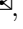# Radiologist Copilot: An Agentic Assistant with Orchestrated Tools for Radiology Reporting with Quality Control

Yongrui Yu[1], Zhongzhen Huang[1], Linjie Mu[1], Shaoting Zhang[1,2]✉, and Xiaofan Zhang[1,3]✉

[1] Shanghai Jiao Tong University, Shanghai, China
[2] SenseTime Research, Shanghai, China
[3] Shanghai Innovation Institute, Shanghai, China
zhangshaoting@sensetime.com; xiaofan.zhang@sjtu.edu.cn

**Abstract.** Radiology reporting is an essential yet time-consuming and error-prone task for radiologists in clinical examinations, especially for volumetric medical images. Rigorous quality control is also critical but tedious, ensuring that the final report meets clinical standards. Existing automated approaches, including radiology report generation methods and medical vision-language models, focus mainly on the report generation phase and neglect the crucial quality control procedure, limiting their capability to provide comprehensive support to radiologists. We propose Radiologist Copilot, an agentic AI assistant equipped with orchestrated tools designed for automated radiology reporting with quality control. Leveraging large language models as the reasoning backbone, the agentic system autonomously selects tools, plans, and executes actions, emulating the behavior of radiologists throughout the holistic radiology reporting process. The orchestrated tools include region localization, think with image paradigm directed region analysis planning, strategic template selection for report generation, quality assessment and feedback-driven adaptive refinement for quality control. Therefore, Radiologist Copilot facilitates accurate, complete, and efficient radiology reporting, assisting radiologists and improving clinical efficiency. Experimental results demonstrate that Radiologist Copilot significantly surpasses other state-of-the-art methods in radiology reporting. The source code will be released upon acceptance.

**Keywords:** Medical Agent · Radiology Reporting · Quality Control.

## 1 Introduction

Radiology reporting is an essential component of radiological examinations, which requires radiologists to meticulously interpret medical images and to produce a detailed report that summarizes key findings and impressions. However, radiology reporting is often time-consuming and prone to errors, particularly for 3D medical imaging modalities such as CT and MRI [7]. In addition, rigorous

quality control is also a critical step in radiology reporting to ensure that the final report is accurate, complete, and free of errors [20]. This verification procedure is tedious, involving detailed checks of content, structure, terminology, and other elements. Therefore, automated radiology reporting with quality control is of great significance for assisting radiologists and improving clinical efficiency.

Previous methods for automated radiology reporting of volumetric medical imaging, such as radiology report generation (RRG) approaches like CT2Rep [7], have made significant progress. More recently, medical vision-language models (VLMs) [22,8,10] and agentic methods [15] have been applied to tasks such as radiology report generation and visual question answering (VQA). However, these approaches cannot fully support radiologists in accomplishing the radiology reporting process, as report generation represents only one phase of the holistic process and lacks the critical quality control procedure.

We propose Radiologist Copilot, an AI agentic assistant with orchestrated tools for automated radiology reporting with quality control. Radiologist Copilot is designed to assist radiologists throughout the entire radiology reporting process, encompassing medical image analysis, report generation, and quality control. The agentic assistant, utilizing large language models (LLMs) as the reasoning backbone, autonomously selects appropriate tools from its tool library, plans and executes actions, thereby emulating the behavior of radiologists to complete the radiology reporting process with quality control. The orchestrated tools comprise a segmentation model for organ and lesion localization, combined with the *Think with Image* paradigm [19] for region analysis planning and investigation of region-of-interest images. We also implement strategic template selection to facilitate report generation, perform quality assessment of the radiology reports, and provide feedback for feedback-driven adaptive refinement. Together, these tools enable the simulation of radiologists' behavior and support automated radiology reporting.

In conclusion, our main contributions are summarized as follows:

- We propose Radiologist Copilot, which assists radiologists in achieving automated radiology reporting with quality control through an agentic framework with orchestrated tools.
- The agentic system, leveraging a reasoning backbone, automatically selects tools, plans, and executes actions, emulating the behavior of radiologists to accomplish the complete radiology reporting process.
- The tool library, consisting of region localization, region analysis planning, strategic template selection, quality assessment, and feedback-driven adaptive refinement, enables comprehensive radiology reporting.
- The experimental results demonstrate that Radiologist Copilot surpasses other state-of-the-art methods in radiology reporting by a large margin.

