

ReXTrust: A Model for Fine-Grained Hallucination Detection in AI-Generated Radiology Reports

Romain Hardy¹, Sung Eun Kim¹, Du Hyun Ro², Pranav Rajpurkar¹

¹Department of Biomedical Informatics, Harvard Medical School

²Department of Orthopedic Surgery, Seoul National University Hospital

romain_hardy@fas.harvard.edu, sungeun_kim2@hms.harvard.edu,

franciu7@snu.ac.kr, pranav_rajpurkar@hms.harvard.edu

Abstract

The increasing adoption of AI-generated radiology reports necessitates robust methods for detecting hallucinations—false or unfounded statements that could impact patient care. We present ReXTrust, a novel framework for fine-grained hallucination detection in AI-generated radiology reports. Our approach leverages sequences of hidden states from large vision-language models to produce finding-level hallucination risk scores. We evaluate ReXTrust on a subset of the MIMIC-CXR dataset and demonstrate superior performance compared to existing approaches, achieving an AUROC of 0.8751 across all findings and 0.8963 on clinically significant findings. Our results show that white-box approaches leveraging model hidden states can provide reliable hallucination detection for medical AI systems, potentially improving the safety and reliability of automated radiology reporting.

1 Introduction

The automated generation of radiology reports using large vision-language models (LVLMs) has recently shown promising results, offering the potential to improve workflow efficiency and standardization in clinical settings (Bannur et al. 2024; Tanida et al. 2023; Hamamci, Er, and Menze 2024; Liu et al. 2024a; Gu et al. 2024; Chen et al. 2024d; Zhou et al. 2024). However, these models sometimes generate hallucinations—statements that are false, unfounded, or inconsistent with the input images. In medical settings, such hallucinations pose significant risks, as false pathological findings could lead to unnecessary interventions or missed diagnoses, directly impacting patient care.

In this paper, we introduce ReXTrust, a white-box model designed to detect hallucinations in LVLM-generated radiology reports. ReXTrust uses a self-attention module trained on LVLM hidden states, enabling fine-grained insights into specific radiological findings and reliable overall hallucination risk scores. By analyzing internal model representations, ReXTrust can identify potential hallucinations during the generation process itself, rather than relying solely on post-hoc analysis. We evaluate ReXTrust on a subset of the MIMIC-CXR dataset, demonstrating its effectiveness across different medical categories and severity levels, with particular attention to clinically significant findings that would impact patient care. Our contributions are threefold:

1. We develop a white-box architecture for radiology finding

hallucination detection that enables both token-level attention analysis and finding-level prediction, making the detection process transparent to clinical users.

2. We demonstrate state-of-the-art performance in detecting hallucinations in radiology reports, particularly for clinically significant findings, achieving superior results compared to existing approaches while maintaining interpretability.
3. We provide empirical evidence that model hidden states contain reliable signals for hallucination detection, suggesting a path toward improved semantic fidelity in medical report generation systems.

2 Background and Related Work

2.1 Hallucination Detection

Hallucination detection is the task of identifying AI-generated content that is false, unfounded, or inconsistent with input data. In medical report generation, hallucinations can manifest as either fictional findings or critical omissions, both of which pose significant risks to patient care. Hallucinations are particularly consequential in radiology, as false findings may both initiate unnecessary medical procedures and mask critical pathologies requiring urgent intervention. Current hallucination detection methods for LLMs and LVLMs can be categorized into three approaches based on their required level of access to model parameters: black-box, gray-box, and white-box approaches. Each category presents different tradeoffs between computational complexity, ease of implementation, and detection accuracy. Our work builds upon these foundations by introducing a white-box approach specifically designed for the radiology domain, where the stakes of hallucination detection are high.

2.2 Black-Box Methods

Black-box methods operate without access to model parameters, relying solely on model outputs. Despite this limitation, these methods have gained prominence due to their broad applicability across proprietary models and APIs. Research in this area has pursued two main directions. The first explores models’ self-evaluation capabilities, with Kadavath et al. (2022) and Lin, Hilton, and Evans (2022) demonstrating that LLMs can identify their own hallucinations with reasonable accuracy when explicitly prompted. The second

direction quantifies model uncertainty by analyzing output diversity. Kuhn, Gal, and Farquhar (2023) and Farquhar et al. (2024) showed that semantic entropy—the uncertainty in the meanings of model-generated text—correlates strongly with hallucination likelihood. Building on these foundations, Friel and Sanyal (2023) proposed using auxiliary LLMs as external validators, demonstrating improved detection performance across both open and closed-domain tasks. More recently, Manakul, Liusie, and Gales (2023) introduced SelfCheck GPT, which evaluates generated sequences through mutual entailment analysis across multiple generations, providing a more robust assessment of content reliability. While these methods provide valuable insights, our work demonstrates that access to model parameters can facilitate robust hallucination detection in the medical domain.

2.3 Gray-Box Methods

Gray-box methods leverage access to token-level probability distributions, enabling more precise analysis than black-box approaches while remaining computationally efficient. Traditional metrics in this category include perplexity, token entropy, and mutual information (Fomicheva et al. 2020; Van der Poel, Cotterell, and Meister 2022). However, token-level uncertainty evaluation presents inherent challenges, as large uncertainty values may reflect the presence of multiple viable continuations rather than potential fabrications. To address this limitation, Fadeeva et al. (2024) developed a claim-conditioned scoring system (CCP) that evaluates model generations according to the probability that a semantically equivalent output would be generated instead, providing a more reliable measure of content veracity. While gray-box methods offer valuable insights into language model uncertainty, our approach moves beyond probability distributions to leverage the rich information contained in hidden states, enabling more nuanced hallucination detection capabilities.

