

Semantic Consistency-Based Uncertainty Quantification for Factuality in Radiology Report Generation

Chenyu Wang¹, Weichao Zhou², Shantanu Ghosh¹, Kayhan Batmanghelich¹, Wenchao Li¹

¹Boston University, ²Massachusetts Institute of Technology
{chyuwang, shawn24, batman, wenchao}@bu.edu, zhouw534@mit.edu

Abstract

Radiology report generation (RRG) has shown great potential in assisting radiologists by automating the labor-intensive task of report writing. While recent advancements have improved the quality and coherence of generated reports, ensuring their factual correctness remains a critical challenge. Although generative medical Vision Large Language Models (VLLMs) have been proposed to address this issue, these models are prone to hallucinations and can produce inaccurate diagnostic information. To address these concerns, we introduce a novel Semantic Consistency-Based Uncertainty Quantification framework that provides both report-level and sentence-level uncertainties. Unlike existing approaches, our method does not require modifications to the underlying model or access to its inner state, such as output token logits, thus serving as a plug-and-play module that can be seamlessly integrated with state-of-the-art models. Extensive experiments demonstrate the efficacy of our method in detecting hallucinations and enhancing the factual accuracy of automatically generated radiology reports. By abstaining from high-uncertainty reports, our approach improves factuality scores by 10%, achieved by rejecting 20% of reports using the Radialog model on the MIMIC-CXR dataset. Furthermore, sentence-level uncertainty flags the lowest-precision sentence in each report with an 82.9% success rate. Our implementation is open-source and available at <https://github.com/BU-DEPEND-Lab/SCUQ-RRG>.

1 Introduction

RRG is gaining importance as healthcare demands grow, placing substantial pressure on radiologists to interpret medical images swiftly and accurately. Automating the report-writing process holds the potential to alleviate this burden, improving both efficiency and diagnostic precision. Vision Large Language Models (VLLMs) have in-

troduced new possibilities in this area by generating detailed and coherent reports from medical images, providing significant assistance to radiologists (Thawkar et al., 2023; Pellegrini et al., 2023). However, despite these advancements, challenges persist—particularly in ensuring the factual accuracy of these generated reports. A notable issue with VLLMs is their tendency to produce “hallucinations”, or information that is ungrounded in the visual data or inconsistent with established medical knowledge. For example, a model might incorrectly generate findings such as a diagnosis of pneumonia when none is present (Hartsock and Rasool, 2024), or fabricate prior medical history that does not exist (Ramesh et al., 2022; Tanida et al., 2023; Hyland et al., 2023). Such hallucinations can lead to inaccurate or misleading diagnostic information, posing significant risks in clinical settings.

Recent studies have explored various methods to address hallucinations in radiology report generation. Ramesh et al. (2022) utilize a GPT-3-based rewriting technique and a BioBERT-based token classification system to remove references to non-existent prior reports. Banerjee et al. (2024) employ Direct Preference Optimization (DPO) to suppress hallucinated prior exams, significantly reducing such errors while maintaining clinical accuracy. However, these methods remain limited in scope, focusing solely on specific hallucinations, namely hallucinated prior exams, and do not enhance the broader factual accuracy of diverse clinical entities critical for dependable diagnostics. Bannur et al. (2023, 2024) tackle hallucinations by integrating current and prior images with detailed report sections, thereby improving the alignment between generated text and visual data to reduce errors and enhance report consistency. These approaches offer a more comprehensive solution than methods targeting specific hallucination types. However, they rely on specialized architectures and additional training resources, limiting their flexibility

and applicability across diverse models.

Addressing the limitations of prior approaches, our framework provides a plug-and-play solution that mitigates hallucinations through uncertainty quantification (UQ), requiring no architectural modifications or additional training. Broadly compatible with diverse VLLM-based RRG models, it emphasizes semantic consistency between generated content and sampled counterparts. Specifically, our UQ framework assesses the consistency of clinical entities within generated reports, assigning high uncertainty to content with low factual precision. We measure this consistency by comparing clinical facts from the original report with those in multiple sampled reports generated from the same query, relying solely on API-level access to broaden applicability. By abstaining from high-uncertainty reports, we enhance the clinical efficacy of generated outputs. Additionally, by flagging high-uncertainty sentences, we guide radiologists to areas needing further validation, reducing their workload and supporting more accurate interventions. In summary, our contributions are as follows:

- (1) We propose a plug-and-play UQ framework that does not require modifications to the internal mechanisms of the model and can be easily integrated with state-of-the-art RRG systems.
- (2) We propose two domain-specific uncertainty quantification methods for report- and sentence-level analysis to identify clinical content with low semantic consistency, improving the factual accuracy of the generated report.
- (3) Our framework improves factuality by abstaining from high-uncertainty reports, achieving a 10% boost in factuality scores by rejecting 20% of reports using the Radiolog model. Additionally, it flags sentences with the highest uncertainty, accurately identifying those with the lowest factual precision at 82.9%.
- (4) We evaluate our framework’s effectiveness in detecting non-existent prior exams and investigate its alignment with factuality across various pathology subgroups.

2 Preliminaries

2.1 RRG with VLLMs

In RRG using VLLMs, the input is a medical image $x \in \mathbb{R}^D$, where D is the dimension of the image, and the output is a generated report $\hat{r} \in \mathcal{V}^*$, with \mathcal{V}^* represents the space of token sequences. To

produce this report, the model processes the image through a series of transformations across three main components.

First, the **image encoder** extracts the visual tokens $f_x = \text{Enc}(x)$ from the image x . Next, an **Alignment Module** generates an embedding $z = g(f_x)$ to map these visual tokens to a text-compatible space. This alignment allows the visual data to be effectively interpreted by the language model (LM). Finally, the aligned embeddings z are passed to a **Large Language Model (LLM)**, \mathcal{M} , which generates the radiology report as a sequence of tokens $\hat{r} \in \mathcal{V}^*$.

The quality of the generated report \hat{r} is evaluated against a reference report r using a correctness function $A(\hat{r}, r)$, which measures lexical or semantic similarity, assessing how well the generated report aligns with the ground truth.

2.2 Rank Calibration

Rank Calibration (Huang et al., 2024) is designed to evaluate the alignment between the uncertainty levels of an LM’s predictions and their actual (in)correctness. An uncertainty measure is considered rank-calibrated if predictions with higher uncertainty are more likely to be incorrect. Given N predictions by the LM, each associated with an uncertainty score u_i for $i = 1, 2, \dots, N$, these scores are evenly partitioned into B intervals $\{\mathcal{I}_b\}_{b=1}^B$, such that each interval contains approximately N/B scores. Using a regression function reg to map the uncertainty score u from any interval \mathcal{I}_b to the accuracy of predictions in that interval, the Empirical Rank-Calibration Error (RCE) assesses the alignment between uncertainty and accuracy as below. Lower Empirical RCE values indicate better calibration.

$$RCE = \frac{1}{B} \sum_{b=1}^B \left| \frac{\sum_{b'=1}^B \mathbf{1}_{b' \neq b} \left[\frac{\sum_{u' \in \mathcal{I}_{b'}} reg(u') \geq \sum_{u \in \mathcal{I}_b} reg(u)}{B-1} \right]}{\sum_{b'=1}^B \mathbf{1}_{b' \neq b} \left[\frac{\sum_{u' \in \mathcal{I}_{b'}} u' \leq \sum_{u \in \mathcal{I}_b} u}{B-1} \right]} \right| \quad (1)$$

2.3 VRO (Variation Ratio for Original Prediction)

The VRO metric (Huang et al., 2023b) measures uncertainty by comparing the model’s original prediction with the predictions generated from multi-