## 2   Related Work

Automated radiology reporting with quality control is crucial for supporting radiologists, since the process for 3D medical imaging modalities is extremely

time-consuming. CT2Rep [7] utilizes an auto-regressive causal transformer architecture and relational memory for radiology report generation of 3D chest CT images. Besides, Reg2RG [4] leverages a region-guided referring and grounding framework for chest CT report generation. Nevertheless, report generation constitutes only one stage of the entire radiology reporting process and lacks the critical quality control procedures.

Medical VLMs have been proposed to address medical tasks such as RRG and VQA. In particular, 3D medical VLMs have shown effectiveness for 3D medical modalities, including CT and MRI. RadFM [22] introduces a generalist vision-language foundation model that unifies both 2D and 3D medical data for radiology. M3D [1] presents M3D-Data, a large-scale multi-modal medical dataset, together with M3D-LaMed, a multi-modal LLM for 3D medical image analysis. Merlin [3] develops a 3D VLM that leverages both structured and unstructured data for abdominal CT interpretation. CT-CHAT [8] is a vision-language foundation model for 3D chest CT volumes, combining the CT-CLIP vision encoder with a pretrained LLM. Med3DVLM [23] is a 3D VLM that enhances multi-modal representations via an efficient encoder, a contrastive learning strategy, and a dual-stream projector. Hulu-Med [10] provides a generalist medical VLM that unifies text, 2D and 3D images, and video within a single architecture.

The advent of AI agents enables autonomous reasoning, planning, and tool usage. Recently, medical agents have been introduced for complex medical tasks. For example, MMedAgent [12] presents a multi-modal medical agent that leverages various tools to handle medical tasks across different modalities, such as radiology report generation for chest X-rays (CXR). MedRAX [6] provides a specialized agent framework for chest X-ray interpretation that integrates state-of-the-art CXR analysis tools. CT-Agent [15] is an anatomy-aware and token-efficient agent for 3D CT VQA, enabling effective reasoning. However, for the RRG task, these agentic methods typically rely on a single tool and lack cooperation among tools. In contrast, our Radiologist Copilot provides a diverse and comprehensive toolset, enables rich inter-tool collaboration, and supports substantially more complex tool planning and execution.

## 3    Methodology

Radiology reporting is essential but time-consuming and error-prone, particularly for 3D medical imaging. Moreover, rigorous quality control is a critical yet tedious step to ensure that the final report is clinically qualified, which is often overlooked in previous studies. To address these challenges, we propose Radiologist Copilot, an AI agentic assistant equipped with an orchestrated set of tools for automated radiology reporting with quality control, to support radiologists and enhance clinical efficiency.

### 3.1    Radiologist Copilot

The proposed Radiologist Copilot, as shown in Fig. 1, is designed to assist radiologists throughout the holistic radiology reporting process, comprising medical
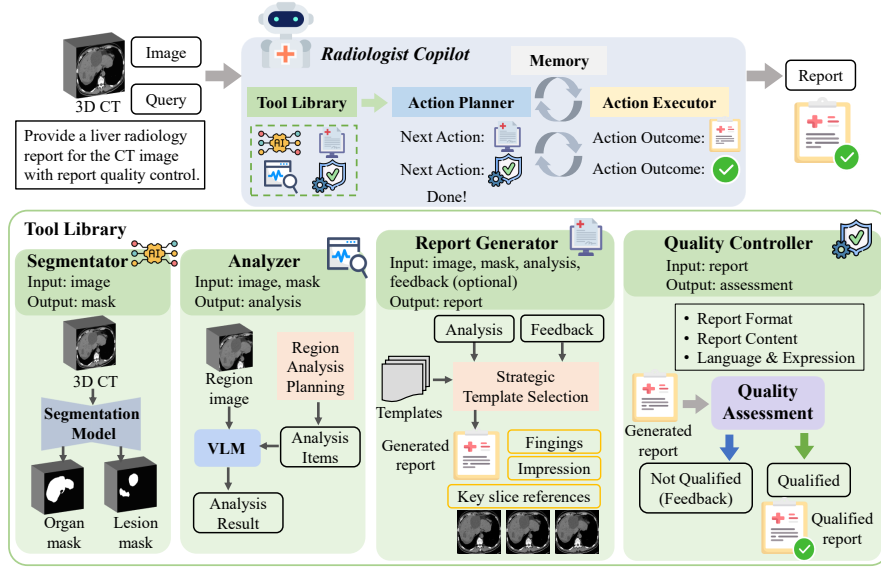
Fig. 1: An overview of the proposed Radiologist Copilot with an agent framework and orchestrated tools for radiology reporting with quality control.