2.4 White-Box Methods

White-box methods require complete access to model weights and typically leverage intermediate representations of the model inputs. The predominant approach involves training supervised classifiers on extracted model activations. For instance, Azaria and Mitchell (2023) and Su et al. (2024) employed feedforward neural networks for post-generation hallucination detection, while Alnuhait et al. (2024) developed similar architectures for pre-generation detection and mitigation. In a significant advancement, Kossen et al. (2024) demonstrated that a logistic regression model trained on hidden activations can effectively predict semantic entropy, suggesting that LLMs encode uncertainty information within their internal representations. Building on this insight, Chen et al. (2024a) introduced the EigenScore metric, which analyzes the consistency of hidden state embeddings across multiple generations to provide a robust measure of output reliability. ReXTrust extends previous white-box approaches by using a self-attention mechanism that analyzes sequences of hidden states, thus achieving high fidelity on finding-level hallucination detection.

2.5 Multimodal Methods

Given that LVLMs incorporate language model components, hallucination detection methods developed for LLMs can be naturally extended to LVLMs (Li et al. 2024b). However, the visual component of LVLMs presents unique opportunities to validate generated content against input images, particularly for verifying visual claims. Recognizing this potential, Yin et al. (2023) developed Woodpecker, which employs auxiliary object detectors and visual question-answering models to verify visual assertions in model-generated text. Similarly, Chen et al. (2024c) approached LVLM hallucination detection through independent visual evidence-gathering modules that systematically confirm or contradict claims extracted from model generations. This approach was further refined by Fei et al. (2024), who employed Siamese networks to compare scene graphs extracted from model generations with those derived from input images.

In the medical domain, specialized LVLM hallucination detection methods have emerged to address domain-specific challenges. Chen et al. (2024b) introduced MediHallDetector for categorizing medical hallucinations, demonstrating superior consistency compared to GPT baselines when evaluating medical content. A significant advance was made by Zhang et al. (2024) with RadFlag, an approach that leverages entailment relationships between findings generated at different sampling temperatures to identify potentially hallucinatory content in LVLM-generated medical reports. ReXTrust builds upon and complements RadFlag by providing fine-grained insights at both the token and finding levels; our empirical results demonstrate the benefits of analyzing model hidden states over post-hoc output analysis.

3 Methodology

In this section, we formalize the task of hallucination detection for radiology report generation. A comparative analysis of black-box, gray-box, and white-box approaches by Liu et al. (2024b) indicated that white-box methods generally demonstrate superior performance. Motivated by this result, we develop ReXTrust as a white-box hallucination detection model.

Throughout our experiments, we utilize the MedVersa model (Zhou et al. 2024) to generate and evaluate candidate reports. Given a sequence of input chest X-rays, MedVersa generates a candidate report R , which we decompose into a set of findings $\{s_i\}_{i=1}^n$. A finding is defined as a single claim in R (e.g., “There is pneumonia”), which usually corresponds to a single sentence. In cases where a sentence contains multiple claims, we split the sentence into individual claims. For instance, the sentence “No evidence of pneumonia, pneumothorax, or pleural effusion” would be split into three distinct findings: “No evidence of pneumonia,” “No evidence of pneumothorax,” and “No evidence of pleural effusion.”

3.1 Model Architecture

For each individual finding s_i , ReXTrust predicts a score in the interval $[0, 1]$, interpreted as the probability that s_i contains hallucinatory content. ReXTrust generates these hallucination risk scores through an end-to-end architecture that

