:

$$\mathcal{L}_{A_j}^t = f_A(f_V(\text{Crop}(V_{A_j}))) = f_A(f_V(\text{Crop}(M_{A_j} \odot \mathbf{V}))), \quad (3)$$

where \odot indicates element-wise multiplication.

Previous works [23]–[27] extract local features by only encoding cropped regions, focusing exclusively on texture features while neglecting the essential geometry features necessary for assessing lesion size and location in medical contexts. In contrast, our method incorporates them as supplementary features. Specifically, we introduce the geometry information by encoding the region mask M_{A_j} , which is uncropped to preserve the original size and position. A lightweight mask encoder f_M is used to extract geometry features $\mathcal{L}_{A_j}^g$, followed by a fully connected layer f_P to project these features for LLM input. This process is formulated as follows:

$$\mathcal{L}_{A_j}^g = f_P(f_M(M_{A_j})). \quad (4)$$

The local features \mathcal{L}_{A_j} are obtained by concatenating the texture features $\mathcal{L}_{A_j}^t$ and geometry features $\mathcal{L}_{A_j}^g$:

$$\mathcal{L}_{A_j} = \text{Concat}(\mathcal{L}_{A_j}^t, \mathcal{L}_{A_j}^g). \quad (5)$$

Typically, the local features \mathcal{L}_A comprise multiple regions \mathcal{L}_{A_j} , each providing specific information for a distinct anatomical area within the CT volumes:

$$\mathcal{L}_A = \{\mathcal{L}_{A_1}, \dots, \mathcal{L}_{A_n}\}. \quad (6)$$

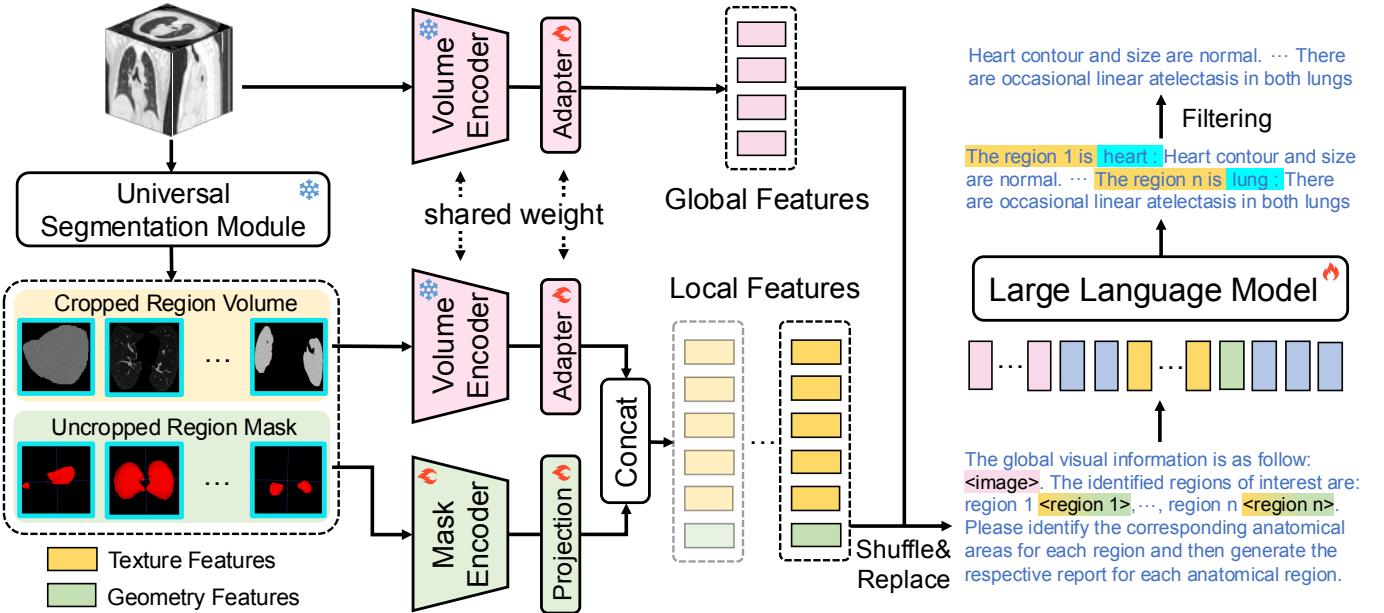


Fig. 2. Overview of the proposed **Reg2RG** framework. It integrates global and local features as visual embeddings for the LLM to generate reports. Global features are encoded from the entire volume, while local features are extracted using segmentation masks to capture lesion details in sub-regions. The local features are decoupled into texture and geometry, where texture is derived from cropped masked volumes and geometry is obtained from the uncropped masks. Shuffling local features across various regions enhances the alignment between visual regions and their corresponding reports. The LLM focuses on each region individually to produce accurate and detailed region-specific reports.

Note that \mathcal{L}_A is not identical to \mathcal{L} in Eq. (1); further details will be provided in Sec. III-D.

C. Global-Local Features Collaboration

In medicine, different regions are interrelated rather than isolated. Therefore, we incorporate global features to provide contextual information. Specifically, the same volume encoder f_V and adapter f_A (as in Sec. III-B) are used to extract global features \mathcal{G} :

$$\mathcal{G} = f_A(f_V(\mathbf{V})). \quad (7)$$

The collaboration of global and local features is achieved by embedding them into the prompt. Our designed prompt \mathcal{P} consists of two parts: $\mathcal{P} = \{\mathcal{I}, \mathcal{T}\}$, where \mathcal{I} denotes the special tokens for the visual embedding and \mathcal{T} represents the text tokens of the instruction. As depicted in Fig. 2, we utilize $\langle\text{image}\rangle$ and $\langle\text{region } i\rangle$ as the special tokens for the global and local features, respectively. These special tokens are replaced by the corresponding features \mathcal{G} and \mathcal{L}_i , which then interact within the LLM to generate \mathbf{R} .

D. Region-Report Alignment Training Strategy

To enhance the explicit link between the referring region and the report, we propose a training strategy that first recognizes the anatomical area of the referring region and then generates the corresponding region report R_i .

