

AUTOMATED STRUCTURED RADIOLOGY REPORT GENERATION WITH RICH CLINICAL CONTEXT

Seongjae Kang^{♣,†} Dong Bok Lee^{♣,†} Juho Jung[♣] Dongseop Kim[♣]

Won Hwa Kim[◇] Sunghoon Joo[♣]

♣VUNO Inc. ♣KAIST ◇POSTECH

{seongjae.kang, juho.jung, dongseop.kim, sunghoon.joo}@vuno.co
markhi@kaist.ac.kr, wonhwa@postech.ac.kr

[†]Equal contribution

ABSTRACT

Automated *structured radiology report generation* (SRRG) from chest X-ray images offers significant potential to reduce workload of radiologists by generating reports in structured formats that ensure clarity, consistency, and adherence to clinical reporting standards. While radiologists effectively utilize available clinical contexts in their diagnostic reasoning, existing SRRG systems *overlook* these essential elements. This fundamental gap leads to critical problems including *temporal hallucinations* when referencing non-existent clinical contexts. To address these limitations, we propose *contextualized SRRG* (C-SRRG) that comprehensively incorporates rich clinical context for SRRG. We curate C-SRRG dataset by integrating comprehensive clinical context encompassing 1) multi-view X-ray images, 2) clinical indication, 3) imaging techniques, and 4) prior studies with corresponding comparisons based on patient histories. Through extensive benchmarking with state-of-the-art multimodal large language models, we demonstrate that incorporating clinical context with the proposed C-SRRG significantly improves report generation quality, as summarized in Fig. 1. We publicly release dataset, code, and checkpoints to facilitate future research for clinically-aligned automated RRG at <https://github.com/vuno/contextualized-srrg>.

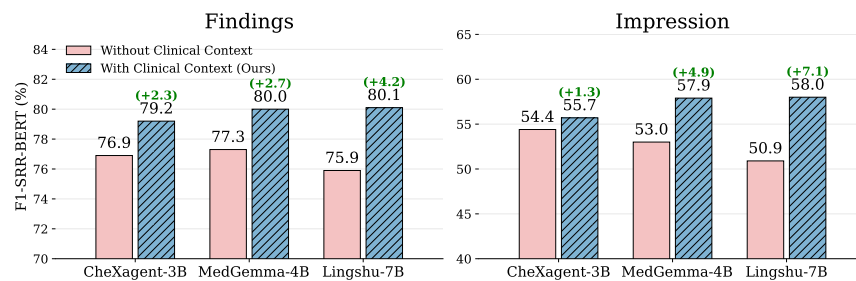


Figure 1: **Clinical context consistently and significantly improves medical MLLMs**—including CheXagent-3B (Chen et al., 2024b), MedGemma-4B (Sellergren et al., 2025), and Lingshu-7B (Team et al., 2025)—on both the findings and impression tasks for SRRG, as measured by F1-SRR-BERT metric (Delbrouck et al., 2025). Clinical context becomes **increasingly critical as MLLMs scale up**, highlighting its importance in RRG.

1 INTRODUCTION

Writing a radiology report requires radiologists to accurately interpret images and synthesize them into two main components: 1) detailed *findings* that systematically document anatomical structures and pathological observations, and 2) concise *impressions* that provide clinical interpretations for subsequent decision-making (Wallis & McCoubrie, 2011; Pahadia et al., 2020; Haygood et al., 2018; Trinh et al., 2019; ESR, 2011). However, generating such comprehensive reports is both cognitively demanding and time-consuming for radiologists. Given the high volume of imaging studies and the time-intensive nature of report writing, there is a critical need for automated systems that can assist radiologists by generating accurate, structured reports while reducing radiologists’ workload and improving diagnostic efficiency (Markotić et al., 2021; Alexander et al., 2022).

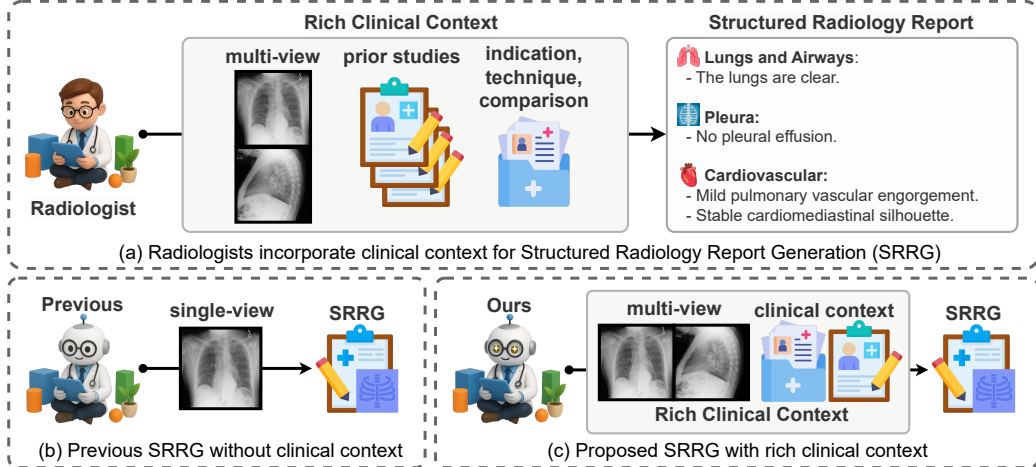


Figure 2: **A conceptual illustration of the proposed C-SRRG.** (a) Radiologists routinely use clinical context, while (b) existing SRRG frameworks do not. Motivated by this gap, (c) C-SRRG leverages multi-view images, indication, technique, and variable-length prior studies/comparisons to generate structured radiology reports.

Automated radiology report generation (RRG) has emerged as a crucial task to address these challenges by assisting radiologists in the diagnostic workflow (Esteva et al., 2019; Sloan et al., 2024; Tanno et al., 2025). Deep learning has accelerated the development of automated RRG frameworks that generate reports directly from medical images (Shin et al., 2016; Jing et al., 2018; Li et al., 2018; Wang et al., 2018; Jing et al., 2020; Chen et al., 2020b; 2022; Wang et al., 2022a). Recent advances in multimodal large language models (MLLMs) further enhanced this capability by integrating vision foundation models with large language models capable of generating coherent and clinically relevant text (Lee et al., 2025; Li et al., 2023a; Hyland et al., 2023; Bannur et al., 2024; Chen et al., 2024b; Sellergren et al., 2025; Team et al., 2025; Wang et al., 2025).