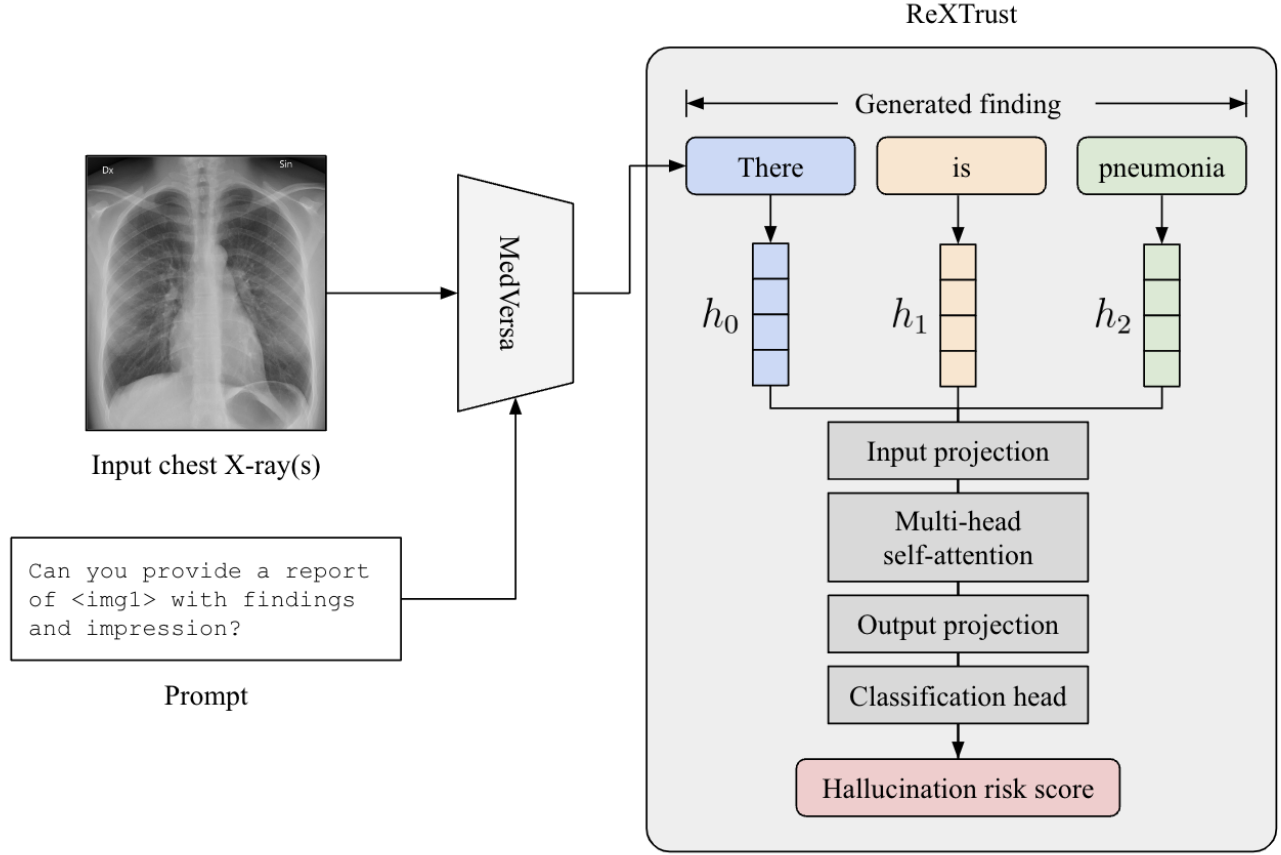


Figure 1: ReXTrust hallucination detection framework. The model processes MedVersa’s hidden states h_i through a self-attention module to produce finding-level hallucination scores.

processes MedVersa’s hidden activation states. The model employs a self-attention module to analyze the sequence of hidden states corresponding to each finding, enabling it to capture both local patterns and broader contextual relationships between tokens when computing the final risk score for s_i . Figure 1 provides a schematic representation of this framework.

For a given finding s_i , ReXTrust first extracts the sequence of hidden states $(h_t^l)_{t \in s_i}$ for all tokens t in s_i , where l denotes the index of a specific hidden layer (we use $l = 16$ in our implementation) and $h_t^l \in \mathbb{R}^d$ ($d = 4096$ for MedVersa). This sequence is projected to a 1024-dimensional latent space through a linear transformation and combined with sinusoidal positional embeddings to encode relative token positions. The model then processes this embedding through three 1024-dimensional linear projections to produce query, key, and value vectors. These vectors are fed through an 8-headed self-attention module with dropout regularization ($p = 0.1$), with each head having dimension 128. The attended outputs are mean-pooled along the sequence dimension, passed through another 1024-dimensional projection, and ultimately fed to a sigmoid classification head that produces a hallucination risk score r .

The advantages of the ReXTrust architecture are two-fold. First, it enables efficient training and inference while maintaining interpretability for clinical users through attention map analysis. Second, the self-attention mechanism naturally captures the full range of token interactions within a finding, allowing the model to simultaneously consider both local features and long-range dependencies when assessing hallucination risk.

3.2 Data

We conduct our experiments using studies from the test set of the MIMIC-CXR database (Johnson et al. 2019). To ensure rigorous evaluation, we randomly partition the subjects into a training set for ReXTrust and a held-out set for final evaluation. The training set contains 231 subjects (corresponding to 1923 studies), while the held-out set contains 58 subjects (424 studies). We maintain strict subject-level separation between all sets to prevent data leakage.

For hallucination labeling at the finding level, we adopt the LLM entailment strategy proposed by Zhang et al. (2024). For each finding generated by MedVersa, we employ OpenAI’s gpt-4o model to classify it as either “completely entailed,” “partially entailed,” or “not entailed” by the cor-

responding ground truth radiologist-written report. We consider a finding to be hallucinated if it is either partially entailed or not entailed by the ground truth report.

3.3 Model Training

ReXTrust is trained using the average binary cross-entropy between predicted hallucination scores and the associated finding-level hallucination labels as the objective function. We employ a weighted random sampler to ensure balanced exposure to positive and negative examples during training, while maintaining equal class weights in the objective function.

The model is optimized for 5 epochs using the AdamW optimizer (Loshchilov 2017) with learning rate 1.0×10^{-4} on a cosine schedule, batch size 128, and standard parameters ($\beta_1 = 0.9$, $\beta_2 = 0.999$) with weight decay 0.01. We employ 5-fold cross-validation on the training set, with model performance monitored via the area under the receiver operating characteristic curve (AUROC) on the validation set. The final weights for ReXTrust are obtained by averaging across folds.

3.4 Finding Severity

Radiological findings exhibit substantial heterogeneity in their clinical significance. For instance, a procedural observation such as “Single portable view of the chest” carries minimal risk if hallucinated, whereas the false positive finding “There is a left apical pneumothorax” could precipitate unnecessary emergency intervention. To systematically analyze this variance in clinical risk, we employ `gpt-4o` to categorize the findings in our evaluation set into four distinct severity tiers:

1. Findings requiring immediate emergency intervention.
2. Findings warranting non-emergency clinical action.
3. Findings with minimal clinical significance.
4. Findings which do not fit into the three previous tiers.

The classification is performed using a standardized prompt (detailed in Appendix A). Our analysis encompasses both the complete evaluation set and a focused evaluation of high-risk findings (tiers 1 and 2), enabling a nuanced assessment of hallucination detection performance in clinically critical scenarios. The reliability of these LLM-generated severity labels is discussed in Section 5.4.

4 Results

4.1 Discriminative Power of ReXTrust

On the complete set of findings, ReXTrust achieves an AUROC of 0.8751 (95% CI: 0.8623, 0.8880). When evaluated exclusively on clinically relevant findings, the AUROC improves to 0.8963 (95% CI: 0.8824, 0.9091). Notably, ReXTrust’s performance remains robust—and even shows improvement—when restricted to clinically significant findings. This suggests that ReXTrust learns generalizable features for hallucination detection rather than overfitting to patterns specific to less critical findings.