While the segmentation module can provide the anatomical area name, this information may introduce bias during report generation, as the model might rely on the name rather than the actual local features. For instance, if the input includes the

name “lung” for a region, the model may focus solely on generating content related to the “lung”, potentially overlooking the real information within the local features. Therefore, we train the model to recognize this information independently from the local features. This strategy helps the model better understand the referred region, making the generated report more reliably grounded. Furthermore, we shuffle the order of the local features \mathcal{L}_A at each step to prevent the model from associating anatomical areas with a fixed sequence. Thus, although the content of \mathcal{L} and \mathcal{L}_A remains unchanged, the order of each local feature varies between them. The generation process is defined as follows:

$$\begin{aligned} \mathbf{R} &= \{(P_1, R_1), \dots, (P_n, R_n)\} \\ &= \text{LLM}(\mathcal{G}, \mathcal{L}_1, \dots, \mathcal{L}_n) \\ &= \text{LLM}(\mathcal{G}, \text{Shuffle}(\mathcal{L}_{A_1}, \dots, \mathcal{L}_{A_n})). \end{aligned} \quad (8)$$

During training, region-level reports R_i of each referring region are used as ground truth. In addition, we add a prefix $P_i = \text{“The region } [i] \text{ is [area name]”}$ before R_i to indicate the anatomical area name. The model learns to recognize the anatomical area name by predicting this prefix. Since anatomical region names are fixed for each specific region, they provide a clear and structured target for the model, enabling it to develop a more precise semantic understanding of the referred regions. The restructured report remains a sequence of text tokens: $\mathbf{R} = \{(P_1, R_1), \dots, (P_n, R_n)\} = \{r_1, r_2, \dots, r_T\}$ (T is the report length). This format integrates seamlessly with the native auto-regressive training process of the LLM. The training process is optimized by minimizing the language-

modeling loss as shown below:

$$\mathcal{L}_{LM} = -\sum_{t=1}^T \log p(r_t | \mathcal{P}, r_1, \dots, r_{t-1}), \quad (9)$$

where \mathcal{P} indicates the prompt, as mentioned in Sec. III-C.

When evaluating the quality of the generated report, the prefix P_i is removed from the R .

IV. EXPERIMENTS AND RESULTS

A. Datasets and Evaluation Metrics

Datasets. We train and evaluate our model alongside comparative methods using two large-scale chest CT datasets, ensuring a more comprehensive assessment.

The RadGenome-ChestCT dataset [16], designed for region-guided 3D chest CT interpretation, is based on the CT-RATE [15] dataset. It consists of 25,692 region-guided CT-report pairs sourced from 21,304 patients. The CT volumes possess a consistent voxel spacing of $1 \text{ mm} \times 1 \text{ mm} \times 3 \text{ mm}$, with anatomical masks generated using the SAT [6] segmentation module. Region-grounded reports are generated with GPT-4 [32] and a named entity recognition model, covering 10 chest anatomical regions: abdomen, bones, breasts, esophagus, heart, lungs, trachea and bronchi, mediastinum, pleura, and thyroid. Following the official data split, 24,128 pairs are allocated for training, while 1,564 pairs are reserved for evaluation.

The CTRG-Chest-548K dataset [17] includes 1,804 CT volume-report pairs. To extract anatomical masks, we employ the segmentation module SAT [6]. Since region-level grounded reports are not provided within this dataset, we utilize Qwen2.5-14B [33] to segment the reports into region-level sections. Regarding the data split, the dataset is randomly partitioned into training and testing sets in an 8:2 proportion.

Evaluation Metrics. In line with prior studies [3], [9], [20], we employ widely recognized natural language generation (NLG) metrics like BLEU-n [34], METEOR [35], and ROUGE-L [36] for evaluation. BLEU-n measures the overlap of n-grams (sequences of n words) between generated and reference reports to gauge word sequence similarity. ROUGE-L assesses by comparing the longest common word subsequence, focusing on textual alignment. METEOR effectively incorporates synonyms and paraphrases, offering a more nuanced assessment of semantic similarity between reports.

NLG metrics focus on word and sentence similarity but neglect the diagnostic accuracy. For example, “The heart is enlarged” and “The heart is not enlarged” can have similar NLG scores despite opposite diagnostic conclusions. Therefore, we also adopt the clinical efficacy (CE) metrics [11]. Following CT-CLIP [15], we employ the RadBERT [37] text classifier to extract 18 types of abnormality labels from chest CT reports. As this classifier is fine-tuned on CT-RATE [15], it ensures high-quality labels for the RadGenome-ChestCT dataset [16], derived from CT-RATE. CE metrics are the precision, recall, and F1 score from comparing abnormality labels of generated and ground-truth reports, providing a clinically meaningful assessment by highlighting significant abnormalities. Since RadBERT is not specifically designed for the CTRG-Chest-548K dataset, the types of abnormalities RadBERT focuses

on do not align with the primary abnormalities in this dataset. Therefore, we only use the CE metrics for the RadGenome-ChestCT dataset.

B. Implementation Details and Baselines

We use the text-prompted segmentation model SAT-Pro [6] to extract anatomical masks when unavailable. The SAT-Pro model is configured with 256 object queries, and its input CT intensities are normalized to the range $[0, 1]$ using min-max normalization. For our input, both global and local volumes are resized to $256 \times 256 \times 64$. We use the pre-trained ViT3D and Perceiver from RadFM [4] as the volume encoder f_V and adapter f_A to extract texture features, each represented by 32 visual embeddings. Given that geometric information is sparser than texture information, a lightweight 3-layer ViT3D serves as the mask encoder f_M to extract geometric features. These features are pooled into a single embedding and projected by f_P to align with the LLM’s embedding space. For the LLM, we choose LLaMA2-7B [38] and utilize LoRA [39] for parameter-efficient fine-tuning. The LoRA configuration is set with a rank $r = 8$, a scaling factor $\alpha = 32$, and a dropout rate of 0.1.

We train the model with the AdamW optimizer [40] at an initial learning rate of 5×10^{-5} , following a constant learning rate schedule with a warmup phase. Training on the RadGenome-Chest dataset takes 48 hours on two RTX 4090 GPUs with PyTorch 2.0, running for 6 epochs with an effective batch size of 16. For the CTRG-Chest-584K dataset, training takes 24 hours for 10 epochs with the same batch size. To optimize memory usage, ZeRO [41] stage 2 is applied alongside gradient checkpointing [42]. We use the checkpoint of the last epoch for evaluation.

For baseline comparisons, we choose three state-of-the-art 3D methods capable of generating CT reports: CT2Rep [3], RadFM [4], and M3D [5]. We also include two 2D methods R2GenGPT [43] and MedVInT [22], which support radiology report generation. For fairness, we employ LLaMA2-7B [38] as the language decoder for all compared methods. Input volumes were resized to $256 \times 256 \times 64$ and represented by 32 visual tokens for all methods. To adapt the 2D methods, we convert the 3D volumes into 2D multi-channel images. Each compared model is initialized with their pre-trained weights and is further fine-tuned on the CT report generation dataset.