image analysis, report generation, and quality control. Given a user query $Q$ and a 3D CT image $I \in \mathbb{R}^{C \times H \times W \times D}$, where $C$, $H$, $W$, and $D$ denote the channel, height, width, and depth, respectively, the Radiologist Copilot, as an agentic assistant, receives these inputs and outputs a qualified radiology report $R$.

The agentic assistant, leveraging a large language model $\text{LLM}_\theta(\cdot)$ as its reasoning backbone, performs inference directly without additional training. With a curated set of tools in its library, the agentic assistant employs an action planner and an action executor to automatically select and coordinate appropriate tools to complete tasks. The orchestrated tools $\mathcal{T} = \{T_i\}_{i=1}^{n}$, where $n$ denotes the number of tools in the tool library, comprise region image localization, region analysis planning for region images, strategic template selection to facilitate report generation, and quality control for radiology reports. These tools cooperatively support automated radiology reporting with quality control, emulating the behaviors of radiologists.

### 3.2 Agentic Assistant

The agentic assistant first initializes the tool library and then iteratively utilizes the action planner and the action executor to plan and execute actions. The memory module records these actions and their outcomes, enabling the verification of task completion.
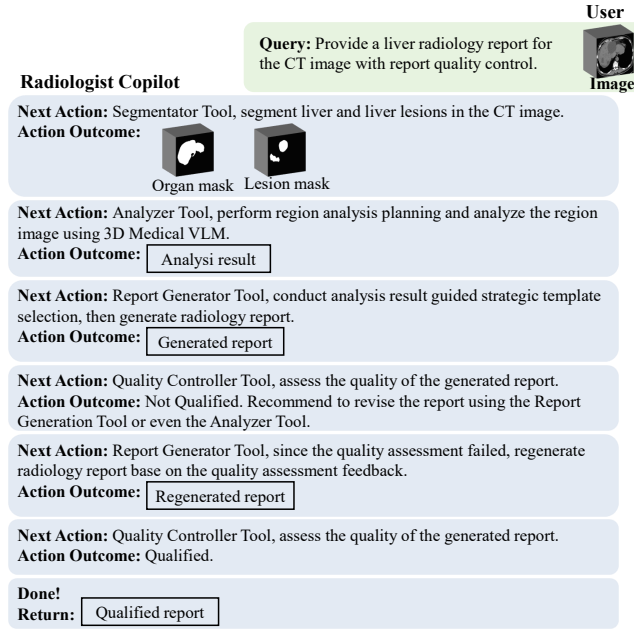
Fig. 2: An illustrative workflow of the Radiologist Copilot.

**Action Planner.** The action planner first analyzes the query task and provides a high-level guideline for subsequent action planning. It then predicts the next action by determining the intermediate goal and selecting the suitable tool $T_i$. At the end of each iteration, the action planner accesses the memory to verify whether the task has been completed, thereby deciding whether to continue or stop.

**Action Executor.** The action executor is responsible for generating and executing commands based on the next action determined by the action planner. For each next action, it generates the corresponding command for the selected tool, carries out the command, and derives the action outcome.

### 3.3 Orchestrated Tools

The tool library comprises the Segmentator Tool for region localization, the Analyzer Tool with the *Think with Image* paradigm [19] for region analysis planning, the Report Generator Tool utilizing strategic template selection, and the Quality Controller Tool, which performs quality assessment of radiology reports and provides feedback for feedback-driven adaptive refinement. These tools support the entire process of automated radiology reporting. Figure 2 provides an illustrative workflow of the Radiologist Copilot.