4.2 Comparison to Other Approaches

We evaluate ReXTrust against a set of baseline methods spanning diverse approaches to hallucination detection. Traditional uncertainty estimation methods include entropy, which measures token-level predictive uncertainty, and CCP scores (Fadeeva et al. 2024), which evaluate generations based on the probability semantically equivalent alternatives. Among general-purpose multimodal detectors, we evaluate UNIHD (Chen et al. 2024c), which employs independent evidence-gathering modules to validate generated content. Since UNIHD was not originally designed for the medical domain, we implement domain-specific adaptations that preserve its core mechanisms while enabling fair comparison on radiology data (detailed in Appendix B). We also compare against EigenScore (Chen et al. 2024a), a white-box method that assesses hallucination likelihood by analyzing similarities between hidden state activations across multiple model-generated sentences (in our case, we use high-temperature samples from MedVersa). Finally, we compare against RadFlag (Zhang et al. 2024), a radiology-specific method that identifies hallucinations by analyzing entailment relationships between findings generated at different sampling temperatures.

Table 1 presents the comparative evaluation using the micro-averaged AUROC, area under the precision-recall curve (AUPRC), and area under the generalized risk coverage curve (AUGRC) metrics (Traub et al. 2024). ReXTrust demonstrates substantial improvements over all baseline methods across all metrics. The strongest competing method, RadFlag, achieves a finding-level AUROC of 0.7999 (95% CI: 0.7844, 0.8152), which is 7.52 points lower than ReXTrust. ReXTrust’s superior performance relative to UNIHD suggests that training a supervised model directly on in-domain hallucination detection using hidden states provides stronger signals than relying on external verification modules. Furthermore, its improvement over RadFlag indicates that analyzing the generation process through hidden states may be more reliable than post-hoc analysis of model outputs.

To validate our architectural choice of self-attention, we evaluate a variant of ReXTrust that processes hidden states independently for each token rather than attending to the full sequence. In this attention-free version, we generate token-level predictions that are then mean-pooled to produce a finding-level hallucination risk score. While this simpler architecture achieves high performance, it is worse than the full self-attention model. This gap suggests that ReXTrust benefits from modeling the contextual relationships between tokens, capturing dependencies that are difficult to identify through token-level analysis alone.

4.3 Performance Across Medical Categories

Following the categorization approach of Zhang et al. (2024), we refine our analysis by partitioning clinically relevant findings into five distinct categories: Lungs, Pleural, Cardiome-diastinal, Musculoskeletal, and Medical Devices. Figure 2 presents the 95% confidence intervals for ReXTrust’s AUGRC scores (red) on the held-out evaluation set across these

Model	AUROC	AUPRC	AUGRC
Entropy	0.7361 (0.7168, 0.7551)	0.6848 (0.6582, 0.7119)	0.2093 (0.1987, 0.2202)
CCP (Fadeeva et al. 2024)	0.6535 (0.6334, 0.6734)	0.5692 (0.5430, 0.5969)	0.1751 (0.1659, 0.1860)
EigenScore (Chen et al. 2024a)	0.6029 (0.5821, 0.6245)	0.5290 (0.5023, 0.5612)	0.1979 (0.1878, 0.2086)
UniHD (Chen et al. 2024c)	0.6569 (0.6401, 0.6733)	0.5710 (0.5475, 0.5970)	0.1638 (0.1528, 0.1719)
RadFlag (Zhang et al. 2024)	0.7999 (0.7844, 0.8152)	0.7429 (0.7184, 0.7653)	0.0988 (0.0907, 0.1067)
ReXTrust (without self-attention)	0.8683 (0.8555, 0.8808)	0.8186 (0.7958, 0.8431)	0.0646 (0.0582, 0.0707)
ReXTrust	0.8751 (0.8623, 0.8880)	0.8310 (0.8097, 0.8529)	0.0637 (0.0571, 0.0701)

Table 1: AUROC, AUPRC, and AUGRC performance of ReXTrust on the held-out evaluation set, compared to baseline hallucination detection methods. ReXTrust significantly outperforms other methods across all metrics. Parentheses show the 95% CIs computed using a 1,000-sample bootstrap.

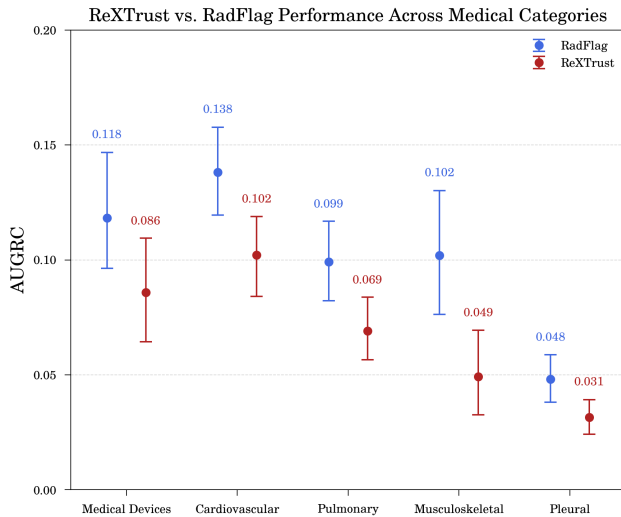


Figure 2: ReXTrust AUGRC performance on five medical categories, compared to RadFlag.

categories, compared against those of RadFlag (blue). ReXTrust demonstrates superior selective classification performance across all categories, as indicated by its lower point estimates. To determine whether ReXTrust’s performance is statistically significantly better than RadFlag’s, we also calculate the confidence intervals on the difference between the AUGRC scores of the two models (Table 2). We find that none of the confidence intervals contain zero, and therefore conclude that the selective classification performance of ReXTrust is statistically significantly better than that of RadFlag at the finding level across all categories.