C. Quantitative Results

1) Natural Language Generation Metrics: Table I presents the NLG results of our model alongside comparisons to other methods. On the RadGenome-ChestCT dataset, our model outperforms all others across all NLG metrics, underscoring its capability to generate high-quality reports. Specifically, we take the MedVInT [22] with second-best results as an example. Our model achieves a relative improvement of 1.1% to 6.7% on BLEU metrics, highlighting enhanced expression similarity to the reference reports. For the METEOR metric considering inflectional variations and synonym matching, our model shows a notable 9.1% improvement, indicating better lexical flexibility and semantic alignment. Our model also

surpasses the second-best by 12.4% in ROUGE-L, highlighting its consistent performance across metrics. A similar pattern is observed on the CTRG-Chest-584K dataset, where our method outperforms the second-best model with a 2.0% to 3.7% improvement in BLEU metrics and a 0.8% gain in METEOR. The lower ROUGE-L score on this dataset may be attributed to the fragmented nature of region-level report generation. This fragmentation affects the coherence between reports for different regions and consequently impacts the evaluation of the longest common sequence. However, this inconsistency does not affect the quality of individual region reports, as evidenced by higher performance on other metrics. These results highlight the effectiveness of our model in generating high-quality reports.

2) Clinical Efficacy Metrics: Table II showcases the CE performance of our model compared to other methods on the RadGenome-ChestCT dataset. Our model surpasses the second-best approach by 3.9%, 22.3%, and 19.3% in precision, recall, and F1 score, respectively. This demonstrates the superiority of our model in generating reports with higher diagnostic accuracy. It is worth noting that there is an inherent trade-off between precision and recall. Achieving higher precision requires minimizing false positives, which often leads the model to adopt a more conservative approach to predicting abnormalities. On the other hand, higher recall necessitates reducing false negatives, encouraging a more aggressive stance in abnormality detection. Therefore, balancing these competing priorities is particularly challenging, as demonstrated by the performance of other models. For instance, MedVInT [22] achieves the second-highest recall but struggles with relatively low precision, while M3D [5] exhibits the reversed trend, favoring precision at the expense of recall. In contrast, our model effectively balances this trade-off, maintaining high performance in both metrics and achieving a significantly improved F1 score. These results underscore our model's ability to maintain clinical relevance and diagnostic reliability while delivering high linguistic quality in the generated reports.

3) Region Recognition Results: Table III presents the region recognition performance of our model on the RadGenome-ChestCT dataset [16]. To enhance the model's interpretability and reliability of generated reports, our approach explicitly requires the model to first identify the anatomical area corresponding to the referring region before generating the associated report. This intermediate recognition step simplifies evaluation and interpretation compared to directly analyzing the generated reports, providing an early-stage validation of both their alignment with the target regions and the reliability of the reports. The results show that our model accurately identifies most anatomical regions except for the lung and pleura, suggesting the generated reports are reliably aligned with the target regions, thereby enhancing interpretability. The subpar results for the lung and pleura are attributed to the performance of the SAT [6] segmentation module, which struggles to produce masks that adequately differentiate between these closely related regions. Once the segmentation model provides correct region masks, our model can effectively identify referring regions, leading to more trustworthy and clinically reliable reports.

4) Region-level Reports Evaluation: Table IV showcases the evaluation results of the region-level reports generated by our model on the RadGenome-ChestCT dataset [16]. We observe that the performance varies for reports from different anatomical regions. The higher occurrence frequency of certain regions in reports contributes to better performance on the abdomen, mediastinum, heart, and lung. In contrast, performance tends to be lower for regions that appear less frequently, such as the breast and thyroid. The subpar CE performance on the bone and trachea & bronchi may be attributed to the RadBERT model from CT-CLIP [15] rarely considering abnormalities in these regions. Investigating approaches to improve the diagnosis of less common regions and abnormalities could serve as a direction for future research.

D. Ablation Study

To demonstrate the effectiveness of each component in our model, we conduct comprehensive ablation studies on the RadGenome-ChestCT dataset [16]. Given that our volume encoder f_V and adapter f_A are initialized using the pre-trained checkpoint from RadFM [4], we take it as the baseline for comparison. First, we examine the contribution of the local feature decoupling (LFD) strategy in Sec. IV-D.1. Next, we investigate the effectiveness of global-local features collaboration by incrementally adding different visual features in Sec. IV-D.2. Additionally, we analyze the performance of our region-report alignment (RRA) training strategy in Sec. IV-D.3. Finally, we validate the necessity of employing large-scale LLMs for region-level referring and grounding in Sec. IV-D.4. The results of these experiments are detailed in Table V, Table VI, Table VII and Table VIII, demonstrating the individual and collective impact of each component.

1) Local Feature Decoupling Strategy: First, we validate the efficacy of the local feature decoupling (LFD) strategy in improving model performance. In the baseline setup where decoupling is not applied, the model uses masked volumes without cropping as local features, combining texture and geometry information without separation. As Table V shows, the model with decoupled features demonstrates superior performance across most metrics. The significant improvement in CE metrics highlights the value of local high-resolution details from decoupled texture information in enhancing diagnostic accuracy. The slight decreases in BLEU-4 and ROUGE-L scores may be due to these metrics relying only on word overlap, which often fails to capture nuanced semantic similarities. By contrast, the higher METEOR score supports this assumption because it takes synonym matching and semantic relationships into consideration.

2) Global-Local Features Collaboration: As shown in Table VI, we assess the effectiveness of global-local features collaboration by incrementally adding different visual features across settings (a) to (c). Without position and size information, the performance of setting (a) falls significantly below the baseline, underscoring the critical role of geometric information in medical imaging for accurately representing local features. Including geometric features in setting (b) leads to notable improvement, surpassing the baseline in most metrics

TABLE I

THE NLG PERFORMANCE OF OUR MODEL COMPARED WITH OTHER SOTA METHODS ON TWO LARGE-SCALE DATASETS. THE BEST AND SECOND-BEST RESULTS ARE IN **BOLD** AND UNDERLINED, RESPECTIVELY. A HIGHER VALUE INDICATES BETTER PERFORMANCE.