**Segmentator.** The Segmentator Tool takes a 3D CT image $I$ as input and outputs segmentation masks $M_{organ}$ and $M_{lesion}$ for target organs and lesions relevant to the user query, using pretrained segmentation models, such as TotalSegmentator [21]. This tool facilitates subsequent region localization and region analysis in accordance with the *Think with Image* paradigm.

**Analyzer.** The Analyzer Tool integrates the CT image $I$ with the organ mask $M_{organ}$ to localize and extract the region-of-interest image, denoted as $I_{region}$. We further propose a novel **Region Analysis Planning (RAP)** module to define specific analysis items for this region, covering organ and lesion characteristics such as size, shape, and density. RAP dynamically determines whether lesion-related characteristics should be analyzed based on the presence of lesions in the lesion mask $M_{lesion}$. These analysis items are then formulated as prompts for 3D medical VLMs, which analyze the region image $I_{region}$ according to these items and generate the corresponding analysis result.

**Report Generator.** In the Report Generator Tool, we introduce **Strategic Template Selection (STS)**. Given some template reports, the LLM selects the most relevant template based on the analysis result. Using the chosen template as a reference, the report is generated according to the analysis result, including findings and impression. Therefore, report generation leverages existing template reports while adapting to the current analysis. The optional feedback information is utilized to direct the revision of previously generated reports that do not satisfy quality standards. Since formal radiology reports often include key slices of the CT image as references, we select the central three slices along the axial plane of the organ for normal cases and of the largest lesion for abnormal cases, based on the segmentation masks. Consequently, the generated report $R_{generated}$ is composed of three components: the findings section, the impression section, and the key slice references.

**Quality Controller.** The Quality Controller is utilized to ensure that the radiology report meets clinical standards. We perform **Quality Assessment** on the generated report $R_{generated}$ using the LLM, and output assessment comments indicating whether it is qualified. If the report is qualified, it is considered the qualified report $R_{qualified}$; if the report is not qualified, feedback information is provided to enable feedback-driven adaptive refinement of $R_{generated}$. The quality assessment performs a thorough validation of the radiology report, covering its format, content, language, and expression. Specifically, the report format is assessed to ensure that the findings provide an objective description of imaging manifestations and the impression presents a diagnostic summary highlighting key conclusions in order of importance. The content is evaluated for anatomical correctness, lesion characterization, and consistency between findings and impression. Language and expression are examined to ensure the use of standardized radiological terminology, correct spelling, clarity, conciseness, and avoidance of redundancy or vague statements.

# 4    Experiments and Results

## 4.1    Experimental Setup

**Datasets.** We leverage the liver radiology reporting task to validate the effectiveness of our proposed Radiologist Copilot. We utilize the publicly available 3D CT dataset AMOS-MM [9] and its medical report generation task, which includes more than 2,000 CT scans with comprehensive radiology reports. This dataset contains 1,287 CT scans with reports for training and 400 CT scans with reports for validation. For the liver radiology reporting task, we filter the reports in the AMOS-MM dataset to identify those containing liver descriptions and extract these descriptions as liver radiology reports. The resulting training set consists of 1,149 CT scans with liver reports, and the validation set contains 367 CT scans with liver reports. For liver reports, the validation set is reserved to evaluate our method, while the training set is used to conduct report analysis. We first obtain report embeddings using BioBERT [11], and then perform K-means clustering on these embeddings to derive template reports. These template reports are summarized into liver analysis items, including liver surface, liver parenchyma, bile ducts, and liver lesions, which are used in the Analyzer Tool.

**Implementation Details.** Our Radiologist Copilot is implemented based on the agentic framework OctoTools [14] using Python 3.10.18, and experiments are conducted on NVIDIA L20 GPUs. The system operates in a training-free manner using pretrained models. We utilize Qwen3-32B [24] as the LLM backbone for agent reasoning. The Segmentator Tool employs TotalSegmentator [21] as the segmentation model, while the Analyzer Tool leverages the 3D medical VLM Hulu-Med [10] for CT analysis. The maximum number of steps is set to 10, and the maximum execution time is limited to 500 seconds per case.