Table 3 evaluates ReXTrust on findings containing commonly occurring keywords. For each keyword, we show ReXTrust’s micro-averaged F_1 score on clinically relevant findings that contain the keyword and on which the model displays high confidence (i.e., the predicted score is either in the top or bottom score quartile for that category). ReXTrust achieves impressive performance on pleural effusion, pneumonia, and edema findings. However, it performs worse on findings related to pneumothorax, atelectasis, and tubes. We

Category	$\text{AUGRC}_{\text{RadFlag}} - \text{AUGRC}_{\text{ReXTrust}}$
Medical Devices	0.0324 (0.0050, 0.0621)
Cardiovascular	0.0359 (0.0151, 0.0551)
Pulmonary	0.0300 (0.0140, 0.0486)
Musculoskeletal	0.0528 (0.0311, 0.0803)
Pleural	0.0166 (0.0082, 0.0253)

Table 2: Difference in AUGRC between RadFlag and ReXTrust on five medical categories. ReXTrust demonstrates statistically significantly better selective classification performance across all categories.

also note that ReXTrust performs reasonably well when detecting hallucinations in findings containing positional keywords despite not containing explicit visual modules.

Keyword	F_1
“pleural effusion”	0.8403 (0.7847, 0.8958)
“pneumothorax”	0.6193 (0.5455, 0.6932)
“consolidation”	0.7000 (0.5857, 0.8000)
“pneumonia”	0.9091 (0.8182, 0.9773)
“edema”	0.8857 (0.8000, 0.9571)
“atelectasis”	0.6000 (0.4500, 0.7500)
“tube”	0.5750 (0.4000, 0.7250)
“right”	0.6981 (0.6038, 0.7830)
“left”	0.6591 (0.5679, 0.7614)

Table 3: Micro-averaged F_1 scores of ReXTrust on selected finding keywords.

5 Discussion

5.1 Qualitative Examples

Figure 3 presents an analysis of ReXTrust’s performance on four representative studies from the held-out evaluation set. ReXTrust’s self-attention mechanism enables identification

of specific semantic components that contribute to the overall hallucination risk assessment. In the first case, ReXTrust assigns high attention weights to “right” and “pneumothorax,” correctly identifying this positive finding as hallucinatory.

High attention weights are not necessarily indicative of hallucination risk, however. For instance, the third model-generated finding indicates the lack of an anomaly; ReXTrust’s attention weights are highest on the “pneumothorax” tokens, but its final prediction indicates that the finding is true. In the fourth finding, it is the converse; ReXTrust focuses on the “consolidation” tokens, but the model falsely predicts that the finding is not hallucinatory.

5.2 Complementarity with RadFlag

To investigate whether ReXTrust and RadFlag capture complementary aspects of hallucination detection, we evaluate a linear ensemble that combines their predictions with weights 0.8 and 0.2, respectively. The ensemble achieves an AUROC of 0.8822 (95% CI: 0.8699, 0.8947), an AUPRC of 0.8471 (95% CI: 0.8276, 0.8670), and an AUGRC of 0.0594 (95% CI: 0.0532, 0.0656), surpassing the performance of both individual models.

The complementary nature of these approaches stems from their fundamentally different detection strategies: ReXTrust analyzes model hidden states to identify potential hallucinations during the generation process, while RadFlag employs temperature-based sampling to detect inconsistencies in model outputs. Furthermore, the success of the ensemble suggests that ReXTrust could be extended to incorporate hidden states from outputs sampled at high temperatures.

5.3 Generalizability Across Architectures

While ReXTrust was trained to detect hallucinations in MedVersa-generated reports using MedVersa hidden states, this framework can be readily adapted to other medical LLMs and LVLs, such as RaDialog (Pellegrini et al. 2023) and LLaVA-Med (Li et al. 2024a), and Maira-2 (Bannur et al. 2024). Although the architecture of ReXTrust would need to be modified to accommodate different hidden state dimensions, the core principle of applying a self-attention module over sequences of hidden states is applicable. As such, models like ReXTrust could serve as valuable tools for assessing the factual reliability of large medical language models. If hallucinations can be reliably predicted from a model’s hidden states, it suggests that modifications to the training methodology may be necessary to reduce the frequency of hallucinatory content.

5.4 Label Reliability

Our study relies on LLM-generated hallucination and severity labels, which inherently contain some degree of noise. Previous work by Zhang et al. (2024) evaluated similar hallucination labels through clinical validation, finding high but imperfect agreement between LLM and clinician assessments. We perform a similar evaluation of the severity labels in collaboration with an expert clinician. From a sample of 60 findings across 39 studies—with 15 findings from each of the four severity categories defined in Section 3.4—we

find that 8 of the 30 findings categorized into tiers 3 and 4 by the LLM were classified as belonging to tiers 1 and 2 by the clinician. Conversely, none of the 30 findings categorized into tiers 1 and 2 by the LLM were reclassified into tiers 3 and 4 by the clinician. These results suggest that we likely underestimate the number of clinically significant findings. However, given ReXTrust’s consistent performance across finding categories, we expect its performance on the true set of clinically significant findings to align with the global results presented in Table 1.

5.5 Limitations and Future Work

We identify three primary limitations of this study. First, ReXTrust depends on supervised labels. Our work leverages expert-written radiology reports from MIMIC-CXR to generate binary hallucination labels for training. In scenarios where high-quality ground truth reports are unavailable, ReXTrust’s white-box approach may be less suitable than unsupervised alternatives. However, given that ReXTrust achieves high performance even when trained on a small dataset, it may still be useful in label-limited settings.

Second, ReXTrust shows suboptimal performance on certain types of radiological findings (see Table 3). While this weakness might be partially mitigated through ensembling with orthogonal approaches, as discussed in Section 5.2, future work should incorporate more explicit visual grounding tools to verify the factuality of assertions in AI-generated findings (He et al. 2024; Zou et al. 2024; Shaaban, Khan, and Yaqub 2024; Bannur et al. 2024).

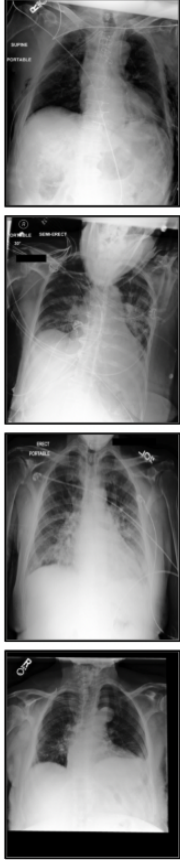
Third, we note that several benchmarks used for comparison in Section 4.2 were not originally designed for radiology (e.g., UNIH). Although we implemented domain-appropriate modifications to ensure fair comparison with our approach, we may underestimate their potential performance in the medical domain. Future work should explore extensive adaptations of general-purpose hallucination detectors to medical imaging tasks.