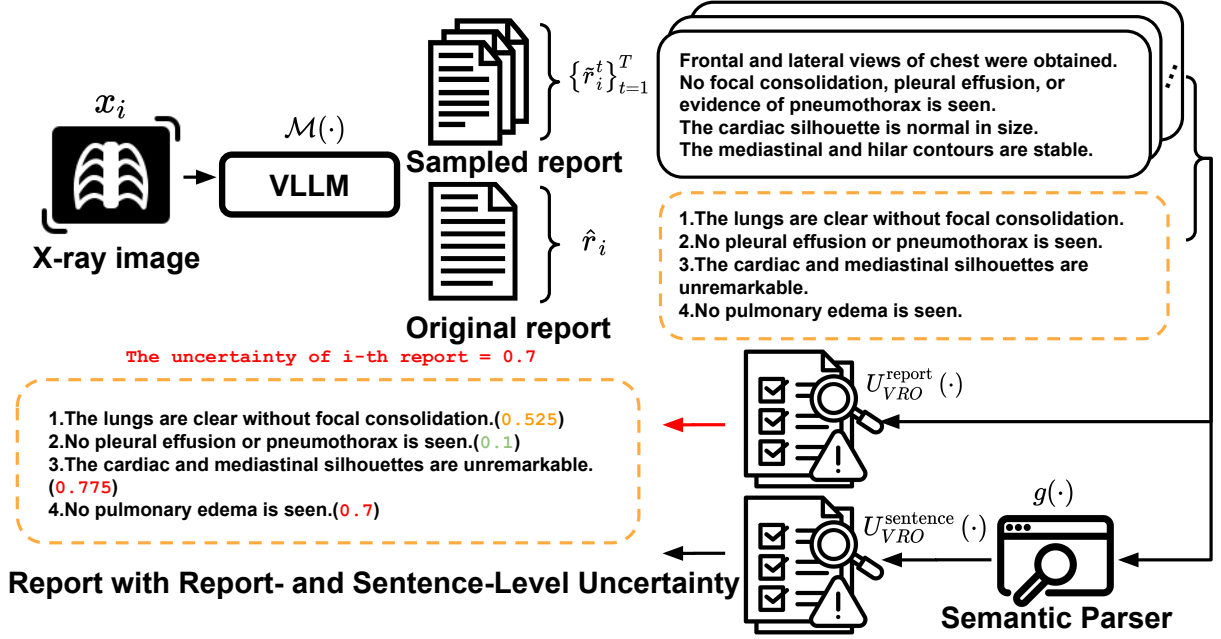


Figure 1: Pipeline of proposed Uncertainty Quantification Framework. Given an X-ray image x_i , the LLM generates an original report \hat{r}_i and sampled reports $\{\tilde{r}_i^t\}_{t=1}^T$. These reports are first processed by a semantic parser g , which extracts entity-label pairs for each sentence in \hat{r}_i . The uncertainty quantification module evaluates semantic consistency at both the report and sentence levels, providing a comprehensive, layered view of uncertainty for the generated report.

ple stochastic inferences. It is calculated as:

$$VRO = 1 - \frac{1}{T} \sum_{i=1}^T (1 - \text{dist}(p_i, p_{LM})) \quad (2)$$

where T is the number of inferences, p_i is the prediction from the i -th inference, and p_{LM} is the original prediction from the model. The function $\text{dist}(\cdot)$ measures the distance between two predictions. Lower VRO values indicate greater consistency between the original and sampled predictions, signifying lower uncertainty in the model’s output.

2.4 Radgraph

RadGraph (Jain et al., 2021) structures chest X-ray reports by extracting clinical entities and their relationships as multiple triplets. Entities include Anatomy (ANAT-DP) and three types of Observation: Definitely Present (OBS-DP), Uncertain (OBSU), and Definitely Absent (OBS-DA). Anatomy refers to body parts like “lung,” while Observations describe features or diagnoses, such as “effusion” or “increased.” Relations between entities are categorized as Suggestive Of, Located At, or Modify, indicating how observations are inferred, located, or modified. To perform this extraction, PubMedBERT (Gu et al., 2021), a pre-trained biomedical language model, was fine-tuned on the RadGraph dataset. It processes radiology report

text to automatically label entities and relations, enabling structured analysis of the clinical content.

2.5 Natural Language Inference based Uncertainty Quantification

A Natural Language Inference (NLI) model takes a pair of sentences (a premise and a hypothesis) as input and outputs logits for the labels—entailment, contradiction, or neutral—indicating the likelihood of each relationship. Kuhn et al. (2023); Lin et al. (2023) leverage these pairwise similarity scores to assess the consistency between response pairs and use them for subsequent uncertainty estimation. Zhang et al. (2024a) use an off-the-shelf DeBERTa-v3-large model (He et al., 2021) to compute NLI-based uncertainty for each sentence s_j in a response. They calculate the probability of “entailment” by normalizing the entailment logit l_e over the sum of entailment and contradiction logits:

$$P(\text{entail} \mid s_j, r') = \frac{\exp(l_e)}{\exp(l_e) + \exp(l_c)} \quad (3)$$

In this way, sentence-report similarity can be calculated to enable UQ.

3 Method

In RRG, given an LLM \mathcal{M} , we use $\tilde{r}_i^t = \mathcal{M}_t(x_i)$ to denote t -th sampled report given a Chest X-

ray image x_i in contrast to $\hat{r}_i = \mathcal{M}(x_i)$ as original report. An uncertainty measure is defined as $U^{\mathcal{M}} : \mathcal{V}^* \times 2^{\mathcal{V}^*} \rightarrow \mathbb{R}$, takes the original report and a set of sampled reports as input and outputs a real value representing the uncertainty. Our core principle is that higher uncertainty should correspond to lower quality in generated outputs. Therefore, we follow the approach as described in Section 2.2 to evaluate the alignment between $U^{\mathcal{M}}$ and the prediction correctness indicated by a clinical metric F via Equation 1 by defining the regression function as $reg(u) = \mathbb{E}[F|U^{\mathcal{M}} = u]$. However, the long-form nature of RRG poses the following challenges in designing the uncertainty measurement U^1 .

- (a) **High Similarity Across Responses:** long texts often yield high similarity across response pairs (Zhang et al., 2024a), limiting UQ methods based on response-level similarity (Kuhn et al., 2023; Lin et al., 2023). Applying similarity at the component level requires extra effort to align corresponding parts, as sampled responses may reorder or omit claims.
- (b) **Lack of Domain-Specific NLI Models:** Zhang et al. (2024a) propose using NLI models for nuanced similarity assessments; however, RRG lacks a specialized NLI model. General NLI models often struggle with the domain’s subtle distinctions, causing error propagation. While Bannur et al. (2024) leverage the in-context learning abilities of GPT-4 and Llama3-70B for entailment verification—potentially making them viable as UQ methods in RRG—these models are impractical for real-time UQ due to high computational demands. See further discussion in Appendix B.
- (c) **Limitations of Self-Evaluation-Based UQ:** Self-evaluation UQ methods (Kadavath et al., 2022; Lin et al., 2023) attempt to verbalize confidence through handcrafted prompts, enabling models to express uncertainty in natural language. However, this approach is currently unavailable for VLLM-based RRG models (Gui et al., 2024), with failure cases demonstrated in the Appendix B.

To overcome these challenges, we propose to quantify uncertainty by evaluating semantic similarity between paired reports with clinical metric F . By focusing on semantic consistency, our method more effectively captures semantic equivalence,

leading to improved uncertainty estimation. In addition, we apply VRO (Huang et al., 2023b), which calculates the similarity between the original and sampled predictions, to enhance computational efficiency. In contrast to previous methods (Zhang et al., 2024a; Kuhn et al., 2023) that require $O(n^2)$ calls to NLI models for pairwise comparisons, our approach reduces this complexity to $O(n)$ calls for consistency measurement while maintaining good performance in UQ with different granularity. In Section 3.1 we will provide details on report-level UQ, while Section 3.2 details sentence-level UQ.

3.1 Report-Level Uncertainty Quantification

Our report-level uncertainty quantification leverages the approach in Equation 2, where we use a factual metric in RRG as the distance function. In this setup, \hat{r}_i is treated as the original prediction, and \tilde{r}_i^t represents the t -th sampled prediction. The uncertainty is computed as:

$$U_{\text{VRO}}^{\text{report}}(\hat{r}_i, \{\tilde{r}_i^t\}_{t=1}^T) = \frac{1}{T} \sum_{t=1}^T (1 - F(\hat{r}_i, \tilde{r}_i^t)) \quad (4)$$

We leverage GREEN (Ostmeier et al., 2024), a state-of-the-art evaluation metric that aligns with radiologist preferences, to implement F . GREEN calculates factual alignment by comparing findings and error counts between reports. Here, the original report serves as the prediction, and the sampled reports are references, effectively capturing semantic equivalence between the original generated and sampled reports. Further details on GREEN are in the Appendix A.