Despite the progress, most automated RRG frameworks overlook *essential clinical context* such as imaging indication, technique, and prior studies that radiologists use to generate reports (Kahn et al., 2009; ESR, 2011). Ignoring the clinical context leads to systematic errors (Liu et al., 2019; Ramesh et al., 2022) as the models fail to capture patient-specific properties and longitudinal changes essential for accurate diagnosis, including *temporal hallucinations* where the models generate references to nonexistent priors or fabricate temporal comparisons (Figs. 22 to 24). Although some work injected partial context—e.g., multi-view images (Yuan et al., 2019; Chen et al., 2022), historical images (Hou et al., 2023a; Zhu et al., 2023b), and prior reports or indications (Miao et al., 2024; Wang et al., 2024)—these are still *limited*, e.g., only consider partial clinical context, rely only on the immediately preceding image or report (Bannur et al., 2024; Liu et al., 2025b;a), with all approaches targeting free-form RRG rather than *structured* RRG (SRRG) (Delbrouck et al., 2025).

To this end, we first present **C-SRRG**, a framework for *contextualized structured radiology report generation* (Fig. 2), built upon the recently introduced SRRG paradigm (Delbrouck et al., 2025). We curate a *large-scale* dataset for *structured* report generation with *rich clinical context*, by leveraging MIMIC-CXR (Johnson et al., 2019) and CheXpert Plus (Chambon et al., 2024). Specifically, our **C-SRRG dataset** provides 1) **multi-view images** (frontal and lateral), 2) clinical **indication**, 3) imaging **technique**, and 4) variable-length **prior studies** with corresponding comparisons, which models can incorporate depending on their architecture.

We evaluate the effectiveness of the proposed C-SRRG with state-of-the-art (SoTA) medical MLLMs—including CheXagent-3B (Chen et al., 2024b), MedGemma-4B (Sellergren et al., 2025), and Lingshu-7B (Team et al., 2025)—and find that incorporating clinical context **substantially and consistently improves report quality** (summarized in Fig. 1 and detailed in Tabs. 3 and 4) measured by various metrics (Papineni et al., 2002; Lin, 2004; Zhang et al., 2019; Delbrouck et al., 2022; 2025). Interestingly, the clinical context **becomes increasingly critical as the models scale up** from 3B to 7B. We also provide a comprehensive analysis, including extensive ablation studies on clinical context (Tabs. 5 to 8), temporal hallucination mitigation (Tab. 9), and organ-level performance (Tab. 10). We will *publicly release* the 1) **dataset**, 2) **code**, and 3) **checkpoints** of benchmarked models to facilitate further research in C-SRRG and benefit the community.

Our contributions and empirical findings are summarized as follows:

- We identify a key limitation of existing SRRG frameworks, *i.e.*, the neglect of essential *clinical context*, which induces systematic errors, most notably *temporal hallucinations* about non-existent prior studies. To address this, we introduce a clinically contextualized SRRG framework (**C-SRRG**) that explicitly integrates clinical context into the generation process.
- We curate *the largest structured* radiology report generation dataset with *rich clinical context*, namely, **C-SRRG dataset**, which includes 1) multi-view images, 2) indication, 3) technique, and 4) prior studies/comparisons, for training and evaluation of the proposed C-SRRG framework.
- As summarized in [Fig. 1](#), we provide a *comprehensive benchmark* of **SoTA MLLM**-based SRRG models, demonstrating that clinical context becomes **increasingly critical as models scale up**—enhancing report quality (*e.g.*, +2.3~4.2/+1.3~7.1 on findings/impression for F1-SRR-BERT) while reducing *temporal hallucinations* ([Tab. 9](#); *e.g.*, 12.2%/18.0% on findings/impression).

2 RELATED WORK

Automated radiology report generation (RRG). Automated RRG has emerged as a promising approach to reduce radiologists’ workload and improve reporting efficiency ([Yang et al., 2023](#); [Sloan et al., 2024](#); [Esteva et al., 2019](#); [Sirshar et al., 2022](#); [Tanno et al., 2025](#); [Singh & Singh, 2025](#)). While early approaches simply combined vision encoders with language decoders for visual feature extraction ([He et al., 2015](#); [Dosovitskiy et al., 2020](#)) and text generation ([Shin et al., 2016](#); [Jing et al., 2018](#); [Li et al., 2018](#); [Wang et al., 2018](#); [Jing et al., 2020](#); [Chen et al., 2020a](#); [Yan & Pei, 2022](#); [Miura et al., 2020](#)), architectural innovations have significantly improved report quality, such as memory-driven transformers ([Chen et al., 2020b](#); [Liu et al., 2024b](#)), specialized architectures for medical domain knowledge ([Yang et al., 2021](#); [Wang et al., 2022b](#); [Kong et al., 2022](#)), cross-modal learning for improved alignment ([Chen et al., 2022](#); [Wang et al., 2022a](#); [Li et al., 2023b](#)), and region-guided frameworks for anatomically relevant features ([Tanida et al., 2023](#); [Li et al., 2023c](#); [Hou et al., 2023b](#)). In this work, we focus on extending the recently proposed structured RRG (SRRG; [Delbrouck et al., 2025](#))—improving clarity, consistency, and interpretability through standardized structure ([Weiss & Langlotz, 2008](#); [Kahn et al., 2009](#); [Bosmans et al., 2012](#); 2015)—by incorporating rich clinical context aligned with radiologists’ workflow.

Multimodal large language models (MLLMs). Building on recent advances in LLMs ([Bai et al., 2023](#); [Achiam et al., 2023](#); [Touvron et al., 2023a;b](#); [Yang et al., 2025](#)), MLLMs have shown strong performance across many domains, including medical applications ([Achiam et al., 2023](#); [Team et al., 2023](#); [Yang et al., 2025](#); [Comanici et al., 2025](#); [Wang et al., 2025](#); [Zhu et al., 2025](#)). They integrate visual understanding with natural-language generation, enabling effective tools for medical image analysis and clinical text generation ([Li et al., 2023a](#); [Chen et al., 2024a](#); [He et al., 2024](#); [Hurst et al., 2024](#); [Lai et al., 2025](#); [Pan et al., 2025](#)). Medical-specific MLLMs further improve performance by incorporating domain knowledge and clinical expertise through specialized training procedures, including CheXagent ([Chen et al., 2024b](#)), MedGemma ([Sellersgren et al., 2025](#)), and Lingshu ([Team et al., 2025](#)). These foundation models are particularly promising for comprehensive RRG frameworks ([Lee et al., 2023](#); [Zhu et al., 2023a](#); [Liu et al., 2024c](#); [Wang et al., 2023](#); [Hyland et al., 2023](#); [Bannur et al., 2024](#); [Lee et al., 2025](#); [Chen et al., 2024b](#); [Dai et al., 2025](#); [Sellersgren et al., 2025](#); [Team et al., 2025](#)), where their ability to accept flexible inputs and produce coherent clinical text is especially valuable. Accordingly, we benchmark medical MLLMs for contextualized SRRG.