Dataset	Method	Year	BL-1	BL-2	BL-3	BL-4	MTR	RG-L
RadGenome-ChestCT	R2GenGPT [43]	2023	43.28	34.11	28.16	24.16	39.85	32.26
	MedVInT [22]	2023	44.28	<u>34.91</u>	<u>28.75</u>	<u>24.60</u>	<u>40.39</u>	32.58
	RadFM [4]	2023	44.20	34.49	28.06	23.65	39.94	31.53
	CT2Rep [3]	2024	<u>44.42</u>	34.43	27.94	23.56	40.16	30.99
	M3D [5]	2024	43.57	34.48	28.54	24.49	39.95	<u>32.61</u>
	Reg2RG	-	47.25	36.49	29.57	24.87	44.07	36.65
CTRG-Chest-584K	R2GenGPT [43]	2023	41.82	36.37	32.70	30.10	47.05	50.93
	MedVInT [22]	2023	47.38	39.60	34.28	30.68	<u>49.32</u>	49.53
	RadFM [4]	2023	<u>48.66</u>	<u>40.28</u>	<u>34.73</u>	<u>30.89</u>	49.18	49.08
	CT2Rep [3]	2024	42.28	36.16	32.08	29.19	47.00	50.17
	M3D [5]	2024	46.27	39.02	34.23	30.86	49.26	<u>50.24</u>
	Reg2RG	-	49.63	41.43	35.91	32.04	49.71	47.76

TABLE II

THE CLINICAL EFFICACY PERFORMANCE OF OUR MODEL AND OTHER SOTA METHODS ON THE RADGENOME-CHESTCT [16] DATASET. THE BEST AND SECOND-BEST RESULTS ARE HIGHLIGHTED IN **BOLD** AND UNDERLINED, RESPECTIVELY. A HIGHER VALUE IMPLIES BETTER PERFORMANCE ACROSS ALL METRICS.

Method	Pre.	Rec.	F1
R2GenGPT [43]	0.340	0.066	0.110
MedVInT [22]	0.377	<u>0.148</u>	<u>0.212</u>
RadFM [4]	0.382	0.131	0.195
CT2Rep [3]	0.317	0.089	0.139
M3D [5]	<u>0.407</u>	0.090	0.148
Reg2RG	0.423	0.181	0.253

and demonstrating the value of spatial and structural details. Further incorporating global features in setting (c) boosts performance across most CE and NLG metrics, highlighting the efficacy of global context in capturing inter-regional relationships and improving report coherence. However, we observe that the improvement in recall comes at the expense of precision. This trade-off may result from the influence of abnormality information within global features, which may misrepresent the diagnosis of individual regions and lead to an increase in false positives across regions.

3) *Region-Report Alignment Training Strategy*: In Table VII, we validate the efficacy of our region-report alignment (RRA) training strategy, which guides the LLM to generate reports based on referring region information for reliable grounding. The results indicate that our model with RRA outperforms the one without it across all CE and NLG metrics, except for ROUGE-L. This improvement demonstrates the effectiveness of the RRA strategy in aligning region-specific features with report generation, ensuring more accurate and clinically relevant outputs. The slightly lower ROUGE-L score can be

TABLE III

THE REGION RECOGNITION RESULTS OF OUR MODEL. THE TRA. & BRO. REGION REFERS TO THE TRACHEA AND BRONCHI.

Region	Pre.	Rec.	F1
Abdomen	0.997	0.996	0.997
Bone	0.996	0.998	0.997
Breast	0.945	0.983	0.964
Esophagus	0.997	0.999	0.998
Heart	0.995	0.996	0.996
Lung	0.443	0.443	0.443
Mediastinum	0.991	0.997	0.994
Pleura	0.442	0.441	0.442
Thyroid	0.991	0.931	0.960
Tra. & Bro.	0.986	0.990	0.988

attributed to its reliance on exact sequence matching, which may overlook the use of synonyms or paraphrased expressions. As a result, ROUGE-L may not fully reflect report quality in cases where the phrasing differs but the semantic content remains consistent. The top METEOR performance supports this, as it better captures semantic similarity by considering synonyms and paraphrasing.

Additionally, the proposed RRA strategy mitigates the precision-recall trade-off in (c). By directing the LLM to reference explicit region information, our model not only ensures that each report is grounded in the correct region but also improves its diagnostic and linguistic quality.

4) *LLM as the Language Decoder*: To validate the necessity of employing a large-scale LLM as the language decoder for region-level referring and grounding, we conduct a comparative study using GPT-2 [44] as the decoder. As demonstrated in Table VIII, the LLaMA2-7B model consistently outperforms the GPT-2 model across all metrics, highlighting the advan-

TABLE IV

THE EVALUATION OF REGION-LEVEL REPORTS GENERATED BY OUR MODEL. THE TRA. & BRO. REGION REFERS TO THE TRACHEA AND BRONCHI. THE **AMOUNT** REPRESENTS THE NUMBER OF REPORTS THAT INCLUDE DESCRIPTIONS FOR THE SPECIFIC REGION.

Region	Amount	CE Metrics			NLG Metrics			
		Precision	Recall	F1-score	BLEU-1	BLEU-4	METEOR	ROUGE-L
Abdomen	1517	0.546	0.288	0.377	39.50	25.13	43.33	38.52
Bone	1509	0.000	0.000	0.000	41.05	28.49	47.15	42.13
Breast	56	0.174	0.211	0.190	22.46	13.63	23.14	21.96
Esophagus	1323	0.222	0.021	0.039	59.58	45.16	56.39	57.49
Heart	1418	0.320	0.135	0.190	39.85	27.44	44.75	43.19
Lung	1514	0.393	0.126	0.191	27.22	15.51	31.74	30.82
Mediastinum	1513	0.557	0.242	0.337	37.65	20.09	41.98	33.68
Pleura	1169	0.333	0.069	0.115	37.20	26.91	38.84	47.21
Thyroid	42	0.125	0.040	0.061	21.83	10.64	19.49	19.47
Tra. & Bro.	1401	0.000	0.000	0.000	48.42	38.93	61.47	56.18

TABLE V

THE EFFECTIVENESS OF LOCAL FEATURE DECOUPLING (LFD) STRATEGY.

LFD	CE Metrics			NLG Metrics			
	Precision	Recall	F1-score	BLEU-1	BLEU-4	METEOR	ROUGE-L
	0.349	0.140	0.200	46.85	25.69	43.74	37.64
✓	0.423	0.181	0.253	47.25	24.87	44.07	36.65

TABLE VI

THE EFFECTIVENESS OF DIFFERENT VISUAL FEATURES. THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**. THE **TXT** AND **GEO** DENOTE THE TEXTURE AND GEOMETRIC INFORMATION OF LOCAL FEATURES, RESPECTIVELY, WHEREAS **GLB** STANDS FOR THE GLOBAL FEATURES.