3.2 Sentence-Level Uncertainty Quantification

While report-level uncertainty quantification is useful, it can obscure variations in certainty across multiple facts within a report, making sentence-level quantification more appropriate. Zhang et al. (2024a) calculates sentence-to-report entailment scores across all sampled reports increases classifier complexity and computational demands, making it inefficient for real-world deployment. To overcome these challenges, we propose a novel method leveraging the RadGraph (Jain et al., 2021) parser. Assume that each report, $\hat{r}_i = \{s_{i1}, s_{i2}, s_{i3} \dots s_{ik_i}\}$, consists of multiple sentences, where k_i indicates the number of sentences within the report. We utilize the RadGraph parser, denoted as $g : \mathcal{V}^* \rightarrow \bar{\mathcal{V}}$, which map sequence(s) to the

¹Will omit \mathcal{M} when the choice of the LM is clear.

set of node-label pairs $\bar{V} = \{(v_k, v_{k_L})\}_{k \in [1..|V|]}$ where each pair represents an entity and its associated label. An entity v_k is a continuous text span (potentially multi-word) that represents either an Anatomy or an Observation. The label v_{k_L} for each entity v_k indicates one of the four possible entity categories describe in Section 2.4. Using this structured output, we calculate an uncertainty value for each sentence s_{ij} , where s_{ij} is the j -th sentence in the generated report \hat{r}_i .

$$U_{\text{VRO}}^{\text{sentence}}(s_{ij}, \{\tilde{r}_i^t\}_{t=1}^T) = \frac{1}{T} \sum_{t=1}^T \left(1 - \frac{|g(s_{ij}) \cap g(\tilde{r}_i^t)|}{|g(s_{ij})|} \right) \\ = \frac{1}{T} \sum_{t=1}^T \left(1 - \frac{|\bar{V}_{ij} \cap \bar{V}_i^t|}{|\bar{V}_{ij}|} \right) \quad (5)$$

where \bar{V}_{ij} represents the set of node-label pairs in the original sentence s_{ij}^j , and \bar{V}_i^t is the set of node-label pairs in the sampled report \tilde{r}_i^t . The term $|\bar{V}_{ij} \cap \bar{V}_i^t|$ denotes the number of entity-label pairs from the original sentence \bar{V}_{ij} that are also present in \bar{V}_i^t . Additionally, $|\bar{V}_{ij}|$ represents the total number of node-label pairs in the original sentence s_{ij}^j , which bounds the uncertainty value between 0 and 1.

4 Experiments

In this section, we aim to answer the following research questions: **RQ1**. How well does our proposed UQ align with the factual correctness of the generated reports? **RQ2**. Can our UQ enhance the radiologist’s intervention process to improve the factual accuracy of generated reports? **RQ3**. Can our UQ detect content referring to non-existent prior information?

4.1 Setup

Datasets. Following previous works, we conduct our experiments on MIMIC-XCR (Johnson et al., 2019). We follow the original train-val-test splits. **Models.** We use RaDialog (Pellegrini et al., 2023) as the base model for our uncertainty quantification experiments. This model was selected due to its clean architecture, strong performance on the ReXRANK (Zhang et al., 2024c) online benchmark, and ease of reproducibility without data restrictions. To further validate our approach, we also apply our method to CheXpertPlus_mimiccxr (Cham-bon et al., 2024), a top-performing model on the MIMIC-CXR benchmark. For this model, we assume only API access to demonstrate the flexibility,

plug-and-play nature, and generalizability of our proposed uncertainty quantification framework to different vision-language model-based radiology report generation systems.

4.2 RRG Evaluation

We evaluate our RRG models using four metric categories from the ReXRANK benchmark (Zhang et al., 2024c), supplemented by the state-of-the-art GREEN evaluation. We use **lexical metrics** such as BLEU and embedding-based BERTScore to assess token-level and semantic similarity. To evaluate pathological and entity-based consistency, we apply **factuality metrics**, including Semb Score and RadGraph Precision, Recall, and F1. We further assess clinical accuracy with **RadCliQ**, which combines RadGraph F1 and BLEU scores, and the **GPT-based evaluator** GREEN, which evaluates clinical accuracy by matching findings and counting errors between generated and reference reports. For detailed descriptions of each metric, see Appendix C.

4.3 UQ Evaluation

In this section, we show how UQ can be evaluated in radiology report generation. In contrast to typical question-answering tasks where the correctness of a model’s prediction is binary, radiology report generation typically involves long-form generation which requires more nuanced evaluation methods for UQ.

Pearson correlation coefficient. The Pearson correlation coefficient can be used to assess how well uncertainty quantification aligns with the factual correctness of generated reports. By measuring the linear relationship between model uncertainty and report quality, Pearson’s coefficient provides insight into whether higher uncertainty corresponds to lower factual accuracy. The Pearson correlation ranges between -1 and 1 , where a negative value indicates an inverse relationship. In our setting, we use Pearson’s coefficient to evaluate this relationship, with a strong negative correlation suggesting that higher uncertainty signals lower report correctness, aligning with the intended behavior of uncertainty quantification.

Rank calibration error. RCE assesses the consistency in ranking, ensuring higher uncertainty corresponds to lower correctness, regardless of a linear relationship. We use the Empirical RCE (Huang et al., 2024), which divides uncertainty values into $B = 20$ bins. For each bin, we calculate the ex-

pected correctness level and average uncertainty. The Empirical RCE is computed by averaging the rank differences between correctness and uncertainty across these bins as Equation 1, offering a principled approach to measure the alignment between uncertainty and correctness without relying on arbitrary thresholds.

Abstention. Abstention allows uncertainty quantification to enhance factual accuracy by rejecting high-uncertainty reports, directing radiologists to focus on certain content. Traditionally, abstention is measured by metrics like AUARC (Huang et al., 2024), which evaluates improvement by abstaining from uncertain cases. However, binary metrics like AUARC are inadequate for the nuanced nature of RRG. To address this, we evaluate abstention at the report level, measuring improvements in factuality scores while balancing the trade-off with coverage. This strategy enables targeted intervention by radiologists, focusing their review on areas where factual accuracy may be compromised.

Uncertainty Precision Alignment. To evaluate sentence-level UQ in RRG, we calculate a factual precision score for each generated sentence using RadGraph (details in Appendix C). We then assess how well high uncertainty scores correspond to sentences with low factual precision within each report. Specifically, we measure the alignment rate between the sentence with the highest uncertainty and the sentence with the lowest factual precision. This alignment metric supports targeted interventions, enabling radiologists to focus on sentences that may require closer review due to potential factual inaccuracies.

4.4 Hallucination Detection

In RRG, references to prior exams are a common form of hallucination (Banerjee et al., 2024). In this section, We empirically investigate whether our UQ can effectively detect and flag these hallucinations by assigning them high uncertainty. Following Banerjee et al. (2024), we define 43 substrings commonly associated with references to prior exams. For report-level uncertainty, we analyze the changes in the percentage of reports with prior exam references and the average number of hallucinated substrings per report before and after applying different levels of abstention.

4.5 UQ Baselines

We compare our method with the previous uncertainty quantification method. Following Kuhn

et al. (2023), we use predictive entropy, length-normalised predictive entropy (Malinin and Gales, 2020) and lexical similarity (Zhang et al., 2024a; Fomicheva et al., 2020). We do not compare with methods involving NLI classifiers and self-evaluation-based UQ due to their unavailability in RRG, as discussed in Appendix B. For all experiments, we use the default temperature value 1 and sample 10 responses to calculate UQ.

5 Results

5.1 Alignment with Factuality (RQ1)

Table 1 demonstrates that our proposed VRO-GREEN exhibits stronger negative Pearson correlations with factuality metrics across both the Radialog Model and the CheXpertPlus_mimiccxr Model when compared to baseline UQ methods. In particular, VRO-GREEN achieves high negative correlations on GREEN (-0.5292 for Radialog, -0.4726 for CheXpertPlus_mimiccxr) and RadCliQ-v0 (-0.4137 for Radialog, -0.3743 for CheXpertPlus_mimiccxr). This indicates VRO-GREEN’s superior capability in aligning uncertainty with factual correctness in radiology report generation. Furthermore, we grouped samples based on the presence of specific pathology findings to examine the correlation between UQ and GREEN for each subgroup for the Radialog Model. Subgroup analysis in Table 3 reveals variation in correlation strength, particularly in the Pneumothorax subgroup, where the correlation is notably weaker at -0.08 , likely due to the underrepresentation of Pneumothorax cases (around 1% of positive cases in the training set).