6 Conclusion

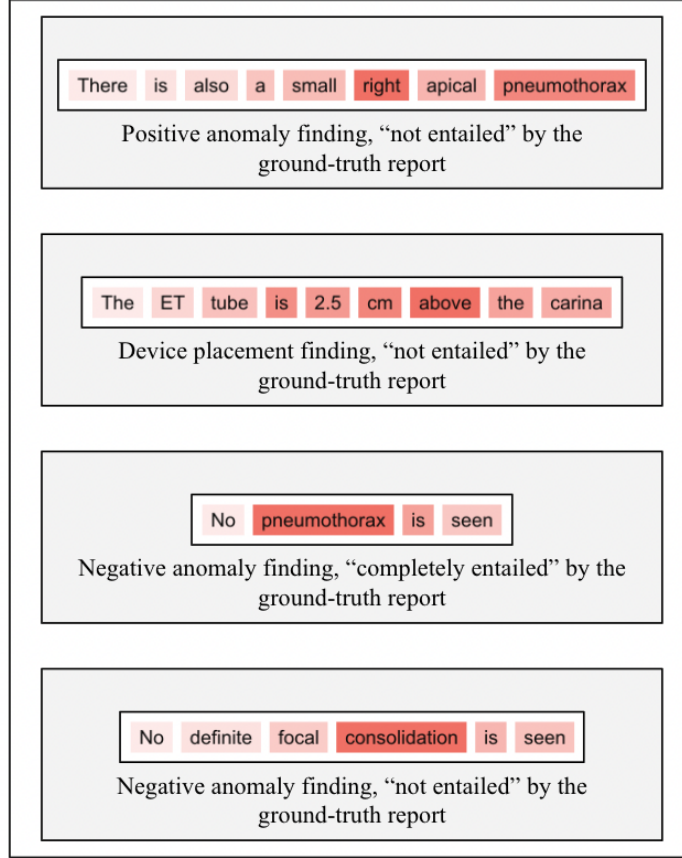
We have presented ReXTrust, a white-box framework for detecting hallucinations in AI-generated radiology reports. Our approach demonstrates superior performance compared to existing hallucination detection methods, achieving an AUROC of 0.8751 across all findings and 0.8963 on clinically significant findings. ReXTrust’s end-to-end architecture enables both granular analysis of hallucination patterns and reliable overall risk assessment.

Our results indicate that model hidden states contain valuable signals about the reliability of generated content that can be effectively leveraged for hallucination detection. The robust performance of ReXTrust on clinically significant findings is particularly noteworthy, as it demonstrates the model’s ability to identify potentially harmful hallucinations that could impact patient care. Our work also identifies important areas for future research, such as improving hallucination detection performance on specific classes of findings and developing methods to reduce dependence on supervised labels. Nevertheless, ReXTrust represents a significant

Chest X-rays



Token-Level Attention Weights



Finding-Level Detection

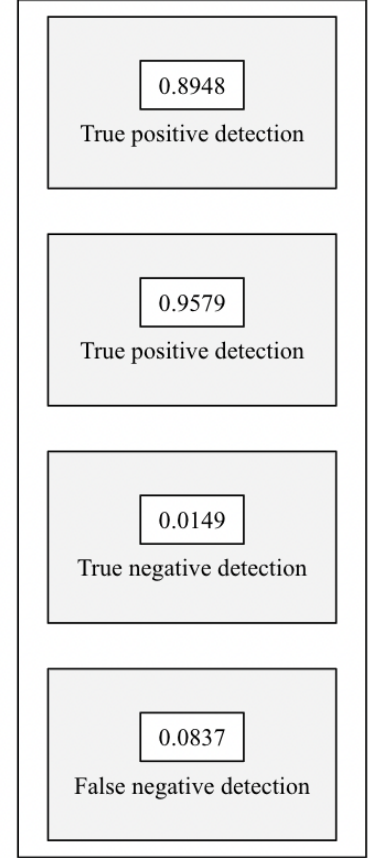


Figure 3: Qualitative examples of ReXTrust on four studies from our held-out evaluation set. Shown on the left are examples of chest X-rays serving as inputs to MedVersa. Shown in the middle are findings generated by MedVersa, as well as token-level attention visualizations from ReXTrust. Shown on the right are the finding-level risk scores output by the classification head of ReXTrust.

advance in ensuring the reliability of medical report generation models. As medical LVLMs continue to evolve and see increased clinical adoption, approaches like ReXTrust will be essential for maintaining high standards of accuracy and safety in medical AI systems.

A Finding Severity Prompt

We present the complete prompt template used to categorize AI-generated findings into severity tiers in Figure 4.

B Adapting UNIHD to Radiology

Noting that UNIHD is not adapted to radiology images and reports, we apply sensible adaptations in order to provide a fair comparison with ReXTrust. Specifically, we modify the prompts of the individual modules to include radiological terminology and examples. We also replace the authors' Grounding Dino object detection module (Liu et al. 2023) with a custom module based on pre-trained models from Cohen et al. (2022).

C Ablation Studies

C.1 Hidden Layer Index

We test the performance of ReXTrust as a function of the hidden layer index of MedVersa. Figure 5 shows the AUROC (blue) and AUPRC (orange) scores of ReXTrust trained on hidden states h_l^l as we let l vary from 0 to 32. We find that the performance of ReXTrust saturates near layer 16.