**Evaluation Metrics.** To thoroughly evaluate the effectiveness of the Radiologist Copilot, we conduct both agent-level and task-level evaluations. For task-level evaluation of radiology reporting, we adopt natural language generation (NLG) and clinical efficacy (CE) metrics. The NLG metrics include BLEU-1 [18], ROUGE-L [13], METEOR [2], and BERTScore [25]. The CE metrics are used to measure clinical accuracy, including F1-RadGraph [5] and GREEN [17]. For agent-level evaluation of agentic tools, we use LLM-as-a-Judge [26], adopting OpenAI GPT-5.1 [16] as the LLM judge. The evaluation assesses the entire process that the agent handles the task, considering four dimensions (Analysis Process, Tool Selection, Action Planning, and Action Execution) with scores ranging from 1 to 5 (Poor, Fair, Moderate, Good, Excellent). Analysis Process evaluates whether the query analysis is thorough, logical, and complete; Tool Selection assesses whether the agent selects tools accurately and appropriately; Action Planning evaluates the reasonableness of the planned actions; and Action Execution assesses whether the agent successfully executes the planned actions.
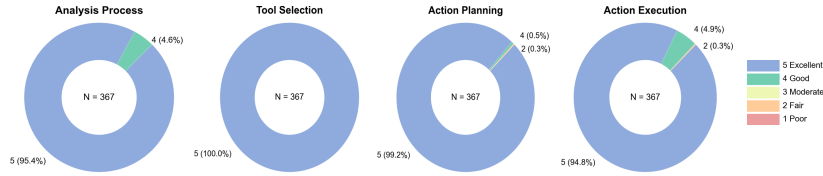
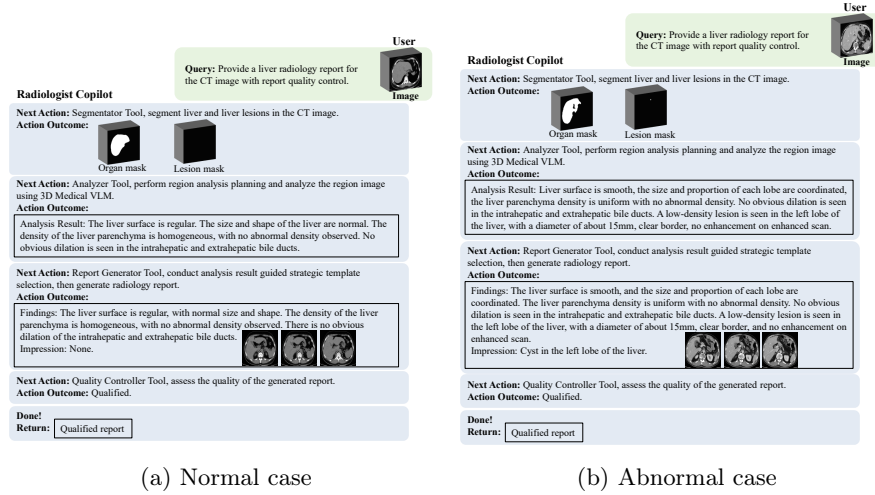Fig. 3: Agent-level evaluation of the Radiologist Copilot using LLM-as-a-Judge.



(a) Normal case                    (b) Abnormal case

Fig. 4: Examples of the Radiologist Copilot workflow.

## 4.2   Experimental Results

To comprehensively evaluate the effectiveness of Radiologist Copilot, we perform both agent-level evaluations of the agentic tools and task-level evaluations of liver radiology reporting, along with ablation studies on its core components.

**Agent-level Evaluation.** To validate the overall process that the agent deals with tasks, Fig. 3 presents the LLM-as-a-Judge assessments. The figure shows the score distribution across four dimensions, with the majority of scores being 5 (Excellent), and most of the remaining scores being 4 (Good). These results demonstrate that the Analysis Process of the agent is thorough, logical, and complete; the Tool Selection is accurate and appropriate; the Action Planning is reasonable; and the Action Execution proceeds successfully. Furthermore, Fig. 4 illustrates the workflow of Radiologist Copilot in handling user queries, showing that the processes of action planning and action execution are coherent and efficient, ultimately generating a qualified radiology report. The entire process simulates the practices of radiologists and enhances clinical efficiency.