Table 2 further validates VRO-GREEN’s alignment effectiveness using Empirical RCE, where it achieves the lowest RCE values on both GREEN (0.015 for Radialog, 0.02 for CheXpertPlus_mimiccxr) and Negative RadCliQ-v0 (0.02 for Radialog, 0.025 for CheXpertPlus_mimiccxr). These results confirm VRO-GREEN’s superior consistency in aligning uncertainty with factual correctness across multiple metrics.

At the sentence level, the Pearson correlation between sentence-level uncertainty (VRO-RadGraph) and factual precision is strong for both models, with -0.52 for the Radialog model and -0.55 for the CheXpertPlus_mimiccxr model, indicating effective alignment with factuality at sentence level.

Uncertainty Method	BLEU Score	BERTScore	Semb Score	RadGraph Recall	RadGraph Precision	RadGraph Combined	GREEN	-RadCliQ-v0	-RadCliQ-v1
Radialog Model									
VRO-GREEN (Ours)	-12.15	-40.92	-30.71	-19.95	-34.15	-27.88	-52.92	-41.37	-39.46
Predictive Entropy	-0.79	-28.77	-21.03	-2.78	-22.35	-12.25	-32.84	-27.72	-25.18
Normalized Entropy	-8.43	-22.32	-19.00	-10.58	-16.78	-13.60	-36.34	-23.19	-22.04
Lexical Similarity	-12.42	-38.13	-26.80	-16.67	-31.46	-24.95	-38.75	-37.21	-35.62
CheXpertPlus Model									
VRO-GREEN (Ours)	-9.39	-34.00	-29.72	-22.70	-30.20	-27.62	-47.26	-37.43	-35.51
Lexical Similarity	-15.17	-31.48	-25.25	-23.30	-27.32	-26.38	-34.91	-32.95	-32.11

Table 1: Pearson correlation values (expressed as percentages) for various metrics with different UQ methods across two models: Radialog and CheXpertPlus_mimicxr. Stronger negative values indicate better performance. Negative RadCliQ metric are used to align with other metrics in Pearson correlation calculations. For CheXpertPlus_mimicxr, we assume API-only access to the model, so only lexical similarity is compared in the table.

UQ Method	RCE(GREEN)	RCE(-RadCliQ-v0)
Radialog Model		
VRO-GREEN (Ours)	0.015	0.02
Predictive Entropy	0.045	0.09
Normalized Entropy	0.045	0.145
Lexical Similarity	0.045	0.04
CheXpertPlus Model		
VRO-GREEN (Ours)	0.02	0.025
Lexical Similarity	0.03	0.03

Table 2: Empirical RCE results for various UQ metrics measured on GREEN and Negative RadCliQ-v0 correctness, with results presented separately for Radialog and CheXpertPlus_mimicxr models.

5.2 Enhancing RRG Intervention(RQ2)

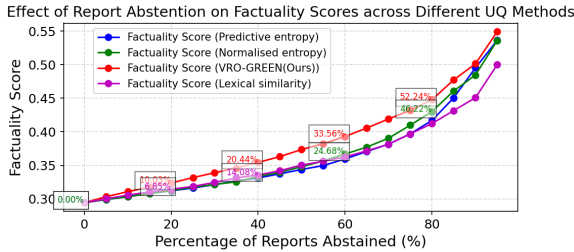


Figure 2: Effect of Report Abstention on Factuality Score across UQ for the RaDialog model. The percentages in boxes represent the improvement(only top-2 visualized) in factuality score after abstention, relative to the initial performance without abstention.

Figure 2 and Figure 5 illustrate the impact of report-level abstention on factuality scores (GREEN) for the Radialog and CheXpertPlus models, respectively. By excluding the top 20% most uncertain reports, our UQ method achieves notable factuality improvements: 10% for Radialog and 9.2% for CheXpertPlus, demonstrating consistent gains across models. These results highlight our method’s effectiveness in enhancing report quality and supporting radiologists in focusing

on more reliable reports.

At the sentence level, uncertainty-precision alignment results reveal that for the Radialog model, the highest-uncertainty sentence aligns with the lowest factual precision at a rate of 82.9%, while the lowest-uncertainty sentence aligns with the highest factual precision at only 59.1%. For the CheXpertPlus model, these rates are 81.2% and 59.6%, respectively, closely mirroring the trend observed in Radialog. This discrepancy indicates that while our sentence-level UQ method effectively flags low-precision sentences with high uncertainty, it performs poorly in cases of low-uncertainty sentences, highlighting the presence of confidently hallucinated sentences that our method struggles to capture. This limitation underscores a key challenge in our current approach and suggests an avenue for future work. More details are discussed in Section 8.

5.3 Detection of Hallucinations of Prior Exams (RQ3)

Figure 3 and Figure 6 demonstrate the effectiveness of our report-level uncertainty quantification in detecting hallucinations of prior exams for the Radialog and CheXpertPlus models, respectively. Rejecting high-uncertainty reports leads to a clear decrease in the percentage of reports with prior references and the average number of prior-related substrings, significantly improving hallucination detection. In contrast, the random baseline, averaged over 5 seeds, shows no reduction in these metrics.

5.4 Qualitative Analysis.

In this section, we analyze the qualitative aspects of our UQ framework for radiology report generation. Specifically, we explore (1) the effect of increasing

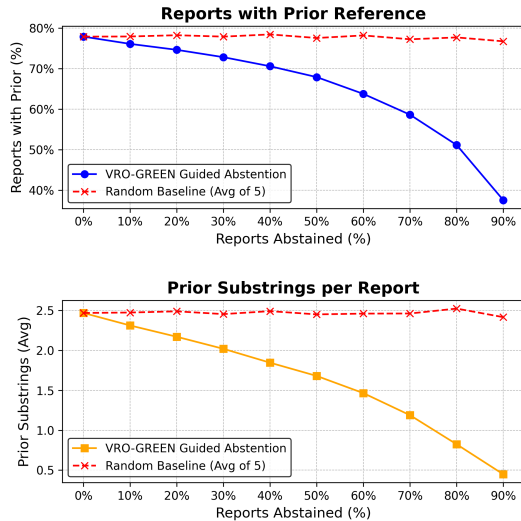


Figure 3: Effect of VRO-GREEN Guided Abstention on Prior References and Substrings for the Radialog model. Solid lines represent VRO-GREEN Guided Abstention, with dashed red lines as the baseline performing random abstention.

the number of sampled reports on UQ performance and (2) a case study showcasing the practical utility of our framework in identifying low factual correctness and guiding radiologist interventions.

Number of sampled reports. Research on short question-answer tasks and long-form generation has shown that increasing the number of sampled responses can enhance the performance of uncertainty quantification. We extend this investigation to radiology report generation, exploring whether the same holds true in this domain. As illustrated in Figure 7, our findings align with previous research, showing that the performance of UQ improves with more samples and converges when using seven samples.

Case Study. Figure 4 illustrates our UQ framework’s ability to identify low factual correctness and assist radiologists in targeted review. In Figure 4(a), report-level UQ assigns the highest uncertainty score to a nonsensical report with zero factual correctness. Figure 4(b) shows sentence-level UQ, ranking sentences by uncertainty to guide radiologist intervention: the high-uncertainty (red) sentence references a non-existent prior exam, the moderate-uncertainty (orange) sentence is partially correct, and the low-uncertainty (green) sentence is fully accurate.

6 Related Work

Multimodal Foundation Models, such as VLLMs, augment large language models (LLMs) with visual inputs (Antropic, 2024; OpenAI, 2023). These models are typically pre-trained on diverse datasets (Erhan et al., 2010; Chen et al., 2020; Li et al., 2022; Lin et al., 2024; Alayrac et al., 2022) before applied to specialized tasks, reducing the requirements for domain-specific data. VLLMs have been evaluated in medical applications such as medical image interpretation and radiology report generation (Litjens et al., 2017; Esteva et al., 2021; Moor et al., 2023; Srivastav et al., 2024), and have demonstrated performance comparable to previous supervised methods (Rajpurkar, 2017; Qin et al., 2018), and in some cases, even rival medical experts (Tiu et al., 2022). However, there are challenges that hinder establishing trust in multimodal foundation models in clinical practice (Truhn et al., 2024; Freyer et al., 2024; Ong et al., 2024). These challenges include ensuring the quality and transparency of the training data (Koçak, 2022; Celi et al., 2022; Chen et al., 2023), effective collaboration between machine learning experts and medical professionals (Cai et al., 2019), and more effective and meaningful evaluation measurements (Wornow et al., 2023). Our paper focuses on the factuality of VLLMs’ predictions in radiology applications where ensuring accuracy and trustworthiness are critical for clinical decision-making (Bates et al., 2021).