Setting	Visual Features			CE Metrics			NLG Metrics			
	TXT	GEO	GLB	Pre.	Rec.	F1	BL-1	BL-4	MTR	RG-L
Baseline			✓	0.382	0.131	0.195	44.20	23.65	39.94	31.53
(a)	✓			0.336	0.089	0.141	43.72	20.77	39.82	35.12
(b)	✓	✓		0.394	0.143	0.210	45.31	23.35	41.79	37.20
(c)	✓	✓	✓	0.372	0.176	0.239	46.21	23.53	42.94	37.01

tages of large-scale LLMs in this context. While the GPT-2 decoder performs competently with only global features in report generation [45], [46], it struggles to effectively refer to the specific regional features and capture the inter-regional relationships required for region-specific reports. Complex region-level referring and grounding need the powerful in-context learning and long-term referencing abilities of large-scale LLMs, underscoring their importance in generating accurate and contextually rich medical reports tailored to specific anatomical regions.

E. Qualitative Results

1) *Report Length Distributions*: Following [11], we analyze the report length distributions of both generated and ground-truth reports. Fig. 3 presents the distributions of report lengths for the ground-truth reports alongside those generated by our method and the SOTA MedVInT [22]. We leverage Kernel

Density Estimation (KDE) to visualize the probability distributions and compute KL divergence to quantify the differences between the distributions of our method and MedVInT relative to the ground-truth reports. The results indicate that our model generates reports with lengths more closely aligned with the ground-truth reports than those generated by MedVInT, as reflected in the lower KL divergence. This suggests that the reports generated by our model are more complete and accurate, whereas MedVInT tends to produce shorter reports, potentially leading to information loss.

2) *Analysis of Generated Reports*: We present two cases from the RadGenome-ChestCT dataset [16] in Fig. 4. Different colors represent the various correctly diagnosed regions described in the reports, while the gray background highlights incorrect diagnoses. Due to the full report length, we focus on the most relevant sections concerning abnormalities in the ground-truth reports.

In the first case, the MedVInT model misdiagnoses multiple

TABLE VII
THE EFFECTIVENESS OF REGION-REPORT ALIGNMENT (RRA) STRATEGY.

RRA	CE Metrics			NLG Metrics			
	Precision	Recall	F1-score	BLEU-1	BLEU-4	METEOR	ROUGE-L
✓	0.372	0.176	0.239	46.21	23.53	42.94	37.01
	0.423	0.181	0.253	47.25	24.87	44.07	36.65

TABLE VIII
THE EFFECTIVENESS OF LARGE-SCALE LLM AS THE LANGUAGE DECODER.

Language Decoder	CE Metrics			NLG Metrics			
	Precision	Recall	F1-score	BLEU-1	BLEU-4	METEOR	ROUGE-L
GPT-2 [44]	0.357	0.169	0.229	35.47	14.72	35.76	25.18
LLaMA2-7B [38]	0.423	0.181	0.253	47.25	24.87	44.07	36.65

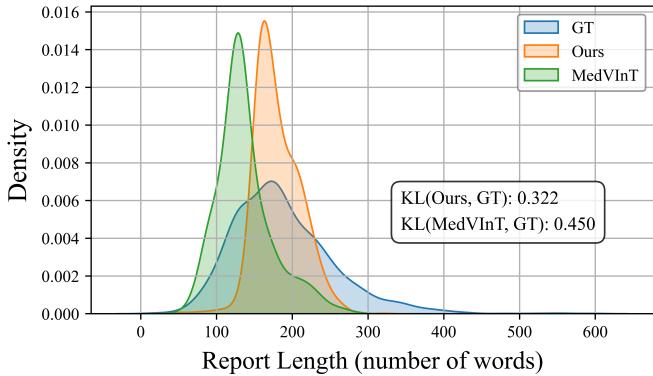


Fig. 3. The report length distributions of the ground-truth reports, along with those generated by our proposed method and the SOTA MedVInT [22]. The Kullback-Leibler (KL) divergence is utilized to quantify the differences between the distributions of our method and MedVInT relative to the ground-truth reports.

regions as normal, missing critical abnormalities. In contrast, our model identifies most abnormalities except the liver lesion. This demonstrates our model’s ability to provide more comprehensive and precise diagnostic information. Notably, it also accurately identifies the location and severity of emphysematous changes in the lungs.

In the second case, our model accurately identifies all abnormalities except for bone degenerative change, whereas MedVInT fails to detect any of them mentioned in the ground-truth report. Although MedVInT points out an abnormality of linear atelectatic changes in the lungs, the diagnosis is likely incorrect as it is not referenced in the ground-truth report. These results demonstrate that our model achieves higher diagnostic accuracy.

3) Region-level Reports with Referring and Grounding: A key advantage of our model is its ability to generate region-level reports explicitly grounded to specific regions. As illustrated in Fig. 5, our generated reports are segmented into distinct sections, each corresponding to a particular anatomical region. We use the same color scheme as in Fig. 4 to represent

the region masks and reports associated with different regions.

Each region-level report begins with the referring region, providing an initial hint about the focus of the LLM. The area recognition result is presented before the report, helping verify the reliability of the generated report. Successfully identifying the referring region indicates that the LLM refers to the correct regional information, enabling the report to be properly grounded in the corresponding region. This enhances the model’s interpretability and reliability of the reports, which is valuable for clinical practice. Conversely, if the referenced region is misidentified, the report becomes unreliable regardless of its content and cannot be definitely linked to any region. This strategy provides a straightforward yet effective mechanism for validating reports and offers reliable reference information to assist radiologists in their interpretation.

4) Segmentation Results on the Two Datasets: We present two segmentation cases from both the RadGenome-ChestCT and CTRG-Chest-584K datasets in Fig. 6. Benefiting from the superior performance of SAT [6], the segmentation results on the RadGenome-ChestCT dataset are highly satisfactory, providing accurate regional information that enables our model to achieve outstanding performance. Regarding the CTRG-Chest-584K dataset, its segmentation results are relatively coarser due to the lower quality of the volumes in this dataset. This may limit the potential of our model, as evidenced by Table I, where its improvement over the second-best method on CTRG-Chest-584K is smaller than that on RadGenome-ChestCT. However, our model still outperforms other methods by effectively leveraging regional information, demonstrating the robustness and effectiveness of our framework.

V. LIMITATIONS AND FUTURE DIRECTIONS

Despite the encouraging outcomes, several limitations of our model warrant further investigation in the future.