Clinical context. Radiologists routinely leverage clinical context when drafting reports, drawing upon patient history, prior studies, and clinical indications ([Kahn et al., 2009](#); [ESR, 2011](#); [Wallis & McCoubrie, 2011](#); [Haygood et al., 2018](#); [Trinh et al., 2019](#); [Pahadia et al., 2020](#); [Castillo et al., 2020](#); [Nguyen et al., 2021](#)), motivating various approaches to integrate such clinical context into automated RRG frameworks. Multi-view image analysis utilizes complementary imaging perspectives, such as frontal and lateral views, to provide comprehensive anatomical coverage ([Yuan et al., 2019](#); [Miao et al., 2024](#); [Chen et al., 2022](#); [Nooralahzadeh et al., 2021](#); [Serra et al., 2023](#); [Liu et al., 2024d](#); [Nicolson et al., 2024](#)). Indication and clinical history integration approaches incorporate patient-specific clinical information ([Hou et al., 2023a](#); [Zhu et al., 2023b](#); [Wang et al., 2024](#); [Liu et al., 2024a](#); [Miao et al., 2024](#)). Previous studies enable temporal comparison and disease progression tracking ([Hou et al., 2023a](#); [Zhu et al., 2023b](#); [Serra et al., 2023](#); [Wang et al., 2024](#); [Liu et al., 2021](#)). Recent works such as MLRG ([Liu et al., 2025b](#)), PriorRG ([Liu et al., 2025a](#)), and MAIRA-2 ([Bannur et al., 2024](#)) have attempted to incorporate clinical context for more comprehensive report generation. However, these approaches have limitations: they either 1) consider only partial clinical

Structured Report (Excerpt):

History: A male patient with hep C cirrhosis and large right pleural effusion status post thoracocentesis.

Comparison: Prior portable AP chest radiograph

Findings:

Pleura:

- Moderate pleural effusion within the right pleural space.
- Moderate right pneumothorax, **new from prior exam.**
- No left pleural effusion or pneumothorax.

Impression:

1. Moderate right-sided pneumothorax.
2. Moderate right pleural effusion.

Hallucination: The phrase “new from prior exam” represents temporal information that cannot be verified from the current study alone, if not with previous history.

Figure 3: **An example of temporal hallucinations.** This report contains “new from prior exam” even though any prior studies are not provided. Please see examples of full structured reports in Figs. 22 to 24.

context, 2) are restricted to specific input configurations, 3) have narrow temporal scope (only the previous prior study), with all 4) focusing on unstructured free-form report generation.

3 METHOD

In this section, we first elaborate on the dataset curation process for contextualized radiology report generation (C-SRRG) in §3.1 and then detail the proposed C-SRRG framework in §3.2.

3.1 CURATION OF CONTEXTUALIZED CLINICAL CONTEXT

Motivation. Our design principle is to reflect the clinical workflow of radiologists that incorporates a diverse diagnostic context such as indication, technique, and comparison (Wallis & McCoubrie, 2011; Trinh et al., 2019; Pahadia et al., 2020; Nguyen et al., 2021), supported by empirical evidence showing improvement in report quality (Castillo et al., 2020; Liu et al., 2021). This emphasis on a comprehensive clinical context aligns with recent work advocating that AI systems must move beyond narrow task-specific approaches that lack the ability to incorporate multimodal data and provide comprehensive interpretation assistance (Dogra et al., 2025). Most importantly, without this context, existing automated systems are prone to *temporal hallucinations*: ground truth reports frequently contain temporal statements such as “new from prior exam” (as shown in Fig. 3), which leads models to hallucinate by referencing nonexistent prior examination (Ramesh et al., 2022). When trained on such data, the SRRG frameworks learn to generate these temporal phrases even when no prior studies are available, as demonstrated in Figs. 25 to 30.

Clinical context. To address this limitation, we incorporate rich clinical context into automatic SRRG frameworks. Specifically, we consider **four clinical elements** that radiologists routinely use:

1. **Multi-view images** (e.g., posteroanterior, anteroposterior, and lateral) provide complementary perspectives from different angles, enabling comprehensive assessment and detection of abnormalities that may be obscured in single views. Multi-view fusion captures richer information through cross-view consistency, improves pathological localization accuracy, and reduces diagnostic uncertainty (Yuan et al., 2019; Miao et al., 2024).
2. **Indication** conveys the clinical rationale for imaging, providing context about patient symptoms, suspected conditions, or clinical questions. This enables models to focus on specific diagnostic questions, tailor findings to physician concerns, and avoid clinically insignificant findings.
3. **Technique** documents examination parameters and limitations including imaging protocols, contrast use, and factors affecting image quality. It helps models note technical caveats, avoid mistaking artifacts for pathology, and prevent duplicate exams.
4. **Prior studies**, when available, enable temporal **comparison** by providing a history to detect disease progression, treatment response, and interval changes. Radiologists routinely consult such a history (Haygood et al., 2018; Liu et al., 2025a), which supports accurate change detection and prevents hallucinations referencing nonexistent prior exams.

Table 1: **Dataset statistics** for C-SRRG-Findings and C-SRRG-Improvement.

Tasks	Train	Valid	Test	Test-reviewed	Total
Findings	181,874	976	1,459	233	184,542
Improvement	405,972	1,505	2,219	231	409,927

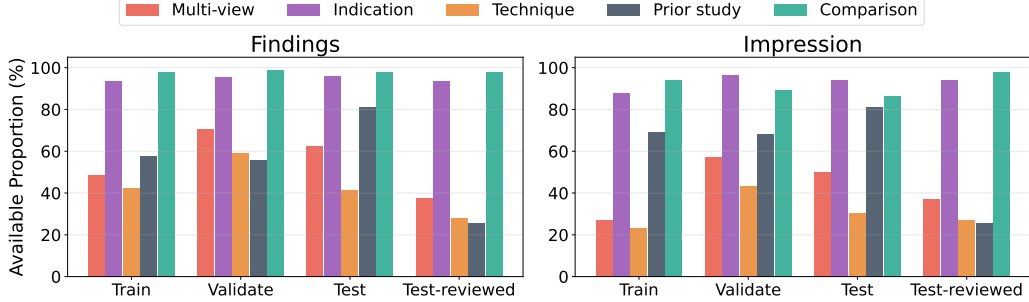


Figure 4: **Available proportion of clinical context** for each split in findings and improvement.