References

- Alnuhait, D.; Kirtane, N.; Khalifa, M.; and Peng, H. 2024. FactCheckmate: Preemptively Detecting and Mitigating Hallucinations in LMs. *arXiv preprint arXiv:2410.02899*.
- Azaria, A.; and Mitchell, T. 2023. The internal state of an LLM knows when it's lying. *arXiv preprint arXiv:2304.13734*.
- Bannur, S.; Bouzid, K.; Castro, D. C.; Schwaighofer, A.; Bond-Taylor, S.; Ilse, M.; Pérez-García, F.; Salvatelli, V.; Sharma, H.; Meissen, F.; et al. 2024. MAIRA-2:

You are an AI radiology assistant helping process reports from chest X-rays. You are given a finding (F) from a chest X-ray report. Your job is to judge and categorize the clinical severity of the finding. For positive findings, here are the definitions for the severity categories you should use:

1. Emergency clinical consequence

- Definition: Findings which could lead to outcomes requiring immediate or urgent medical intervention to prevent significant harm, deterioration, or life-threatening complications. This includes scenarios where rapid diagnosis and management are critical to stabilize the patient but excludes conditions that can be managed electively.
- Examples: Tension pneumothorax, acute myocardial infarction, Endotracheal tube misplacement compromising airway function.

2. Non-emergency but actionable clinical consequence

- Definition: Findings which require medical attention or follow-up but do not necessitate immediate or urgent intervention to prevent significant harm. These findings are important for patient health management and should be addressed through routine clinical care, elective procedures, or scheduled follow-ups.
- Examples: Newly diagnosed malignancy without acute complications, chronic infection like tuberculosis, significant but non-acute organ abnormalities such as a thyroid nodule needing further investigation.

3. Clinically insignificant consequence

- Definition: Findings that pose no health risk or adverse outcome for the patient and require no clinical intervention or monitoring. These findings are typically benign, normal variations, or imaging artifacts that do not impact patient management or treatment plans.
- Examples: Normal anatomical variations (e.g., azygos fissure), incidental imaging artifacts, findings that are clinically irrelevant to patient health and do not necessitate any follow-up.

4. Other

- Definition: Any findings that do not fit into the above categories.

For negative findings, you should evaluate the severity of the negations of those findings. For example, if you are given the finding, "There is no pneumothorax," you should evaluate the severity of the negation, "There is pneumothorax," and return the corresponding severity category (which would be "emergency clinical consequence" in this case).

Return a JSON of the format: {"severity_category": <category>, "reason": <reason>}, where <category> is one of "emergency clinical consequence", "non-emergency but actionable clinical consequence", "clinically insignificant consequence", or "other", and <reason> is a string explaining your reasoning for the severity category.

Figure 4: Prompt template for radiological finding severity classification. The prompt provides detailed definitions and examples for each severity category to ensure consistent classification across all experiments.

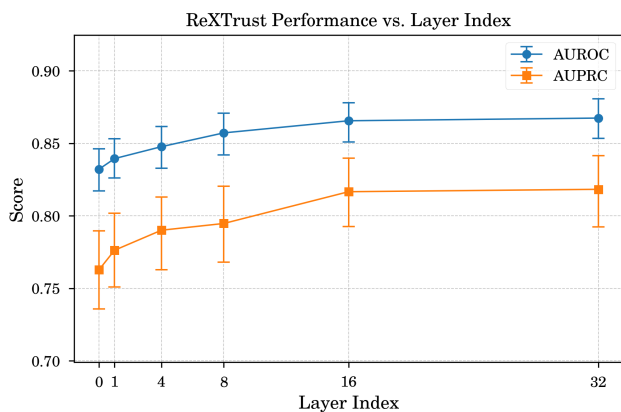


Figure 5: AUROC and AUPRC performance of ReXTrust as a function of the layer index of the MedVersa hidden states used as inputs. Model performance saturates near layer 16. Error bars indicate 95% confidence intervals.

Grounded Radiology Report Generation. *arXiv preprint arXiv:2406.04449*.

Chen, C.; Liu, K.; Chen, Z.; Gu, Y.; Wu, Y.; Tao, M.; Fu, Z.; and Ye, J. 2024a. INSIDE: LLMs’ Internal States Retain the Power of Hallucination Detection. *arXiv preprint arXiv:2402.03744*.

Chen, J.; Yang, D.; Wu, T.; Jiang, Y.; Hou, X.; Li, M.; Wang, S.; Xiao, D.; Li, K.; and Zhang, L. 2024b. Detecting and Evaluating Medical Hallucinations in Large Vision Language Models. *arXiv preprint arXiv:2406.10185*.

Chen, X.; Wang, C.; Xue, Y.; Zhang, N.; Yang, X.; Li, Q.; Shen, Y.; Liang, L.; Gu, J.; and Chen, H. 2024c. Unified hallucination detection for multimodal large language models. *arXiv preprint arXiv:2402.03190*.

Chen, Z.; Varma, M.; Delbrouck, J.-B.; Paschali, M.; Blanke-meier, L.; Van Veen, D.; Valanarasu, J. M. J.; Youssef, A.; Cohen, J. P.; Reis, E. P.; et al. 2024d. Chexagent: Towards a foundation model for chest x-ray interpretation. *arXiv preprint arXiv:2401.12208*.

Cohen, J. P.; Viviano, J. D.; Bertin, P.; Morrison, P.; Torabian, P.; Guarrera, M.; Lungren, M. P.; Chaudhari, A.; Brooks, R.; Hashir, M.; and Bertrand, H. 2022. TorchXRyVision: A library of chest X-ray datasets and models. In *Medical Imaging with Deep Learning*.

Fadeeva, E.; Rubashevskii, A.; Shelmanov, A.; Petrakov, S.; Li, H.; Mubarak, H.; Tsybalov, E.; Kuzmin, G.; Panchenko, A.; Baldwin, T.; et al. 2024. Fact-checking the output of large language models via token-level uncertainty quantification. *arXiv preprint arXiv:2403.04696*.