Hallucination in Foundation Models (Rawte et al., 2023; Ji et al., 2023) can lead to non-factual predictions, and this issue persists regardless of the model’s size (Lee et al., 2023; Jeblick et al., 2024; Xu et al., 2024; Zhang et al., 2024b). Studies (Zheng et al., 2023; Lu et al., 2023) have shown that a lack of domain-specific knowledge in assigned tasks can cause foundation models to produce hallucinated outputs, a behavior that is often difficult for the models to correct on their own (Huang et al., 2023a). A large body of research has focused on addressing this issue by filling the knowledge gaps with additional oracle labels (Kim et al., 2024; Shinn et al., 2024; Gou et al., 2024; Banerjee et al., 2024; Bannur et al., 2024). However, as Feng et al. (2024) points out, the knowledge gaps will always exist because knowledge is continually evolving. Moreover, in radiology, filling these gaps requires expertise-intensive labeling of data such as medical imaging data (Koyyada and Singh, 2023; Kim et al., 2022) and Electronic Health Records (EHR)



X-ray image	Generated Report	Ground-Truth Report
	<p>PA and lateral views of the chest at 10:30 are submitted.</p> <p>Factual correctness = 0 Report-level uncertainty = 1</p>	<p>In comparison with study of ____, there is little change and no evidence of acute cardiopulmonary disease. Cardiac silhouette is mildly prominent, but there is no vascular congestion, pleural effusion, or acute focal pneumonia.</p>
(a) Report-level UQ case study 1		
	<p>The hallucination of prior exams</p> <ol style="list-style-type: none"> 1. In comparison with the study of ____, there is little overall change. 2. Again there is extensive opacification at the right base consistent with pleural effusion and compressive atelectasis. 3. The left lung is relatively clear. <p>Partially correct</p> <p>Accurate sentence</p>	<p>There is a very tiny right apical pneumothorax following removal of the right-sided chest tube. There is persistent elevation of the right hemidiaphragm with atelectasis at the right lung base and a right-sided pleural effusion. A rounded opacity is seen in the right suprahilar region and is stable. The left lung is relatively clear aside from atelectasis at the left lung base and a small left-sided pleural effusion.</p>
(b) Sentence-level UQ case study 2		

Figure 4: Two separate analyses of report- and sentence-level UQ in radiology report generation using MIMIC-CXR data. (a) The report-level UQ study assigns an uncertainty score to the entire report. (b) The sentence-level UQ study ranks individual sentences by uncertainty, with red (1.0) indicating high uncertainty, orange (0.75) indicating moderate uncertainty, and green (0.47) indicating low uncertainty. This color-coded ranking helps inform radiologists on which sentences may require closer attention.

data (Mc Cord et al., 2019).

Uncertainty Quantification (UQ) has been extensively studied in conventional ML (Gupta et al., 2006; Shafer and Vovk, 2008; Vaicenavicius et al., 2019; Tibshirani et al., 2019; Abdar et al., 2021), and is receiving increasing attention due to its potential to mitigate hallucinations in foundation models (Xiao and Wang, 2021; Fadeeva et al., 2024). Straightforward methods include querying models about their confidence (Xiong et al., 2023; Joshi et al., 2017), and using Perplexity score (Jelinek et al., 1977). Recent research on semantic uncertainty (Kuhn et al., 2023; Zhang et al., 2024a) draw insights from the coherence of model predictions by using an additional NLI (MacCartney, 2009) classifier. Calibration methods and conformal prediction techniques (Liu and Wu, 2024; Quach et al., 2024; Gui et al., 2024) can offer statistical guarantees on the factuality of the outputs, provided that there is a held-out dataset for extracting necessary information. UQ in radiology report generation poses unique challenges to existing methods by involving image inputs and long-form radiologist reports as outputs (Koçak, 2022; Jeblick et al., 2024; Smit et al., 2020). Our method is related to RadGraph (Jain et al., 2021), which structures radiology reports by extracting pre-defined clinical entities and their relations. Prior works that used RadGraph for UQ involve an additional reinforcement learning step (Delbrouck et al., 2022). However, our approach does not involve such step.

7 Conclusion

In this paper, we tackle the challenge of hallucinations in RRG through a novel UQ approach. Our plug-and-play framework introduces both report-level and sentence-level UQ to detect low-factuality reports and identify non-existent prior hallucinations, supporting more effective radiologist intervention. Applied to the MIMIC-CXR dataset, our method achieved a 10% improvement in factuality by rejecting 20% of high-uncertainty reports using the Radialog model. Additionally, sentence-level UQ flagged sentences with the lowest factual precision at 82.9% accuracy, enabling targeted intervention. Future work will focus on exploring supervised uncertainty measures to improve factuality, particularly addressing cases where the UQ framework assigns low uncertainty to hallucinated predictions generated by VLLMs. Additionally, integrating uncertainty directly into the generation process could guide models toward more factual outputs by conditioning generation on uncertainty thresholds, thus enhancing both the reliability of UQ and overall model trustworthiness.

8 Acknowledgments

This work was supported in part by the U.S. National Science Foundation under grant CCF-2340776. Additional support was provided by the NIH (Award Number 1R01HL141813-01) and the Pennsylvania Department of Health. We also gratefully acknowledge the computational resources made available by Pittsburgh Super Computing (grant number TGASC170024).

Limitations

In this section, we outline the limitations of our work and potential areas for improvement.

First, while we demonstrate the effectiveness of our method across different model architectures using the MIMIC-CXR dataset, our evaluation is limited to this dataset. Expanding our experiments to other datasets, such as IU X-Ray (Demner-Fushman et al., 2016) or the recently published CheXpert Plus (Chambon et al., 2024), could further validate the generalizability of our approach.

Second, due to challenges outlined in Section 3, we were only able to compare our method against three relatively simple baselines. As UQ techniques continue to evolve within this domain, the development of domain-specific models, such as tailored NLI models for RRG, could enable a more comprehensive comparison in future work.

Third, while our sentence-level uncertainty quantification effectively aligns high-uncertainty sentences with low factual precision, it struggles to align low-uncertainty sentences with high factual precision, revealing a gap in detecting confidently hallucinated sentences. This limitation suggests the need for enhanced UQ techniques and the potential benefit of incorporating a fact-checking module to improve reliability and distinguish factual inaccuracies.

Finally, our current sentence-level UQ is designed with intervention in mind, focusing solely on the factual precision of generated reports. However, this approach overlooks factual completeness, meaning it does not account for important factual information that may be omitted from the generated report. Future work could address this by designing UQ methods that consider both factual precision and completeness, providing a more balanced evaluation of report quality.

These limitations highlight opportunities for further refinement and experimentation in UQ methodologies for radiology report generation

References

- Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. 2021. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information fusion*, 76:243–297.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736.
- Antropic. 2024. [The claude 3 model family: Opus, sonnet, haiku](#).
- Oishi Banerjee, Hong-Yu Zhou, Subathra Adithan, Stephen Kwak, Kay Wu, and Pranav Rajpurkar. 2024. Direct preference optimization for suppressing hallucinated prior exams in radiology report generation. *arXiv preprint arXiv:2406.06496*.
- Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Anton Schwaighofer, Sam Bond-Taylor, Maximilian Ilse, Fernando Pérez-García, Valentina Salvatelli, Harshita Sharma, Felix Meissen, et al. 2024. Maira-2: Grounded radiology report generation. *arXiv preprint arXiv:2406.04449*.
- Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Perez-Garcia, Maximilian Ilse, Daniel C Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, et al. 2023. Learning to exploit temporal structure for biomedical vision-language processing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15016–15027.
- David W Bates, David Levine, Ania Syrowatka, Masha Kuznetsova, Kelly Jean Thomas Craig, Angela Rui, Gretchen Purcell Jackson, and Kyu Rhee. 2021. The potential of artificial intelligence to improve patient safety: a scoping review. *NPJ digital medicine*, 4(1):54.
- Carrie J Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. "hello ai": uncovering the onboarding needs of medical practitioners for human-ai collaborative decision-making. *Proceedings of the ACM on Human-computer Interaction*, 3(CSCW):1–24.
- Leo Anthony Celi, Jacqueline Cellini, Marie-Laure Charpignon, Edward Christopher Dee, Franck Deroncourt, Rene Eber, William Greig Mitchell, Lama Moukheiber, Julian Schirmer, Julia Situ, et al. 2022. Sources of bias in artificial intelligence that perpetuate healthcare disparities—a global review. *PLOS Digital Health*, 1(3):e0000022.
- Pierre Chambon, Jean-Benoit Delbrouck, Thomas Sounack, Shih-Cheng Huang, Zhihong Chen, Maya Varma, Steven QH Truong, Chu The Chuong, and Curtis P Langlotz. 2024. Chexpert plus: Augmenting a large chest x-ray dataset with text radiology reports, patient demographics and additional image formats. *arXiv preprint arXiv:2405.19538*.
- Richard J Chen, Judy J Wang, Drew FK Williamson, Tiffany Y Chen, Jana Lipkova, Ming Y Lu, Sharifa Sahai, and Faisal Mahmood. 2023. Algorithmic fairness in artificial intelligence for medicine and healthcare. *Nature biomedical engineering*, 7(6):719–742.

- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- Jean-Benoit Delbrouck, Pierre Chambon, Christian Bluethgen, Emily Tsai, Omar Almusa, and Curtis P Langlotz. 2022. Improving the factual correctness of radiology report generation with semantic rewards. *arXiv preprint arXiv:2210.12186*.
- Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. 2016. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310.
- Dumitru Erhan, Yoshua Bengio, Aaron Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. 2010. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(19):625–660.
- Andre Esteva, Katherine Chou, Serena Yeung, Nikhil Naik, Ali Madani, Ali Mottaghi, Yun Liu, Eric Topol, Jeff Dean, and Richard Socher. 2021. Deep learning-enabled medical computer vision. *NPJ digital medicine*, 4(1):5.
- Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, et al. 2024. Fact-checking the output of large language models via token-level uncertainty quantification. *arXiv preprint arXiv:2403.04696*.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. 2024. [Don’t hallucinate, abstain: Identifying LLM knowledge gaps via multi-LLM collaboration](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14664–14690, Bangkok, Thailand. Association for Computational Linguistics.
- Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Oscar Freyer, Isabella Wiest, Jakob Kather, and Stephen Gilbert. 2024. [A future role for health applications of large language models depends on regulators enforcing safety standards](#). *The Lancet. Digital health*, 6:e662–e672.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2024. [CRITIC: Large language models can self-correct with tool-interactive critiquing](#). In *The Twelfth International Conference on Learning Representations*.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23.
- Yu Gui, Ying Jin, and Zhimei Ren. 2024. Conformal alignment: Knowing when to trust foundation models with guarantees. *arXiv preprint arXiv:2405.10301*.
- Hoshin V Gupta, Keith J Beven, and Thorsten Wagener. 2006. Model calibration and uncertainty estimation. *Encyclopedia of hydrological sciences*.
- Iryna Hartsock and Ghulam Rasool. 2024. Vision-language models for medical report generation and visual question answering: A review. *arXiv preprint arXiv:2403.02469*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023a. Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*.
- Xinmeng Huang, Shuo Li, Mengxin Yu, Matteo Sesia, Hamed Hassani, Insup Lee, Osbert Bastani, and Edgar Dobriban. 2024. Uncertainty in language models: Assessment through rank-calibration. *arXiv preprint arXiv:2404.03163*.
- Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2023b. Look before you leap: An exploratory study of uncertainty measurement for large language models. *arXiv preprint arXiv:2307.10236*.
- Stephanie L Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C Castro, Mercy Ranjit, Anton Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, et al. 2023. Maira-1: A specialised large multimodal model for radiology report generation. *arXiv preprint arXiv:2311.13668*.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597.
- Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven QH Truong, Du Nguyen Duong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew P Lungren, Andrew Y Ng, et al. 2021. Radgraph: Extracting clinical entities and relations from radiology reports. *arXiv preprint arXiv:2106.14463*.

- Katharina Jeblick, Balthasar Schachtner, Jakob Dexl, Andreas Mittermeier, Anna Theresa Stüber, Johanna Topalis, Tobias Weber, Philipp Wesp, Bastian Oliver Sabel, Jens Rieke, et al. 2024. Chatgpt makes medicine easy to swallow: an exploratory case study on simplified radiology reports. *European radiology*, 34(5):2817–2825.
- Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. 1977. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1):S63–S63.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. 2019. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Doyun Kim, Joowon Chung, Jongmun Choi, Marc D. Succi, John Conklin, Maria Gabriela Figueiro Longo, Jeanne B. Ackman, Brent P. Little, Milena Petranovic, Mannudeep K. Kalra, Michael H. Lev, and Synho Do. 2022. [Accurate auto-labeling of chest x-ray images based on quantitative similarity to an explainable ai model](#). *Nature Communications*, 13(1):1867.
- Geunwoo Kim, Pierre Baldi, and Stephen McAleer. 2024. Language models can solve computer tasks. *Advances in Neural Information Processing Systems*, 36.
- Burak Koçak. 2022. Key concepts, common pitfalls, and best practices in artificial intelligence and machine learning: focus on radiomics. *Diagnostic and Interventional Radiology*, 28(5):450.
- Shiva prasad Koyyada and Thipendra P. Singh. 2023. [An explainable artificial intelligence model for identifying local indicators and detecting lung disease from chest x-ray images](#). *Healthcare Analytics*, 4:100206.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*.
- Peter Lee, Sebastien Bubeck, and Joseph Petro. 2023. [Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine](#). *New England Journal of Medicine*, 388(13):1233–1239.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR.
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024. Vila: On pre-training for visual language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26689–26699.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*.
- Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen A.W.M. van der Laak, Bram van Ginneken, and Clara I. Sánchez. 2017. [A survey on deep learning in medical image analysis](#). *Medical Image Analysis*, 42:60–88.
- Terrance Liu and Zhiwei Steven Wu. 2024. Multi-group uncertainty quantification for long-form text generation. *arXiv preprint arXiv:2407.21057*.
- Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. 2023. Are emergent abilities in large language models just in-context learning? *arXiv preprint arXiv:2309.01809*.
- Bill MacCartney. 2009. *Natural language inference*. Stanford University.
- Andrey Malinin and Mark Gales. 2020. Uncertainty estimation in autoregressive structured prediction. *arXiv preprint arXiv:2002.07650*.
- Kimberly A. Mc Cord, Hannah Ewald, Aviv Ladanie, Matthias Briel, Benjamin Speich, Heiner C. Bucher, and Lars G. Hemkens. 2019. [Current use and costs of electronic health records for clinical trial research: a descriptive study](#). *Canadian Medical Association Open Access Journal*, 7(1):E23–E32.
- Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. 2023. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265.
- Jasmine Chiat Ling Ong, Shelley Yin-Hsi Chang, Wasswa William, Atul J Butte, Nigam H Shah, Lita Sui Tjien Chew, Nan Liu, Finale Doshi-Velez, Wei Lu, Julian Savulescu, et al. 2024. Ethical and regulatory challenges of large language models in medicine. *The Lancet Digital Health*, 6(6):e428–e432.