Our framework utilizes the segmentation results to acquire local features. The inaccurate or incomplete masks, caused by segmentation failures, could potentially affect the performance of our model to some extent. Therefore, strengthening the robustness of our method against segmentation errors is

	Ground-truth There are several millimetric nonspecific nodules in both lungs. The largest of the described metastatic lesions is observed in the diaphragmatic dome localization at the junction of segment 7-8, and its longest diameter was 17 mm. Intervertebral disc distances are narrowed. There are osteophytes in the vertebral corpus corners. There are minimal emphysematous changes and occasional linear atelectasis in both lungs. Hypodense lesions were observed in both lobes of the liver. The neural foramina are narrowed. There are atheromatous plaques in the aorta and coronary arteries	MedViNT When examined in the lung parenchyma window; No mass lesion was observed in both lung parenchyma. No lytic-destructive lesions were detected in bone structures. Upper abdominal organs are included in the study partially and evaluated as suboptimal. There are minimal emphysematous changes in both lungs. There are calcific atheromatous plaques in the aortic arch and descending aorta	Ours No space-occupying lesion was detected in the liver that entered the cross-sectional area. Degenerative changes were observed in the bone structure. Calcific atheroma plaques were observed in the thoracic aorta and coronary arteries. When examined in the lung parenchyma window; Mild emphysematous changes were observed in both lungs. A 1 cm diameter subpleural nonspecific parenchymal nodule was observed in the middle lobe of the right lung.
Abdomen Bone Breast Esophagus Heart Lung Mediastinum Pleura Thyroid Trachea and Bronchie			
	Ground-truth Minimal calcified atheroma plaques are observed on the walls of the aortic arch and coronary vascular structures. There are osteophytic degenerative changes in the vertebra corpus corners that tend to merge in the right anterolateral. An increase in thoracic kyphosis is observed. There are centriacinar emphysematous changes in both lungs. There are sequela parenchymal changes in the upper lobe of the left lung, the inferior lingular segment, and the medial segment of the middle lobe of the right lung in the lower lobe.	MedViNT As far as can be seen, the mediastinal main vascular structures are normal. Vertebral corpus heights are preserved. Bone structures in the study area are natural. When examined in the lung parenchyma window; Linear atelectatic changes were observed in the left lung upper lobe inferior lingular and right lung middle lobe medial segment. No mass lesion-active infiltration with distinguishable borders was detected in the lung parenchyma.	Ours There are calcified atheromatous plaques in the walls of the thoracic aorta and coronary vascular structures. No lytic or destructive lesions were observed in the bone structures within the image, and the vertebral corpus heights were preserved. There are minimal emphysematous changes in both lungs. No active infiltration or mass lesion was detected in both lungs. There are sequela parenchymal changes in the apex of both lungs.
Abdomen Bone Breast Esophagus Heart Lung Mediastinum Pleura Thyroid Trachea and Bronchie			

Fig. 4. Case studies of our model and the SOTA MedViNT [22]. The different colors represent distinct anatomical areas, as shown at the bottom of each example. The gray background highlights incorrect diagnoses.

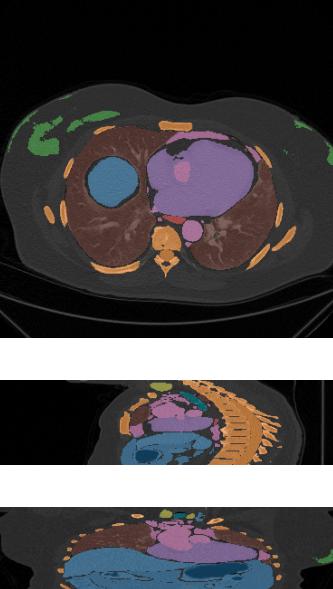
	Our Region-level Reports The region 0 is abdomen: There are calcified atheromatous plaques in the walls of the thoracic aorta and coronary vascular structures. In the upper abdominal sections within the image, no solid mass was detected as far as it can be observed within the borders of non-contrast CT. The region 1 is bone: No lytic or destructive lesions were observed in the bone structures within the image, and the vertebral corpus heights were preserved. The region 2 is breast: As far as can be observed, the left breast was not observed (operated). The region 3 is esophagus: No pathological increase in wall thickness is observed in the thoracic esophagus. The region 4 is heart: Calibration of vascular structures, heart contour and size are natural. Mediastinal vascular structures and cardiac examination could not be evaluated optimally because of the lack of IV contrast. The region 5 is lung: There are minimal emphysematous changes in both lungs. No active infiltration or mass lesion was detected in both lungs. There are sequela parenchymal changes in the apex of both lungs. The region 6 is mediastinum: No lymph node was detected in the mediastinum in pathological size and appearance. There are calcified atheromatous plaques in the walls of the thoracic aorta and coronary vascular structures. Mediastinal vascular structures and cardiac examination could not be evaluated optimally because of the lack of IV contrast. The region 7 is pleura: No pericardial, pleural effusion or thickness increase was observed. The region 8 is thyroid: In the thyroid gland, there is a hypodense nodule with a diameter of 12 mm in the right lobe and 14 mm in the left lobe. The region 9 is trachea and bronchie: Trachea, both main bronchi are open and no occlusive pathology is detected.
---	--

Fig. 5. Region-level reports generated by our model. Each regional report refers to a specific region and is grounded in the anatomical area depicted in the left figure. The different colors correspond to distinct anatomical regions.

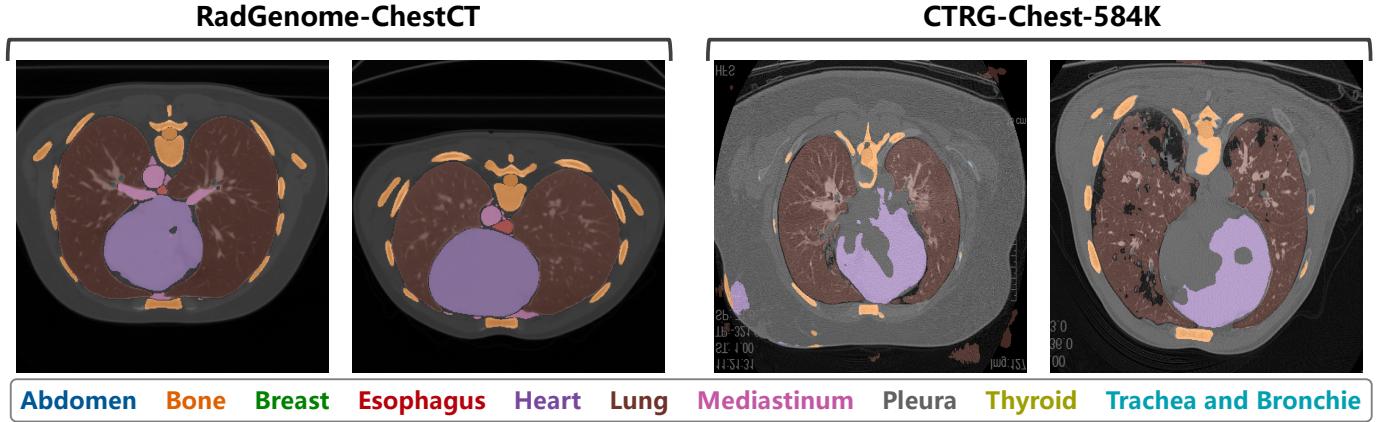


Fig. 6. Segmentation results on the RadGenome-ChestCT and CTRG-Chest-584K datasets. The different colors represent distinct anatomical areas, as shown at the bottom.