Farquhar, S.; Kossen, J.; Kuhn, L.; and Gal, Y. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017): 625–630.

Fei, H.; Luo, M.; Xu, J.; Wu, S.; Ji, W.; Lee, M.-L.; and Hsu, W. 2024. Fine-grained Structural Hallucination Detection for Unified Visual Comprehension and Generation in

Multimodal LLM. In *Proceedings of the 1st ACM Multimedia Workshop on Multi-modal Misinformation Governance in the Era of Foundation Models*, 13–22.

Fomicheva, M.; Sun, S.; Yankovskaya, L.; Blain, F.; Guzmán, F.; Fishel, M.; Aletras, N.; Chaudhary, V.; and Specia, L. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8: 539–555.

Friel, R.; and Sanyal, A. 2023. Chainpoll: A high efficacy method for llm hallucination detection. *arXiv preprint arXiv:2310.18344*.

Gu, T.; Liu, D.; Li, Z.; and Cai, W. 2024. Complex organ mask guided radiology report generation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 7995–8004.

Hamamci, I. E.; Er, S.; and Menze, B. 2024. Ct2rep: Automated radiology report generation for 3d medical imaging. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 476–486. Springer.

He, J.; Li, P.; Liu, G.; and Zhong, S. 2024. Parameter-Efficient Fine-Tuning Medical Multimodal Large Language Models for Medical Visual Grounding. *arXiv preprint arXiv:2410.23822*.

Johnson, A. E.; Pollard, T. J.; Berkowitz, S. J.; Greenbaum, N. R.; Lungren, M. P.; Deng, C.-y.; Mark, R. G.; and Horng, S. 2019. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1): 317.

Kadavath, S.; Conerly, T.; Askell, A.; Henighan, T.; Drain, D.; Perez, E.; Schiefer, N.; Hatfield-Dodds, Z.; DasSarma, N.; Tran-Johnson, E.; et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.

Kossen, J.; Han, J.; Razzak, M.; Schut, L.; Malik, S.; and Gal, Y. 2024. Semantic entropy probes: Robust and cheap hallucination detection in llms. *arXiv preprint arXiv:2406.15927*.

Kuhn, L.; Gal, Y.; and Farquhar, S. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*.

Li, C.; Wong, C.; Zhang, S.; Usuyama, N.; Liu, H.; Yang, J.; Naumann, T.; Poon, H.; and Gao, J. 2024a. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36.

Li, Q.; Lyu, C.; Geng, J.; Zhu, D.; Panov, M.; and Karray, F. 2024b. Reference-free Hallucination Detection for Large Vision-Language Models. *arXiv preprint arXiv:2408.05767*.

Lin, S.; Hilton, J.; and Evans, O. 2022. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*.

Liu, C.; Tian, Y.; Chen, W.; Song, Y.; and Zhang, Y. 2024a. Bootstrapping Large Language Models for Radiology Report Generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 18635–18643.

Liu, L.; Pan, Y.; Li, X.; and Chen, G. 2024b. Uncertainty Estimation and Quantification for LLMs: A Simple Supervised Approach. *arXiv preprint arXiv:2404.15993*.

Liu, S.; Zeng, Z.; Ren, T.; Li, F.; Zhang, H.; Yang, J.; Li, C.; Yang, J.; Su, H.; Zhu, J.; et al. 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.

Loshchilov, I. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Manakul, P.; Liusie, A.; and Gales, M. J. 2023. Self-checkgpt: Zero-resource black-box hallucination detection for generative large language models. *arXiv preprint arXiv:2303.08896*.

Pellegrini, C.; Özsoy, E.; Busam, B.; Navab, N.; and Keicher, M. 2023. RaDialog: A large vision-language model for radiology report generation and conversational assistance. *arXiv preprint arXiv:2311.18681*.

Shaaban, M. A.; Khan, A.; and Yaqub, M. 2024. Med-PromptX: Grounded Multimodal Prompting for Chest X-ray Diagnosis. *arXiv preprint arXiv:2403.15585*.

Su, W.; Wang, C.; Ai, Q.; Hu, Y.; Wu, Z.; Zhou, Y.; and Liu, Y. 2024. Unsupervised real-time hallucination detection based on the internal states of large language models. *arXiv preprint arXiv:2403.06448*.

Tanida, T.; Müller, P.; Kaissis, G.; and Rueckert, D. 2023. Interactive and explainable region-guided radiology report generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7433–7442.

Traub, J.; Bungert, T. J.; Lüth, C. T.; Baumgartner, M.; Maier-Hein, K. H.; Maier-Hein, L.; and Jaeger, P. F. 2024. Overcoming common flaws in the evaluation of selective classification systems. *arXiv preprint arXiv:2407.01032*.

Van der Poel, L.; Cotterell, R.; and Meister, C. 2022. Mutual information alleviates hallucinations in abstractive summarization. *arXiv preprint arXiv:2210.13210*.

Yin, S.; Fu, C.; Zhao, S.; Xu, T.; Wang, H.; Sui, D.; Shen, Y.; Li, K.; Sun, X.; and Chen, E. 2023. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045*.

Zhang, S.; Sambara, S.; ; Banerjee, O.; Acosta, J.; Fahrner, J.; and Rajpurkar, P. 2024. RadFlag: A Black-Box Hallucination Detection Method for Medical Vision Language Models. *arXiv preprint arXiv:2411.00299*.

Zhou, H.-Y.; Adithan, S.; Acosta, J. N.; Topol, E. J.; and Rajpurkar, P. 2024. A Generalist Learner for Multifaceted Medical Image Interpretation. *arXiv preprint arXiv:2405.07988*.

Zou, K.; Bai, Y.; Chen, Z.; Zhou, Y.; Chen, Y.; Ren, K.; Wang, M.; Yuan, X.; Shen, X.; and Fu, H. 2024. MedRG: Medical Report Grounding with Multi-modal Large Language Model. *arXiv preprint arXiv:2404.06798*.