- OpenAI. 2023. [Gpt-4v\(ision\) system card](#).
- Sophie Ostmeier, Justin Xu, Zhihong Chen, Maya Varma, Louis Blankemeier, Christian Bluethgen, Arne Edward Michalson, Michael Moseley, Curtis Langlotz, Akshay S Chaudhari, et al. 2024. Green: Generative radiology report evaluation and error notation. *arXiv preprint arXiv:2405.03595*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Chantal Pellegrini, Ege Özsoy, Benjamin Busam, Nassir Navab, and Matthias Keicher. 2023. Radialog: A large vision-language model for radiology report generation and conversational assistance. *arXiv preprint arXiv:2311.18681*.
- Chunli Qin, Demin Yao, Yonghong Shi, and Zhijian Song. 2018. Computer-aided detection in chest radiography based on artificial intelligence: a survey. *Biomedical engineering online*, 17:1–23.
- Victor Quach, Adam Fisch, Tal Schuster, Adam Yala, Jae Ho Sohn, Tommi S. Jaakkola, and Regina Barzilay. 2024. [Conformal language modeling](#). In *The Twelfth International Conference on Learning Representations*.
- P Rajpurkar. 2017. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *ArXiv abs/1711*, 5225.
- Vignav Ramesh, Nathan A Chi, and Pranav Rajpurkar. 2022. Improving radiology report generation systems by removing hallucinated references to non-existent priors. In *Machine Learning for Health*, pages 456–473. PMLR.
- Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*.
- Glenn Shafer and Vladimir Vovk. 2008. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3).
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2024. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36.
- Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew P Lungren. 2020. Chexbert: combining automatic labelers and expert annotations for accurate radiology report labeling using bert. *arXiv preprint arXiv:2004.09167*.
- Shaury Srivastav, Mercy Ranjit, Fernando Pérez-García, Kenza Bouzid, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Harshita Sharma, Maximilian Ilse, Valentina Salvatelli, Sam Bond-Taylor, Fabian Falck, Anja Thieme, Hannah Richardson, Matthew P. Lungren, Stephanie L. Hyland, and Javier Alvarez-Valle. 2024. [MAIRA at RRG24: A specialised large multimodal model for radiology report generation](#). In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, pages 597–602, Bangkok, Thailand. Association for Computational Linguistics.
- Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. 2023. Interactive and explainable region-guided radiology report generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7433–7442.
- Omkar Thawkar, Abdelrahman Shaker, Sahal Shaji Mullaipilly, Hisham Cholakkal, Rao Muhammad Anwer, Salman Khan, Jorma Laaksonen, and Fahad Shahbaz Khan. 2023. Xraygpt: Chest radiographs summarization using medical vision-language models. *arXiv preprint arXiv:2306.07971*.
- Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candès, and Aaditya Ramdas. 2019. Conformal prediction under covariate shift. *Advances in neural information processing systems*, 32.
- Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, Andrew Y Ng, and Pranav Rajpurkar. 2022. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nature Biomedical Engineering*, 6(12):1399–1406.
- Daniel Truhn, Jan-Niklas Eckardt, Dyke Ferber, and Jakob Nikolas Kather. 2024. Large language models and multimodal foundation models for precision oncology. *NPJ Precision Oncology*, 8(1):72.
- Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Schön. 2019. Evaluating model calibration in classification. In *The 22nd international conference on artificial intelligence and statistics*, pages 3459–3467. PMLR.
- Michael Wornow, Yizhe Xu, Rahul Thapa, Birju S. Patel, Ethan H. Steinberg, S. Fleming, Michael A. Pfeffer, Jason A. Fries, and Nigam H. Shah. 2023. [The shaky foundations of large language models and foundation models for electronic health records](#). *NPJ Digital Medicine*, 6.
- Yijun Xiao and William Yang Wang. 2021. On hallucination and predictive uncertainty in conditional language generation. *arXiv preprint arXiv:2103.15025*.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*.
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817*.

Caiqi Zhang, Fangyu Liu, Marco Basaldella, and Nigel Collier. 2024a. Luq: Long-text uncertainty quantification for llms. *arXiv preprint arXiv:2403.20279*.

Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. 2024b. [How language model hallucinations can snowball](#). In *Forty-first International Conference on Machine Learning*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Xiaoman Zhang, Hong-Yu Zhou, Xiaoli Yang, Oishi Banerjee, Julián N Acosta, Josh Miller, Ouwen Huang, and Pranav Rajpurkar. 2024c. Rexrank: A public leaderboard for ai-powered radiology report generation. *arXiv preprint arXiv:2411.15122*.

Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang. 2023. Why does chatgpt fall short in providing truthful answers? *arXiv preprint arXiv:2304.10513*.

A GREEN

Given a generated report as the hypothesis and a ground truth report as the reference, the GREEN evaluation framework assesses clinical accuracy by analyzing both error counts across various clinically significant categories and counts of matched findings. Specifically, GREEN categorizes errors as follows:

- (a) False report of a finding in the candidate.
- (b) Missing a finding present in the reference.
- (c) Misidentification of a finding’s anatomic location/position.
- (d) Misassessment of the severity of a finding.
- (e) Mentioning a comparison that isn’t in the reference.
- (f) Omitting a comparison detailing a change from a prior study. The GREEN score is then calculated as:

$$\text{GREEN} = \frac{\# \text{ matched findings}}{\# \text{ matched findings} + \sum_{i=(a)}^{(f)} \# \text{ error}_i}$$

B Challenges in Applying Other UQ Methods to RRG

NLI-based UQ

The lack of domain-specific NLI models in RRG makes this approach infeasible. Although [Bannur et al. \(2024\)](#) leverage in-context learning with large models like GPT-4 and Llama3 for entailment verification in RRG, they are mainly designed to evaluate generated reports against ground-truth reports.

Their study reports that RadFact’s entailment verification with Llama3-70B requires a single compute node with four A100 GPUs, taking approximately 17 seconds per comparison, while GPT-4, hosted on Microsoft Azure, takes around 27 seconds. Considering the computational requirements discussed in Section 3, a single report-level UQ with GPT-4 would require around 675 seconds for five sampled reports, and Llama3’s GPU needs make it too costly for UQ applications. All of the above highlights the challenges of applying NLI-based methods for UQ in radiology report generation RRG. Therefore, we call for the development of RRG-tailored NLI models to better support UQ in this domain.

Pathology Finding	Pearson Correlation (%)
No Finding	-68.44
Enlarged Cardiomeastinum	-42.37
Cardiomegaly	-35.07
Lung Opacity	-38.84
Lung Lesion	-41.47
Edema	-34.01
Consolidation	-43.60
Pneumonia	-38.26
Atelectasis	-42.91
Pneumothorax	-8.09
Pleural Effusion	-27.82
Pleural Other	-40.04
Fracture	-46.36
Support Devices	-42.27

Table 3: Pearson correlation (as percentages) between UQ and GREEN (The overall correlation is -0.52) across various subgroups of pathology findings.

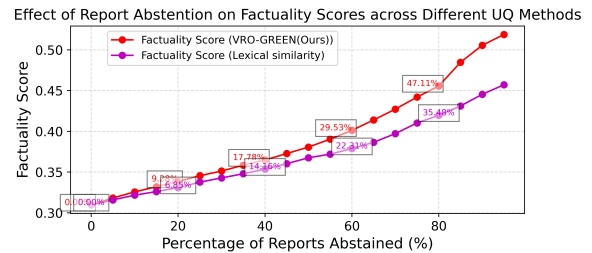


Figure 5: Effect of Report Abstention on Factuality Score across UQ for the CheXpertPlus_mimiccxr model. The percentages in boxes represent the improvement in factuality score after abstention, relative to the initial performance without abstention. We assume API-only access to the model, so only lexical similarity is compared in the figure.

Self-Evaluation-Based UQ

We demonstrate failure as shown in Figure 8 cases when applying Self-Evaluation-Based UQ to the RaDialog model in RRG. These limitations likely stem from the model’s smaller size, making it less

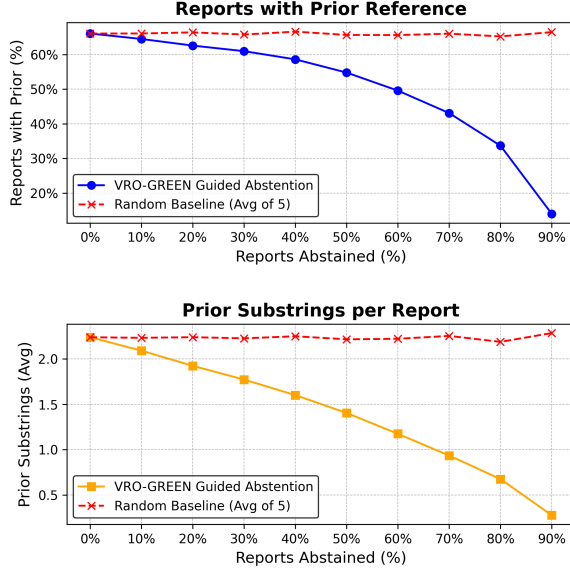


Figure 6: Effect of VRO-GREEN Guided Abstention on Prior References and Substrings for the CheXpertPlus_mimiccxr model. Solid lines represent VRO-GREEN Guided Abstention, with dashed red lines as a baseline performing random abstention.

capable of self-probing compared to larger models like GPT-4.

C Evaluation

RRG Evaluation

These metrics are organized into four categories: **Lexical Metrics.** We apply traditional Natural Language Processing (NLP) metrics such as BLEU (Papineni et al., 2002) to measure token-level similarity between the generated and ground-truth reports. In addition, we leverage the embedding-based similarity metric BertScore (Zhang et al., 2019) to capture more nuanced relationships between the texts.