valuable for improving the reliability of the generated reports. Given that the relative positions of anatomical regions are generally consistent across patients, incorporating anatomical spatial priors can help refine the anatomical masks. This prior anatomical knowledge enhances the model’s robustness to segmentation failures, thereby ensuring the quality of the extracted local features. Moreover, we plan to explore integrating segmentation model optimization into the training process to further enhance the quality of the anatomical masks, thereby better supporting diagnosis.

Another limitation is the insufficiency of detailed lesion information. The current local features in our framework are restricted to organ-level information, which is inadequate for precise lesion detection. As shown in Fig. 4, the model may miss certain lesions or misidentify their locations. To enhance the model’s ability to detect and characterize lesions, we plan to explore integrating lesion segmentation or detection into the framework. Incorporating lesion-specific information will enable the model to capture more detailed insights into abnormalities, thereby facilitating the generation of more accurate and clinically relevant reports.

VI. CONCLUSION

In this study, we propose the **Reg2RG** framework for CT report generation. Unlike existing methods relying only on global features, our approach integrates local features with global features, improving the model’s ability to identify detailed lesions in sub-regions, while also enhancing the model’s interpretability and the reliability of the reports. Specifically, we use anatomical masks from a universal segmentation model to capture local features of referring regions. To retain high-resolution local details with low computational cost, a local feature decouple strategy (LFD) is introduced to decouple local features into two parts. Texture features capture fine-grained details of cropped region volumes, while geometric features encode position and size information lost during cropping. The global features are also incorporated to achieve a holistic volume representation. Through the collaboration of global and local features, the model effectively captures the

inter-regional relationships while preserving detailed insights within each region. To improve referring and grounding, we propose a training strategy RRA that uses region recognition to guide region-specific report generation. This strategy enhances interpretability and reliability by ensuring reports are grounded in the correct regions. Extensive experiments on two large-scale chest CT datasets demonstrate the superiority of our model over compared methods in both NLG and CE metrics.

Our work propels CT report generation with a region-guided mechanism, enhancing the trustworthiness of the generated reports. In the future, we plan to extend to additional imaging modalities and anatomical regions while integrating more fine-grained information, aiming to provide a versatile and efficient solution for automated report generation in radiology.

REFERENCES

- [1] S. Hussain, I. Mubeen, N. Ullah, S. S. U. D. Shah, B. A. Khan, M. Zahoor, R. Ullah, F. A. Khan, and M. A. Sultan, “Modern diagnostic imaging technique applications and risk factors in the medical field: a review,” *BioMed research international*, vol. 2022, no. 1, p. 5164970, 2022.
- [2] S. K. Goergen, F. J. Pool, T. J. Turner, J. E. Grimm, M. N. Appleyard, C. Crock, M. C. Fahey, M. F. Fay, N. J. Ferris, S. M. Liew *et al.*, “Evidence-based guideline for the written radiology report: Methods, recommendations and implementation challenges,” *Journal of medical imaging and radiation oncology*, vol. 57, no. 1, pp. 1–7, 2013.
- [3] I. E. Hamamci, S. Er, and B. Menze, “Ct2rep: Automated radiology report generation for 3d medical imaging,” *arXiv preprint arXiv:2403.06801*, 2024.
- [4] C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, “Towards generalist foundation model for radiology,” *arXiv preprint arXiv:2308.02463*, 2023.
- [5] F. Bai, Y. Du, T. Huang, M. Q.-H. Meng, and B. Zhao, “M3d: Advancing 3d medical image analysis with multi-modal large language models,” *arXiv preprint arXiv:2404.00578*, 2024.
- [6] Z. Zhao, Y. Zhang, C. Wu, X. Zhang, Y. Zhang, Y. Wang, and W. Xie, “One model to rule them all: Towards universal segmentation for medical images with text prompts,” *arXiv preprint arXiv:2312.17183*, 2023.
- [7] J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang, “Segment anything in medical images,” *Nature Communications*, vol. 15, no. 1, p. 654, 2024.
- [8] R. Stock, Y. Kirchhoff, M. R. Rokuss, A. Ravindran, and K. Maier-Hein, “Segment anything in medical images with mnunet,” in *CVPR 2024: Segment Anything In Medical Images On Laptop*.
- [9] C. Liu, Z. Wan, Y. Wang, H. Shen, H. Wang, K. Zheng, M. Zhang, and R. Arcucci, “Benchmarking and boosting radiology report generation for 3d high-resolution medical images,” *arXiv preprint arXiv:2406.07146*, 2024.