Figure 5: **Distribution of the number of prior studies** available per sample for findings and improvement.

Dataset curation. We build on the recently proposed SRRG dataset (Delbrouck et al., 2025), which includes both MIMIC (Johnson et al., 2019) and CheXpert Plus (Chambon et al., 2024). We employ dataset-specific approaches to extract the necessary clinical context. When multiple views are available, we integrate multi-view images using `ViewPosition` for MIMIC and `frontal_lateral`, `ap_pa` for CheXpert Plus. For MIMIC, each patient is identified by a unique `subject_id` with associated `StudyDate` and `StudyTime` fields. We group patients by `subject_id`, then use temporally-ordered `StudyDate` and `StudyTime` to establish chronological sequences. For CheXpert Plus, each study contains a `deid_patient_id` and `patient_report_date_order` field. We group studies by patient and use the order of report dates to form longitudinal sequences. For other clinical contexts (indication, technique, comparison), we use SRRG components, parsed from free-form reports using GPT-4 (Achiam et al., 2023).

Dataset analysis. Accordingly, we curate two C-SRRG tasks that mirror clinical practice—**C-SRRG-Findings** and **C-SRRG-Improvement** with **train**, **valid**, **test**, and **test-reviewed**¹ splits, as summarized in Tab. 1. The splits enforce the strict separation of patients between training and evaluation to prevent data leakage and properly assess generalization, as in Fig. 10. The availability of a clinical context varies across splits and tasks (Fig. 4). The availability of prior studies follows a

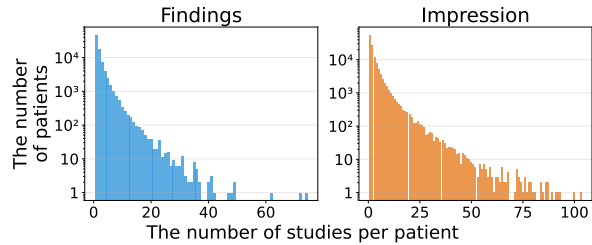


Figure 6: **Distribution of the number of studies per patient.**

¹The *test-reviewed* split reports are reviewed by board-certified radiologists (Delbrouck et al., 2025).

SYSTEM PROMPT:

You are an expert radiologist.

USER PROMPT:

Analyze the current chest X-ray images and compare them with the previous studies to write the IMPRESSION section of a radiology report. Provide a concise clinical summary and diagnosis, noting any changes from the prior studies, focusing on the most recent comparisons. Consider the available clinical contexts when formulating your impression.

Figure 7: **Prompt format for IMPRESSION** for C-SRRG.

=== CURRENT CLINICAL HISTORY/INDICATION ===

Male with end-stage renal disease on hemodialysis, multiple orthopedic hardware, bioprosthetic aortic valve replacement, coronary artery disease status post coronary artery bypass grafting, permanent pacemaker for sick sinus syndrome, admitted for MRSA bacteremia, transferred to the critical care unit for hypotension during anesthesia induction. Patient intubated for respiratory status assessment.

=== CURRENT TECHNIQUE ===

Chest single view

=== CURRENT COMPARISON ===

Prior imaging at an unspecified time.

=== CURRENT IMAGES ===

[Images: Current chest X-ray images]

Figure 8: **An example of current study** for C-SRRG.

ASSISTANT RESPONSE:**IMPRESSION:**

1. Worsened left basilar opacity, which may represent atelectasis or pneumonitis in the appropriate clinical context.
2. Mild improvement in pulmonary vascularity.

Figure 9: **An example of ground-truth assistant responses** in the C-SRRG-Impression dataset.

long-tailed distribution (Fig. 5), alongside the long-tailed counts of studies per patient (Fig. 6)—from no history to extensive longitudinal sequences. This variability reflects *real-world clinical practice* and requires models to handle missing information while leveraging available context.

3.2 CONTEXTUALIZED RADIOLOGY REPORT GENERATION (C-SRRG)

Prompt design. We construct prompt templates for four core settings: 1) findings prompts with prior studies (Fig. 11), 2) findings prompts without prior studies (Fig. 12), 3) impression prompts with prior studies (Fig. 7), and 4) impression prompts without prior studies (Fig. 13). Each prompt consists of the clinical context (*i.e.*, indication, technique, and comparison for the current study), and the associated images (Fig. 8). When available, it incorporates prior studies that also include indication, technique, comparison, and reports on findings or impression (Fig. 14). These structured components are concatenated to form a *single multimodal token sequence*. As shown in Figs. 9 and 15, the response format is standardized for the generation of structured reports. Detailed examples of prompt structures and integration of clinical context can be found in §D and §E.

Training and inference. We fine-tune medical MLLMs on these contextualized prompt–response pairs for both findings and impression tasks. Models receive prompt and clinical context to form a unified multimodal input sequence. The training objective is then the next-token prediction task under an autoregressive language modeling loss: $\frac{1}{T} \sum_{t=1}^T -\log p_{\theta}(y_t|x, y_{<t})$, where x is the multimodal token sequence comprising the prompt (Figs. 7 and 11 to 13), the clinical context of the current study (Fig. 8) and any prior studies (Fig. 14), and $y_{1:T}$ is the target token sequence (*e.g.*, reports on findings or impression; Figs. 9 and 15). Here, p_{θ} denotes the MLLM parameterized by θ . We minimize the negative log-likelihood with respect to θ , *i.e.*, standard cross-entropy over the vocabulary. If prior studies are available, they are inserted into designated slots; otherwise, the model receives only the clinical context of the current study. This design allows the model to adapt to *heterogeneous clinical contexts* (Figs. 4 to 6), to produce context-aware reports when there is prior information, and to avoid hallucinated temporal comparisons when not.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUPS

Implementation details. We evaluate **CheXagent-3B** (Chen et al., 2024b), **MedGemma-4B** (Søllergren et al., 2025), and **Lingshu-7B** (Team et al., 2025). We first train baseline models without clinical context, generating reports directly from single image. When training with C-SRRG, we use all available clinical context (e.g., indication, prior studies). The only exception is CheXagent-3B on the C-SRRG-Impression, where we use only *indication* due to training failure (detailed in §B). We consider the two most recent prior studies with limited number of images (2/3/2 for CheXagent-3B/MedGemma-4B/Lingshu-7B) due to computational constraints. We apply LoRA (Hu et al., 2022) for fine-tuning, optimizing with Adam (Kingma & Ba, 2014), and use vLLM (Kwon et al., 2023) for inference. Greedy decoding is adopted for *reproducibility* consistent with benchmarking purpose. All experiments run on a single NVIDIA H100 GPU. Detailed hyperparameter settings are in Tab. 2.