Factuality Metrics. To assess factual consistency between the generated and ground-truth reports, we use two key approaches. First, Semb Score is calculated by passing both reports through the CheXbert model (Smit et al., 2020), which extracts present/absent/uncertain labels for 14 CheXpert pathological observations (Irvin et al., 2019). Cosine similarity between the resulting embeddings is then computed. Additionally, we evaluate RadGraph Precision, Recall, and F1 using the RadGraph model (Jain et al., 2021), which parses reports into graphs of clinical entities (anatomical references and observations) and their relations, followed by calculating the overlap

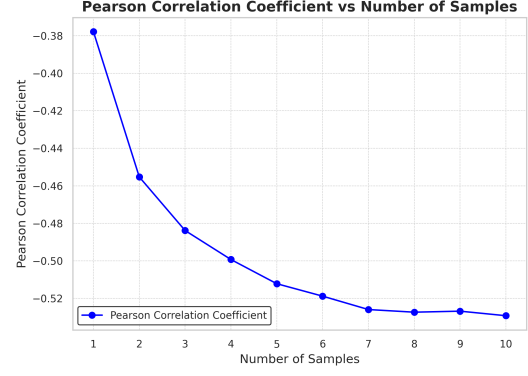


Figure 7: The effect of different number of samples on the Pearson correlation values with VRO-GREEN.

between the entities and relations in the generated and reference reports.

RadCliQ Metrics. RadCliQ integrates RadGraph F1 and BLEU scores to estimate the clinical error rate, providing a holistic quality assessment of the generated report. This metric closely aligns with radiologists’ evaluations of report quality. For evaluation, we report both version 0 and version 1 of RadCliQ.

GPT-based Evaluator. The GREEN metric is an open-source evaluation framework for radiology report generation. It calculates matching findings and error counts between the generated report and the ground-truth report, providing clinical accuracy scores, while using a smaller language model for efficiency.

Sentence Factual Precision

We use RadGraph to compute the sentence-level factual precision score for UQ evaluation. Following the notation from Section 3, we calculate the precision for each sentence s_{ij} , where s_{ij} represents the j -th sentence in the generated report \hat{r}_i .

$$\begin{aligned} P_{\text{sentence}}(s_{ij}, r_i) &= \frac{|g(s_{ij}) \cap g(r_i)|}{|g(s_{ij})|} \\ &= \frac{|\bar{V}_{ij} \cap \bar{V}_i|}{|\bar{V}_{ij}|} \end{aligned} \quad (6)$$

For edge cases in our experimented LLM model, such as the sentence "The patient is status post sternotomy," RadGraph fails to parse and returns an empty pair set. We flag these cases with a precision score of negative one², indicating the lowest

²Correspondingly, we assign an uncertainty score of $u=1$ for these sentences during uncertainty calculation.

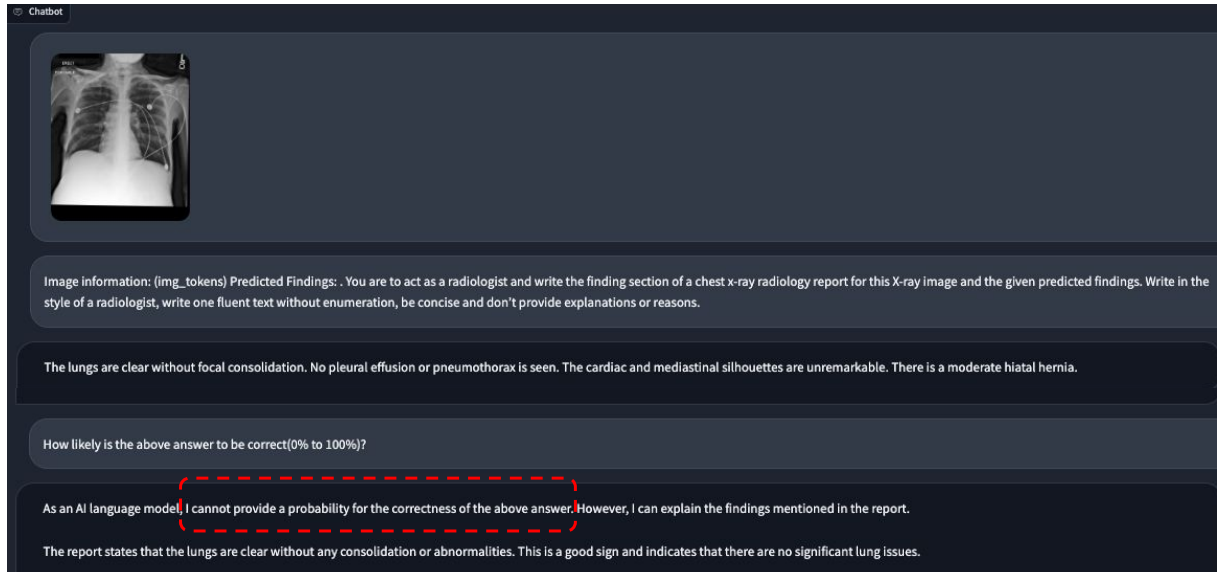


Figure 8: Example of Failure Cases in Applying Self-Evaluation-Based UQ to the RaDialog Model

Sentence Removal	BLEU Score↑	BERTScore↑	Semb Score↑	RadGraph Recall↑	RadGraph Precision↑	RadGraph Combined↑	GREEN↑	RadCliQ-v0↓	RadCliQ-v1↓
Original Report	0.1836	0.3995	0.4075	0.1801	0.2150	0.1874	0.2942	3.3838	1.1520
3% Removed	0.1814(0.1800)	0.3994(0.3915)	0.4082(0.4014)	0.1787(0.1751)	0.2166(0.2154)	0.1875(0.1842)	0.2936	3.3810(3.4217)	1.1513(1.1752)
5% Removed	0.1792(0.1763)	0.3990(0.3893)	0.4092(0.3982)	0.1778(0.1720)	0.2185(0.2154)	0.1878 (0.1822)	0.2929	3.3782(3.4348)	1.1506 (1.1845)
7% Removed	0.1758(0.1722)	0.3983(0.3869)	0.4091(0.3943)	0.1763(0.1681)	0.2202(0.2156)	0.1876(0.1795)	0.2924	3.3786(3.4508)	1.1522(1.1957)
9% Removed	0.1713(0.1693)	0.3971(0.3857)	0.4106(0.3942)	0.1753(0.1655)	0.2237(0.2169)	0.1877(0.1782)	0.2904	3.3758(3.4545)	1.1524(1.1994)
11% Removed	0.1678(0.1653)	0.3963(0.3827)	0.4116(0.3863)	0.1736(0.1611)	0.2263(0.2153)	0.1875(0.1749)	0.2889	3.3742(3.4805)	1.1530(1.2162)
13% Removed	0.1650(0.1615)	0.3959(0.3805)	0.4121 (0.3843)	0.1717(0.1584)	0.2287 (0.2165)	0.1871(0.1736)	0.2882	3.3729 (3.4905)	1.1536(1.2236)

Table 4: Comparison of various metrics across different levels of pruning guided by sentence uncertainty. Values in **bold** indicate metrics that improve the most compared to the original reports. An up arrow (↑) signifies that a higher value is better for the metric, while a down arrow (↓) indicates that a lower value is preferable. Values in parentheses show results from a baseline of random pruning at the same level.

possible precision. This approach is reasonable, as such sentences typically refer to non-existent prior exams.

Sentence Abstention

Given the positive results of report-level UQ in improving factual accuracy through report abstention, we extended this approach to sentence-level abstention to assess whether removing high-uncertainty sentences across the dataset could improve performance across various metrics. Table 4 presents the results.

As the percentage of sentence abstention increases, we observe a drop in lexical scores such as BLEU and RadGraph recall, while GREEN remains consistent. Notably, RadGraph precision and RadCliQ metrics demonstrate improvement, indicating that selectively removing high-uncertainty sentences leads to higher factual precision in these aspects. Compared to the baseline of random sentence removal at the same levels, our uncertainty-guided abstention consistently yields superior re-

sults, indirectly demonstrating that our sentence-level UQ effectively identifies sentences with low factual accuracy.

However, unlike report-level abstention, sentence-level abstention may be less practical due to the complex nature of radiology reports, where sentences often contain multiple clinical claims. Removing entire sentences risks omitting relevant information, making this approach too coarse for practical application. In future work, we aim to integrate sentence-level UQ into the generation process itself, enabling more granular control to enhance factual accuracy without the need for sentence removal.