# Clinically Grounded Agent-based Report Evaluation: An Interpretable Metric for Radiology Report Generation

Radiological imaging plays a central role in diagnosis, treatment planning, and clinical decision-making. The emergence of vision-language foundation models has generated enormous interest in the use of artificial intelligence to read imaging and generate reports describing imaging findings. However, automated radiology report generation (RRG) suffers from a difficulty in evaluating generated reports for clinical accuracy, which presents a major challenge to safe deployment. Existing metrics often rely on surface-level similarity and/or behave as black boxes and lack clinical interpretability.

To address this challenge, we developed **ICARE**(**I**nterpretable and **C**linically-grounded **A**gent-based **R**eport **E**valuation), an interpretable evaluation framework that uses large language model agents and dynamic multiple-choice question answering (MCQA). Two agents, assigned either the ground-truth or generated report, generate clinically meaningful MCQA and then quiz each other. The answer agreement on these questions reflects whether critical findings are preserved or consistent, acting as interpretable proxies for clinical precision and recall. By linking evaluation scores to specific clinical question–answer pairs, **ICARE** offers transparent and interpretable assessment of report quality while preserving semantic understanding and scalability.

**ICARE** demonstrates strong alignment with expert radiologists' preferences across report comparisons, outperforming existing metrics such as BLEU, BERTScore, and RadGraph in capturing clinically relevant differences. Extensive clinician validation studies confirm the quality and relevance of generated questions, while controlled perturbation experiments show **ICARE** is sensitive to clinical content and robust across evaluation seeds. Cluster-level analysis further reveals interpretable model error patterns (e.g., omission vs. hallucination of findings), offering insight beyond a single aggregate score.

**ICARE** is fully automated, scalable to large datasets, adaptable to varied imaging modalities and other medical text generation tasks (e.g., pathology reports or discharge summaries), and facilitates post-deployment monitoring even without ground-truth data. By making evaluation transparent and clinically meaningful, **ICARE** supports the development of safer systems in healthcare and strengthens trust in automated radiology report generation.
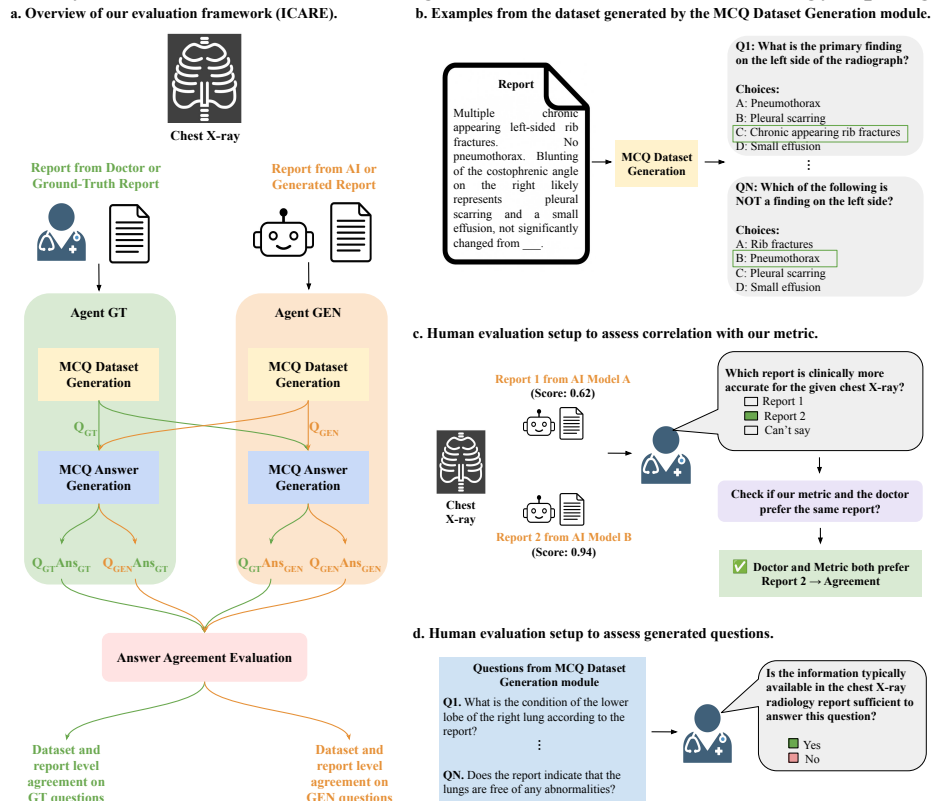


Figure 1: Overview of our evaluation framework and human validation process.