Evaluation metrics. We use standard metrics, such as **BLEU** (Papineni et al., 2002), **ROUGE-L** (Lin, 2004), and **BERTScore** (Zhang et al., 2019), to assess text quality. For clinical accuracy, we report **F1-RadGraph** (Delbrouck et al., 2022) and SRRG-specific metrics (Delbrouck et al., 2025): **F1-SRRG-BERT**, built on CXR-BERT (Boecking et al., 2022) for structured evaluation, and **Category Score** (only for findings) for the correctness of organ-section headers.

Table 2: Summary of hyperparameters.

Name	Value
<i>LoRA</i>	
Rank r	32
α	64
Dropout p	0.1
<i>Training</i>	
Batch size	128
Optimizer	Adam
Epochs	1
Learning rate	2e-4
LR scheduler	Cosine
Warmup ratio	3%
<i>Inference</i>	
Package	vLLM
Strategy	Greedy

4.2 EXPERIMENTAL RESULTS

Table 3: Results on the C-SRRG-Findings. Clinical context is incorporated with our C-SRRG framework.

Model	Clinical Context	Split	Traditional Metrics				F1-SRR-BERT			Category Score		
			BLEU	ROUGE-L	BERT Score	F1-RadGraph	Precision	Recall	F1-Score	Precision	Recall	F1-Score
CheXagent-3B	✗	Valid	1.97	20.63	30.33	13.07	44.67	45.16	43.46	73.61	81.17	75.54
		Test	2.08	20.09	31.91	12.99	43.73	42.54	41.70	74.47	85.26	77.74
		Test-reviewed	2.13	20.38	32.73	12.96	44.94	42.78	42.31	72.84	87.35	77.55
	✓	Valid	2.31	23.01	33.46	15.76	48.73	48.20	46.79	77.58	83.46	78.73
		Test	1.89	20.92	33.28	13.58	45.18	44.10	43.07	75.79	85.69	78.82
		Test-reviewed	1.98	21.64	34.32	14.05	47.50	45.09	44.59	76.08	88.87	79.93
MedGemma-4B	✗	Valid	1.51	20.95	30.83	13.98	42.93	45.50	42.12	78.48	78.00	76.26
		Test	1.58	19.69	31.52	13.30	42.32	41.38	40.19	76.31	82.36	77.44
		Test-reviewed	1.60	20.11	32.61	13.42	44.49	42.94	41.92	75.39	86.56	78.24
	✓	Valid	4.98	27.22	37.87	20.44	50.52	49.68	48.42	80.38	83.73	80.35
		Test	3.05	23.17	35.65	15.91	45.84	44.24	43.43	78.28	84.67	79.59
		Test-reviewed	4.29	24.37	36.60	17.01	47.90	45.17	44.96	76.73	87.84	80.04
Lingshu-7B	✗	Valid	1.42	17.68	27.20	10.56	40.15	41.45	39.29	74.37	75.97	73.57
		Test	1.40	17.71	29.65	11.14	40.60	39.41	38.65	75.86	81.47	76.84
		Test-reviewed	1.60	18.62	31.09	12.09	42.85	40.82	40.37	74.32	85.20	77.39
	✓	Valid	6.02	28.70	38.85	21.67	51.16	50.50	49.20	81.97	83.03	80.87
		Test	3.16	23.53	35.60	16.07	45.96	44.42	43.63	79.80	83.20	79.68
		Test-reviewed	4.42	23.70	35.76	16.09	47.48	44.80	44.54	77.57	86.71	79.83

Results on the C-SRRG-Findings. Tab. 3 demonstrates **substantial improvements** achieved by C-SRRG on the C-SRRG-Findings across all evaluation metrics, except for slight BLEU decreases for CheXagent-3B on the test/test-reviewed splits (-0.19/-0.15). For example, F1-SRR-BERT scores improve by **+3.33/+1.37/+2.28** (CheXagent-3B), **+6.30/+3.24/+3.04** (MedGemma-4B), and **+9.91/+4.98/+4.17** (Lingshu-7B), with **larger models consistently showing greater gains**. Category Score performance likewise improves by **+3.19/+0.99/+2.38** (CheXagent-3B), **+4.09/+2.15/+1.80** (MedGemma-4B), and **+7.30/+2.84/+2.44** (Lingshu-7B).

Results on the C-SRRG-Impression. Tab. 4 also shows **significant gains** achieved by our C-SRRG on the C-SRRG-Impression. F1-SRR-BERT improves by **+0.8/+3.12** (CheXagent-3B), **+5.5/+4.43/+4.69** (MedGemma-4B), and **+7.42/+7.68/+6.16** (Lingshu-7B) except for CheXagent-3B on the valid split (-0.06). CheXagent-3B also exhibits similar BLEU score decreases on the C-SRRG-Findings (Tab. 3), indicating that rich clinical context may compromise text generation fluency in smaller models. Importantly, while performance consistently drops as **models scale up without clinical context** from 3B to 7B parameters, it **improves substantially with context**, suggesting the **critical importance of clinical context** in scaling up MLLMs for SRRG.

Table 4: **Results on the C-SRRG-Impression.** Clinical context is incorporated with our C-SRRG framework.