Table 1: Task-level evaluation on the liver radiology reporting of our proposed 3D medical agent compared with 3D medical VLMs.

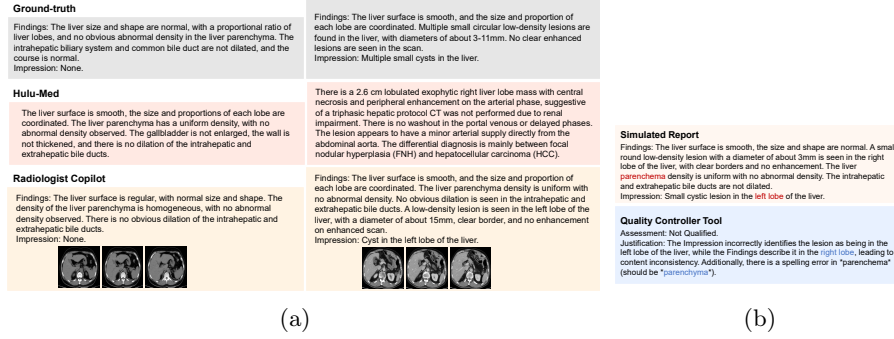| Method | NLG Metrics | | | | CE Metrics | |
|---|---|---|---|---|---|---|
| | BLEU-1 | ROUGE-L | METEOR | BERTScore | F1-RadGraph | GREEN |
| **3D Medical VLM** | | | | | | |
| RadFM [22] | 0.1492 | 0.1415 | 0.2340 | 0.5541 | 0.0686 | 0.0353 |
| M3D [1] | 0.1775 | 0.1302 | 0.1220 | 0.5359 | 0.0475 | 0.0209 |
| Merlin [3] | 0.0015 | 0.0908 | 0.0569 | 0.5119 | 0.1617 | 0.1024 |
| CT-CHAT [8] | 0.2440 | 0.2012 | 0.2599 | 0.6127 | 0.1196 | 0.0390 |
| Med3DVLM [23] | 0.1967 | 0.1422 | 0.1847 | 0.5608 | 0.0660 | 0.0539 |
| Hulu-Med [10] | 0.1867 | 0.1723 | 0.2380 | 0.5947 | 0.1209 | 0.2163 |
| **3D Medical Agent** | | | | | | |
| Ours | **0.4025** | **0.3222** | **0.4560** | **0.7024** | **0.2585** | **0.4379** |



Fig. 5: (a) Case study of Hulu-Med and Radiologist Copilot. (b) The validation of the Quality Controller Tool.

**Task-level Evaluation.** We compare the liver radiology reporting performance of our proposed method with state-of-the-art (SOTA) approaches. The results illustrated in Table 1 demonstrate that Radiologist Copilot substantially outperforms other SOTA methods on both NLG and CE metrics, highlighting the practical effectiveness of this training-free agentic assistant. Fig. 5a shows generated liver reports, including normal and abnormal cases. The results indicate that the reports generated by Radiologist Copilot are highly consistent with the ground-truth in both format and content, significantly surpassing Hulu-Med [10].

**Ablation Studies.** We perform ablation studies on the key components of Radiologist Copilot, with results presented in Table 2. The results highlight the importance of region analysis planning and strategic template selection, as removing either component leads to a noticeable degradation of generated liver reports. Although removing strategic template selection slightly improves the

Table 2: Ablation study on the components of our proposed Radiologist Copilot. RAP, STS, and QC denote region analysis planning, strategic template selection, and quality control, respectively.

| Method | BLEU-1 | ROUGE-L | METEOR | BERTScore | F1-RadGraph | GREEN |
|---|---|---|---|---|---|---|
| Ours | **0.4025** | **0.3222** | **0.4560** | **0.7024** | **0.2585** | 0.4379 |
| Ours (w/o RAP) | 0.3600 | 0.2588 | 0.3610 | 0.6469 | 0.1675 | 0.2269 |
| Ours (w/o STS) | 0.2983 | 0.2698 | 0.3915 | 0.6553 | 0.2429 | **0.4582** |
| Ours (w/o QC) | 0.3998 | 0.3149 | 0.4462 | 0.6944 | 0.2545 | 0.4360 |