- [10] S. Li, P. Qiao, L. Wang, M. Ning, L. Yuan, Y. Zheng, and J. Chen, “An organ-aware diagnosis framework for radiology report generation,” *IEEE Transactions on Medical Imaging*, 2024.
- [11] Z. Chen, Y. Song, T.-H. Chang, and X. Wan, “Generating radiology reports via memory-driven transformer,” *arXiv preprint arXiv:2010.16056*, 2020.
- [12] Z. Chen, Y. Shen, Y. Song, and X. Wan, “Cross-modal memory networks for radiology report generation,” *arXiv preprint arXiv:2204.13258*, 2022.
- [13] H. Jin, H. Che, Y. Lin, and H. Chen, “Promptmrg: Diagnosis-driven prompts for medical report generation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 3, 2024, pp. 2607–2615.
- [14] T. Tanida, P. Müller, G. Kaassis, and D. Rueckert, “Interactive and explainable region-guided radiology report generation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7433–7442.
- [15] I. E. Hamamci, S. Er, F. Almas, A. G. Simsek, S. N. Esirgun, I. Dogan, M. F. Dasdelen, B. Wittmann, E. Simsar, M. Simsar *et al.*, “A foundation model utilizing chest ct volumes and radiology reports for supervised-level zero-shot detection of abnormalities,” *arXiv preprint arXiv:2403.17834*, 2024.
- [16] X. Zhang, C. Wu, Z. Zhao, J. Lei, Y. Zhang, Y. Wang, and W. Xie, “Radgenome-chest ct: A grounded vision-language dataset for chest ct analysis,” *arXiv preprint arXiv:2404.16754*, 2024.
- [17] Y. Tang, H. Yang, L. Zhang, and Y. Yuan, “Work like a doctor: Unifying scan localizer and dynamic generator for automated computed tomography report generation,” *Expert Systems with Applications*, vol. 237, p. 121442, 2024.
- [18] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *International conference on machine learning*. PMLR, 2022, pp. 12 888–12 900.
- [19] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” *Advances in neural information processing systems*, vol. 36, 2024.
- [20] Z. Chen, L. Luo, Y. Bie, and H. Chen, “Dia-llama: Towards large language model-driven ct report generation,” *arXiv preprint arXiv:2403.16386*, 2024.
- [21] C.-Y. Li, K.-J. Chang, C.-F. Yang, H.-Y. Wu, W. Chen, H. Bansal, L. Chen, Y.-P. Yang, Y.-C. Chen, S.-P. Chen *et al.*, “Towards a holistic framework for multimodal large language models in three-dimensional brain ct report generation,” *arXiv preprint arXiv:2407.02235*, 2024.
- [22] X. Zhang, C. Wu, Z. Zhao, W. Lin, Y. Zhang, Y. Wang, and W. Xie, “Pmc-vqa: Visual instruction tuning for medical visual question answering,” *arXiv preprint arXiv:2305.10415*, 2023.
- [23] W. Wang, Z. Chen, X. Chen, J. Wu, X. Zhu, G. Zeng, P. Luo, T. Lu, J. Zhou, Y. Qiao *et al.*, “Visionllm: Large language model is also an open-ended decoder for vision-centric tasks,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [24] S. Zhang, P. Sun, S. Chen, M. Xiao, W. Shao, W. Zhang, Y. Liu, K. Chen, and P. Luo, “Gpt4roi: Instruction tuning large language model on region-of-interest,” *arXiv preprint arXiv:2307.03601*, 2023.
- [25] C. Ma, Y. Jiang, J. Wu, Z. Yuan, and X. Qi, “Groma: Localized visual tokenization for grounding multimodal large language models,” *arXiv preprint arXiv:2404.13013*, 2024.
- [26] Q. Guo, S. De Mello, H. Yin, W. Byeon, K. C. Cheung, Y. Yu, P. Luo, and S. Liu, “Regionopt: Towards region understanding vision language model,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 13 796–13 806.
- [27] H. You, H. Zhang, Z. Gan, X. Du, B. Zhang, Z. Wang, L. Cao, S.-F. Chang, and Y. Yang, “Ferret: Refer and ground anything anywhere at any granularity,” *arXiv preprint arXiv:2310.07704*, 2023.
- [28] S. Bannur, K. Bouzid, D. C. Castro, A. Schwaighofer, S. Bond-Taylor, M. Ilse, F. Pérez-García, V. Salvatelli, H. Sharma, F. Meissen *et al.*, “Maira-2: Grounded radiology report generation,” *arXiv preprint arXiv:2406.04449*, 2024.
- [29] A. Alkhaldi, R. Alnajim, L. Alabdullatef, R. Alyahya, J. Chen, D. Zhu, A. Alsinan, and M. Elhoseiny, “Minigpt-med: Large language model as a general interface for radiology diagnosis,” *arXiv preprint arXiv:2407.04106*, 2024.
- [30] X. Huang, H. Huang, L. Shen, Y. Yang, F. Shang, J. Liu, and J. Liu, “A refer-and-ground multimodal large language model for biomedicine,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2024, pp. 399–409.
- [31] L. Wang, H. Wang, H. Yang, J. Mao, Z. Yang, J. Shen, and X. Li, “Interpretable bilingual multimodal large language model for diverse biomedical tasks,” *arXiv preprint arXiv:2410.18387*, 2024.
- [32] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [33] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, G. Dong, H. Wei, H. Lin, J. Tang, J. Wang, J. Yang, J. Tu, J. Zhang, J. Ma, J. Xu, J. Zhou, J. Bai, J. He, J. Lin, K. Dang, K. Lu, K. Chen, K. Yang, M. Li, M. Xue, N. Ni, P. Zhang, P. Wang, R. Peng, R. Men, R. Gao, R. Lin, S. Wang, S. Bai, S. Tan, T. Zhu, T. Li, T. Liu, W. Ge, X. Deng, X. Zhou, X. Ren, X. Zhang, X. Wei, X. Ren, Y. Fan, Y. Yao, Y. Zhang, Y. Wan, Y. Chu, Y. Liu, Z. Cui, Z. Zhang, and Z. Fan, “Qwen2 technical report,” *arXiv preprint arXiv:2407.10671*, 2024.
- [34] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.
- [35] S. Banerjee and A. Lavie, “Meteor: An automatic metric for mt evaluation with improved correlation with human judgments,” in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.
- [36] C.-Y. Lin, “Rouge: A package for automatic evaluation of summaries,” in *Text summarization branches out*, 2004, pp. 74–81.
- [37] A. Yan, J. McAuley, X. Lu, J. Du, E. Y. Chang, A. Gentili, and C.-N. Hsu, “Radbert: adapting transformer-based language models to radiology,” *Radiology: Artificial Intelligence*, vol. 4, no. 4, p. e210258, 2022.
- [38] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [39] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [40] I. Loshchilov, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [41] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, “Zero: Memory optimizations toward training trillion parameter models,” in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2020, pp. 1–16.
- [42] T. Chen, B. Xu, C. Zhang, and C. Guestrin, “Training deep nets with sublinear memory cost,” *arXiv preprint arXiv:1604.06174*, 2016.
- [43] Z. Wang, L. Liu, L. Wang, and L. Zhou, “R2gengpt: Radiology report generation with frozen llms,” *Meta-Radiology*, vol. 1, no. 3, p. 100033, 2023.
- [44] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [45] A. Nicolson, J. Dowling, and B. Koopman, “Improving chest x-ray report generation by leveraging warm starting,” *Artificial intelligence in medicine*, vol. 144, p. 102633, 2023.
- [46] O. Alfarghaly, R. Khaled, A. Elkorany, M. Helal, and A. Fahmy, “Automated radiology report generation using conditioned transformers,” *Informatics in Medicine Unlocked*, vol. 24, p. 100557, 2021.