Model	Clinical Context	Split	Traditional Metrics				F1-SRR-BERT		
			BLEU	ROUGE-L	BERT Score	F1-RadGraph	Precision	Recall	F1-Score
CheXagent-3B	✗	Valid	9.44	34.03	61.82	19.30	63.80	63.48	59.10
		Test	7.83	29.40	59.82	16.13	57.18	59.18	54.27
		Test-reviewed	7.42	28.60	58.35	13.71	51.32	56.34	49.74
	✓	Valid	7.52	32.99	60.93	17.90	66.28	61.75	59.04
		Test	7.03	29.18	59.66	16.07	59.69	58.42	55.07
		Test-reviewed	6.92	29.04	58.91	14.84	55.42	58.26	52.86
MedGemma-4B	✗	Valid	8.92	41.24	60.94	17.80	62.19	60.77	56.81
		Test	7.15	37.84	59.09	15.35	56.27	57.01	52.69
		Test-reviewed	7.57	35.91	58.35	14.57	51.69	54.42	49.51
	✓	Valid	11.76	46.26	64.28	24.25	65.78	66.77	62.31
		Test	10.58	41.92	61.85	19.23	59.45	61.89	57.12
		Test-reviewed	11.21	40.15	61.12	19.16	55.02	60.71	54.20
Lingshu-7B	✗	Valid	8.15	32.17	59.15	17.23	63.82	57.10	55.06
		Test	6.65	27.27	57.18	13.87	56.03	51.55	49.33
		Test-reviewed	7.04	27.70	57.37	13.49	52.34	52.85	48.37
	✓	Valid	11.77	38.46	64.82	25.29	69.42	63.57	62.48
		Test	10.58	32.86	62.07	19.85	63.04	58.39	57.01
		Test-reviewed	11.61	33.66	62.04	21.28	57.48	58.80	54.53

Table 5: **Effect of clinical context for train/eval on the C-SRRG-Findings** using MedGemma-4B.

Clinical Context		Split	F1-SRR-BERT		
Train	Eval		Precision	Recall	F1-Score
✗	✗	Valid	42.93	45.50	42.12
		Test	42.32	41.38	40.19
		Test-reviewed	44.49	42.94	41.92
✓	✗	Valid	47.00	47.09	45.35
		Test	42.76	41.55	40.64
		Test-reviewed	44.79	43.12	42.45
✗	✓	Valid	45.28	45.25	43.56
		Test	43.02	40.94	40.44
		Test-reviewed	44.40	41.50	41.36
✓	✓	Valid	50.52	49.68	48.42
		Test	45.84	44.24	43.43
		Test-reviewed	47.90	45.17	44.96

Table 7: **Ablation study on clinical context for the C-SRRG-Findings** using MedGemma-4B.

Configuration	Split	F1-SRR-BERT		
		Precision	Recall	F1-Score
Single-view	Valid	47.00	47.09	45.35
	Test	42.76	41.55	40.64
	Test-reviewed	44.79	43.12	42.45
Multi-view	Valid	47.46	47.21	45.80
	Test	44.44	42.57	41.92
	Test-reviewed	45.49	42.69	42.39
+ Indication	Valid	46.95	46.06	44.85
	Test	44.62	42.56	41.98
	Test-reviewed	45.92	43.14	42.79
+ Technique	Valid	50.35	49.27	48.24
	Test	45.50	43.89	43.15
	Test-reviewed	47.48	44.60	44.42
+ Comparison + Prior studies	Valid	50.52	49.68	48.42
	Test	45.84	44.24	43.43
	Test-reviewed	47.90	45.17	44.96

Table 6: **Effect of clinical context for train/eval on the C-SRR-Impression** using MedGemma-4B.

Clinical Context		Split	F1-SRR-BERT		
Train	Eval		Precision	Recall	F1-Score
✗	✗	Valid	62.19	60.77	56.81
		Test	56.27	57.01	52.69
		Test-reviewed	51.69	54.42	49.51
✓	✗	Valid	63.87	61.86	58.45
		Test	54.42	56.53	51.64
		Test-reviewed	51.45	57.86	51.17
✗	✓	Valid	62.60	64.23	59.10
		Test	53.34	59.11	52.59
		Test-reviewed	49.35	58.68	50.66
✓	✓	Valid	65.78	66.77	62.31
		Test	59.45	61.89	57.12
		Test-reviewed	55.02	60.71	54.20

Table 8: **Ablation study on clinical context for C-SRRG-Impression** using MedGemma-4B.

Configuration	Split	F1-SRR-BERT		
		Precision	Recall	F1-Score
Single-view	Valid	63.87	61.86	58.45
	Test	54.42	56.53	51.64
	Test-reviewed	51.45	57.86	51.17
Multi-view	Valid	65.74	62.92	59.89
	Test	55.70	58.11	53.36
	Test-reviewed	51.78	59.37	52.25
+ Indication	Valid	66.91	65.02	61.67
	Test	58.47	59.11	55.00
	Test-reviewed	53.32	60.47	52.86
+ Technique	Valid	65.91	65.65	61.78
	Test	58.65	60.24	55.88
	Test-reviewed	54.66	60.05	53.39
+ Comparison + Prior studies	Valid	65.78	66.77	62.31
	Test	59.45	61.89	57.12
	Test-reviewed	55.02	60.71	54.20

Effect of clinical context on training/evaluation. We next conduct ablation studies with four settings: 1) train+eval without context (baseline); 2) train with, eval without; 3) train without, eval with; and 4) train+eval with context. **Tabs. 5 and 6** report F1-SRR-BERT on C-SRRG-Findings and C-SRRG-Impression using MedGemma-4B. We find incorporating clinical context in only one phase provides limited improvement or slight degradation: *e.g.*, +3.23/+0.45/+0.53 (train ✓), or +1.44/+0.25/-0.56 (test ✓) in findings and +1.65/-1.05/+1.66 (train ✓), or +2.29/-0.1/+1.15 (test ✓) in impression, which shows the **benefit of using context in both phases** for SRRG performance.

Impact of Each Clinical Context Component. We ablate four clinical-context components, 1) multi-view images, 2) indication, 3) technique, and 4) prior studies with comparison, to isolate their contributions on the performance. **Tabs. 7 and 8** report F1-SRR-BERT for each variant on the C-SRRG-Findings, C-SRRG-Impression, respectively. **All components contribute incrementally** to both tasks (except for few cases), with **performance being highest when using all available context**, which shows the importance of incorporating clinical context in SRRG.

Mitigation of temporal hallucinations.