Table 3: Ablation study on Radiologist Copilot equipped with different VLMs.

| Method | BLEU-1 | ROUGE-L | METEOR | BERTScore | F1-RadGraph | GREEN |
|---|---|---|---|---|---|---|
| RadFM | 0.1492 | 0.1415 | 0.2340 | 0.5541 | 0.0686 | 0.0353 |
| Ours (RadFM) | 0.3215 | 0.2465 | 0.3633 | 0.6510 | 0.1994 | 0.3719 |
| CT-CHAT | 0.2440 | 0.2012 | 0.2599 | 0.6127 | 0.1196 | 0.0390 |
| Ours (CT-CHAT) | 0.3671 | 0.2784 | 0.3738 | 0.6558 | 0.2312 | 0.1485 |
| Hulu-Med | 0.1867 | 0.1723 | 0.2380 | 0.5947 | 0.1209 | 0.2163 |
| Ours (Hulu-Med) | 0.4025 | 0.3222 | 0.4560 | 0.7024 | 0.2585 | 0.4379 |

GREEN metric, possibly due to the inclusion of more information, other metrics decrease significantly. From the results in the table, the impact of quality control is less prominent compared to other components, as reports generated through the entire process of Radiologist Copilot generally meet the required standards. Nevertheless, the quality control procedure plays an important role, offering practical assurance that the generated reports adhere to clinical standards. Therefore, we provide a separate assessment of the Quality Controller Tool, as shown in Fig. 5b, which effectively detects issues in the simulated report, including content inconsistencies and spelling errors.

In addition, we conduct an ablation study on the Radiologist Copilot equipped with different medical VLMs in the Analyzer Tool. The results shown in Table 3 indicate that even when equipped with different VLMs, Radiologist Copilot consistently maintains strong radiology reporting capabilities, demonstrating the superiority of our proposed agentic assistant.

## 5   Discussion and Conclusion

We validate our proposed Radiologist Copilot on the CT dataset, demonstrating strong performance in generating liver radiology reports. It is feasible to generate chest or abdomen radiology reports, since the agent integrates 3D VLMs that enable CT image analysis across different anatomical regions. In addition, the agent leverages LLMs as the reasoning backbone, flexibly selecting appro-

priate LLMs based on task complexity, model capability, and resource cost to accomplish tasks efficiently.

In conclusion, Radiologist Copilot introduces an agentic AI assistant with orchestrated tools for radiology reporting with quality control. The agentic system automatically plans and executes actions using its toolset, which includes region localization, region analysis planning, strategic template selection, quality assessment, and feedback-driven adaptive refinement for the holistic radiology reporting process. Experimental results demonstrate that Radiologist Copilot enables accurate, comprehensive, and efficient report generation, highlighting its potential to assist radiologists and improve clinical efficiency.

# References

1. Bai, F., Du, Y., Huang, T., Meng, M.Q.H., Zhao, B.: M3d: Advancing 3d medical image analysis with multi-modal large language models. arXiv preprint arXiv:2404.00578 (2024)
2. Banerjee, S., Lavie, A.: 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. pp. 65–72 (2005)
3. Blankemeier, L., Cohen, J.P., Kumar, A., Van Veen, D., Gardezi, S.J.S., Paschali, M., Chen, Z., Delbrouck, J.B., Reis, E., Truyts, C., et al.: Merlin: A vision language foundation model for 3d computed tomography. Research Square pp. rs–3 (2024)
4. Chen, Z., Bie, Y., Jin, H., Chen, H.: Large language model with region-guided referring and grounding for ct report generation. IEEE Transactions on Medical Imaging (2025)
5. Delbrouck, J.B., Chambon, P., Chen, Z., Varma, M., Johnston, A., Blankemeier, L., Van Veen, D., Bui, T., Truong, S., Langlotz, C.: Radgraph-xl: A large-scale expert-annotated dataset for entity and relation extraction from radiology reports. In: Findings of the Association for Computational Linguistics: ACL 2024. pp. 12902–12915 (2024)
6. Fallahpour, A., Ma, J., Munim, A., Lyu, H., Wang, B.: Medrax: Medical reasoning agent for chest x-ray. arXiv preprint arXiv:2502.02673 (2025)
7. Hamamci, I.E., Er, S., Menze, B.: Ct2rep: Automated radiology report generation for 3d medical imaging. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 476–486. Springer (2024)
8. Hamamci, I.E., Er, S., Wang, C., Almas, F., Simsek, A.G., Esirgun, S.N., Dogan, I., Durugol, O.F., Hou, B., Shit, S., et al.: Developing generalist foundation models from a multimodal dataset for 3d computed tomography. arXiv preprint arXiv:2403.17834 (2024)
9. Ji, Y., Bai, H., Ge, C., Yang, J., Zhu, Y., Zhang, R., Li, Z., Zhanng, L., Ma, W., Wan, X., et al.: Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. Advances in neural information processing systems **35**, 36722–36732 (2022)
10. Jiang, S., Wang, Y., Song, S., Hu, T., Zhou, C., Pu, B., Zhang, Y., Yang, Z., Feng, Y., Zhou, J.T., et al.: Hulu-med: A transparent generalist model towards holistic medical vision-language understanding. arXiv preprint arXiv:2510.08668 (2025)

11. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: Biobert: a pre-trained biomedical language representation model for biomedical text mining. Bioinformatics **36**(4), 1234–1240 (2020)
12. Li, B., Yan, T., Pan, Y., Luo, J., Ji, R., Ding, J., Xu, Z., Liu, S., Dong, H., Lin, Z., et al.: Mmedagent: Learning to use medical tools with multi-modal agent. arXiv preprint arXiv:2407.02483 (2024)
13. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out. pp. 74–81 (2004)
14. Lu, P., Chen, B., Liu, S., Thapa, R., Boen, J., Zou, J.: Octotools: An agentic framework with extensible tools for complex reasoning. arXiv preprint arXiv:2502.11271 (2025)
15. Mao, Y., Xu, W., Qin, Y., Gao, Y.: Ct-agent: A multimodal-llm agent for 3d ct radiology question answering. arXiv preprint arXiv:2505.16229 (2025)
16. OpenAI: Gpt-5.1: A smarter, more conversational chatgpt. https://openai.com/zh-Hans-CN/index/gpt-5-1/ (2025)
17. Ostmeier, S., Xu, J., Chen, Z., Varma, M., Blankemeier, L., Bluethgen, C., Md, A.E.M., Moseley, M., Langlotz, C., Chaudhari, A.S., et al.: Green: Generative radiology report evaluation and error notation. In: Findings of the association for computational linguistics: EMNLP 2024. pp. 374–390 (2024)
18. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)
19. Su, Z., Xia, P., Guo, H., Liu, Z., Ma, Y., Qu, X., Liu, J., Li, Y., Zeng, K., Yang, Z., et al.: Thinking with images for multimodal reasoning: Foundations, methods, and future frontiers. arXiv preprint arXiv:2506.23918 (2025)
20. Warr, H., Ibrahim, Y., McGowan, D.R., Kamnitsas, K.: Quality control for radiology report generation models via auxiliary auditing components. In: International Workshop on Uncertainty for Safe Utilization of Machine Learning in Medical Imaging. pp. 70–80. Springer (2024)
21. Wasserthal, J., Breit, H.C., Meyer, M.T., Pradella, M., Hinck, D., Sauter, A.W., Heye, T., Boll, D.T., Cyriac, J., Yang, S., et al.: Totalsegmentator: robust segmentation of 104 anatomic structures in ct images. Radiology: Artificial Intelligence **5**(5), e230024 (2023)
22. Wu, C., Zhang, X., Zhang, Y., Hui, H., Wang, Y., Xie, W.: Towards generalist foundation model for radiology by leveraging web-scale 2d&3d medical data. Nature Communications **16**(1), 7866 (2025)
23. Xin, Y., Ates, G.C., Gong, K., Shao, W.: Med3dvlm: An efficient vision-language model for 3d medical image analysis. arXiv preprint arXiv:2503.20047 (2025)
24. Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al.: Qwen3 technical report. arXiv preprint arXiv:2505.09388 (2025)
25. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675 (2019)
26. Zheng, L., Chiang, W.L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al.: Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in neural information processing systems **36**, 46595–46623 (2023)