To quantify temporal hallucinations, we train MedGemma-4B under two conditions: **without** clinical context (**baseline**) and **with** clinical context (**C-SRRG**). We evaluate both on evaluation sets **without clinical context**, and count reports that contain one of the following 33 indicators: 1) time references ('new', 'newly', 'recent', 'recently', 'previous', 'prior', 'interval', 'compared to', 'since', 'from prior'), 2) stability indicators ('unchanged', 'stable', 'persistent', 'persisting'), and 3) change indicators ('improved', 'improvement', 'worsened', 'worsening', 'increased', 'decreased', 'enlarging', 'reducing', 'progression', 'regression', 'evolving', 'evolve', 'developing', 'developed', 'resolving', 'resolved', 'temporal change', 'compare', 'comparison'). By this, we can detect whether the generated reports **contained hallucinations** by identifying inappropriate temporal references in the **absence of clinical context**. Tab. 9 shows that **clinical context substantially mitigates hallucinations**: Findings drop from 22.9% to 10.7% (−12.2%) and Impression from 43.8% to 25.8% (−18.0%). This shows that C-SRRG effectively handles *heterogeneous clinical context availability*, i.e., the absence of clinical context, while successfully mitigating temporal hallucinations.

Table 9: Mitigation effect of temporal hallucination.

Task	Split	Temporal Hallucination Rate		Mitigation
		Baseline (X)	C-SRRG (✓)	
Findings	Valid	146/976 (15.0%)	70/976 (7.2%)	-7.8%
	Test	416/1459 (28.5%)	194/1459 (13.3%)	-15.2%
	Test-reviewed	49/233 (21.0%)	21/233 (9.0%)	-12.0%
	Overall	611/2668 (22.9%)	285/2668 (10.7%)	-12.2%
Impression	Valid	630/1505 (41.9%)	364/1505 (24.2%)	-17.7%
	Test	1012/2219 (45.6%)	599/2219 (27.0%)	-18.6%
	Test-reviewed	92/231 (39.8%)	58/231 (25.1%)	-14.7%
	Overall	1734/3955 (43.8%)	1021/3955 (25.8%)	-18.0%

Anatomical region analysis. We compare organ-level performance for findings task against the baseline using the SRRG anatomical categories (Delbrouck et al., 2025). Tab. 10 reports the Category Score on the validation split using MedGemma-4B, with abbreviations: P = pleura, A = abdominal, H/M = hila/mediastinum, O = Other, L/A = lungs/airways, C = cardiovascular, M/C = musculoskeletal/chest wall, T/C/S = tubes/catheters/support devices. We observe that incorporating clinical context with proposed C-SRRG **improves the performance across all the anatomical regions**, except for H/M.

Table 10: Organ-level Category Score on the Valid split.

Region	Baseline (X)			C-SRRG (✓)		
	Precision	Recall	F1-Score	Precision	Recall	F1-Score
P	47.69	39.62	41.72	58.2	52.39	52.77
A	8.0	8.0	8.0	17.86	17.86	17.86
H/M	37.91	36.55	36.97	34.25	33.07	33.41
O	9.21	7.94	8.2	12.17	10.13	10.6
L/A	40.67	64.67	45.68	57.25	62.24	55.69
C	71.77	68.6	68.99	73.28	70.26	70.51
M/C	26.22	25.26	25.55	43.2	41.97	42.3
T/C/S	57.04	64.18	58.45	59.81	64.39	60.35

5 CONCLUSION, LIMITATIONS, AND FUTURE WORK

We introduced *contextualized structured radiology report generation (C-SRRG)*, a framework that aligns with radiologists’ diagnostic workflow by integrating rich clinical context including multi-view images, indication, imaging technique, and prior studies with comparisons. Through comprehensive evaluation of state-of-the-art medical MLLMs, we demonstrate that clinical context integration consistently enhances text quality, diagnostic accuracy, and reduces temporal hallucinations. Importantly, our findings reveal that clinical context becomes **increasingly critical** as models scale up, suggesting that larger foundation models require more sophisticated contextual integration to achieve optimal performance. We will publicly release our dataset, code, and model checkpoints to foster further research in C-SRRG and benefit the broader community.

Limitations. The C-SRRG dataset relies on synthetic LLM annotations from reformulated reports, which may introduce biases and subtle hallucinations. Our supervised fine-tuning approach with greedy decoding may limit the full capture of clinical reasoning processes. Computational and architectural constraints limited our evaluation to 7B parameter models with restricted multiple image processing capabilities (e.g., Lingshu-7B and CheXagent-3B limited to 2 images) and context windows that constrain comprehensive longitudinal history integration. Additionally, our recency-based selection strategy for prioritizing recent studies, while capturing clinically relevant temporal information, may occasionally omit important historical context.

Future work should explore scaling to larger foundation models with extended-context capabilities, developing intelligent clinical context selection policies through learned strategies or retrieval-augmented approaches, and incorporating preference learning techniques with radiologists’ feedback. Expanding to comprehensive clinical modalities also presents promising avenues for enhanced diagnostic accuracy, with detailed discussions provided in §A.

REPRODUCIBILITY STATEMENT

We ensure reproducibility by building C-SRRG entirely from publicly available MIMIC-CXR (Johnson et al., 2019) and CheXpert-Plus (Chambon et al., 2024) datasets, with detailed documentation of our data processing pipeline including longitudinal patient history extraction and multi-view image integration in §3.1. All experimental configurations are specified in §4.1, including model hyperparameters (learning rate $2e-4$, batch size 128, LoRA rank 32), exact data splits with patient-level separation, and standard evaluation metrics. We commit to publicly releasing our complete codebase, the C-SRRG dataset with clinical context annotations, trained model checkpoints, and documentation for dataset recreation. All experiments use reproducible libraries (Hugging Face PEFT (Mangrulkar et al., 2022), vLLM (Kwon et al., 2023)) on a single NVIDIA H100 GPU with fixed random seeds.

ETHICS STATEMENT

Our work presents no new ethical concerns as C-SRRG is built entirely from existing de-identified public datasets (MIMIC-CXR and CheXpert Plus) that have undergone rigorous de-identification and received appropriate IRB approvals. No additional patient data was collected for this work, and all privacy protections from the source datasets are maintained. We acknowledge that automated report generation systems may produce hallucinations when referencing non-existent prior studies, which our work specifically addresses by incorporating comprehensive clinical context. Our dataset and models are intended solely for research purposes and should not be used for clinical decision-making without appropriate validation and regulatory approval. The computational requirements are modest (single GPU training), minimizing environmental impact while maintaining research accessibility.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Robert Alexander, Stephen Waite, Michael A Bruno, Elizabeth A Krupinski, Leonard Berlin, Stephen Macknik, and Susana Martinez-Conde. Mandating limits on workload, duty, and speed in radiology. *Radiology*, 304(2):274–282, 2022.
- Katherine P Andriole. Picture archiving and communication systems: past, present, and future. *Journal of Medical Imaging*, 10(6):061405–061405, 2023.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Shruthi Bannur, Kenza Bouzid, Daniel C. Castro, Anton Schwaighofer, Sam Bond-Taylor, Maximilian Ilse, Fernando P’erez-Garc’ia, Valentina Salvatelli, Harshita Sharma, Felix Meissen, Mercy Prasanna Ranjit, Shaury Srivastav, Julia Gong, Fabian Falck, Ozan Oktay, Anja Thieme, Matthew P. Lungren, Maria Teodora Wetscherek, Javier Alvarez-Valle, and Stephanie L. Hyland. Maira-2: Grounded radiology report generation. *ArXiv*, abs/2406.04449, 2024. URL <https://api.semanticscholar.org/CorpusID:270357817>.
- Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel Coelho de Castro, Anton Schwaighofer, Stephanie L. Hyland, Maria Teodora Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan Oktay. Making the most of text semantics to improve biomedical vision-language processing. In *European Conference on Computer Vision*, 2022. URL <https://api.semanticscholar.org/CorpusID:248300135>.
- Jan ML Bosmans, Lieve Peremans, Matías Menni, Arthur M. De Schepper, Philippe Duyck, and Paul M. Parizel. Structured reporting: if, why, when, how—and at what expense? results of a focus group meeting of radiology professionals from eight countries. *Insights into Imaging*, 3:295–302, 2012. URL <https://api.semanticscholar.org/CorpusID:14138260>.

-
- Jan ML Bosmans, Emanuele Neri, Osman Ratib, and Charles E Kahn Jr. Structured reporting: a fusion reactor hungry for fuel. *Insights into imaging*, 6(1):129–132, 2015.
- Chelsea Castillo, Tom Steffens, L. Sim, and Liam J. Caffery. The effect of clinical information on radiology reporting: A systematic review. *Journal of Medical Radiation Sciences*, 68:60 – 74, 2020. URL <https://api.semanticscholar.org/CorpusID:221405468>.
- Pierre Chambon, Jean-Benoit Delbrouck, Thomas Sounack, Shih-Cheng Huang, Zhihong Chen, Maya Varma, Steven QH Truong, Curtis P Langlotz, et al. Chexpert plus: Hundreds of thousands of aligned radiology texts, images and patients. *arXiv e-prints*, pp. arXiv–2405, 2024.
- Junying Chen, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, Guangjun Yu, Xiang Wan, and Benyou Wang. Huatuoqpt-vision, towards injecting medical visual knowledge into multimodal llms at scale. *ArXiv*, abs/2406.19280, 2024a. URL <https://api.semanticscholar.org/CorpusID:270764495>.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. *ArXiv*, abs/2010.16056, 2020a. URL <https://api.semanticscholar.org/CorpusID:226222210>.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. *arXiv preprint arXiv:2010.16056*, 2020b.
- Zhihong Chen, Yaling Shen, Yan Song, and Xiang Wan. Cross-modal memory networks for radiology report generation. *arXiv preprint arXiv:2204.13258*, 2022.
- Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier, Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes Reis, Emily Bao Tsai, Andrew Johnston, Cameron Olsen, Tanishq Mathew Abraham, Sergios Gatidis, Akshay S. Chaudhari, and Curtis P. Langlotz. Chexagent: Towards a foundation model for chest x-ray interpretation. *ArXiv*, abs/2401.12208, 2024b. URL <https://api.semanticscholar.org/CorpusID:279241090>.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Wei Dai, Peilin Chen, Chanakya Ekbote, and Paul Pu Liang. Qoq-med: Building multimodal clinical foundation models with domain-aware grpo training. *arXiv preprint arXiv:2506.00711*, 2025.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022.
- Jean-Benoit Delbrouck, Pierre Chambon, Christian Blüthgen, Emily B Tsai, Omar Almusa, and Curt P. Langlotz. Improving the factual correctness of radiology report generation with semantic rewards. In *Conference on Empirical Methods in Natural Language Processing*, 2022. URL <https://api.semanticscholar.org/CorpusID:253098780>.
- Jean-Benoit Delbrouck, Justin Xu, Johannes Moll, Alois Thomas, Zhihong Chen, Sophie Ostmeier, Asfandiyar Azhar, Kelvin Zhenghao Li, Andrew Johnston, Christian Bluethgen, et al. Automated structured radiology report generation. *arXiv preprint arXiv:2505.24223*, 2025.
- Siddhant Dogra, Xiaoman Zhang, Ezequiel Silva, and Pranav Rajpurkar. The financial, operational, and clinical advantages of generalist radiology ai. *Radiology*, 316(3), 2025. doi: 10.1148/radiol.242362. URL <https://pubs.rsna.org/doi/10.1148/radiol.242362>. Published Online: Sep 9 2025.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929, 2020. URL <https://api.semanticscholar.org/CorpusID:225039882>.

-
- ESR. Good practice for radiological reporting. guidelines from the european society of radiology (esr). *Insights into Imaging*, 2:93 – 96, 2011. URL <https://api.semanticscholar.org/CorpusID:13036653>.
- Andre Esteva, Alexandre Robicquet, Bharath Ramsundar, Volodymyr Kuleshov, Mark DePristo, Katherine Chou, Claire Cui, Greg Corrado, Sebastian Thrun, and Jeff Dean. A guide to deep learning in healthcare. *Nature medicine*, 25(1):24–29, 2019.
- Tamara Miner Haygood, Barry Mullins, Jia Sun, Behrang Amini, Priya R Bhosale, Hyunseon Christine Kang, Tara L. Sagebiel, and Bilal Mujtaba. Consultation and citation rates for prior imaging studies and documents in radiology. *Journal of Medical Imaging*, 5, 2018. URL <https://api.semanticscholar.org/CorpusID:21674385>.
- Kristiina Häyrynen, Kaija Saranto, and Pirkko Nykänen. Definition, structure, content, use and impacts of electronic health records: a review of the research literature. *International journal of medical informatics*, 77(5):291–304, 2008.
- Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015. URL <https://api.semanticscholar.org/CorpusID:206594692>.
- Sunan He, Yuxiang Nie, Hongmei Wang, Shu Yang, Yihui Wang, Zhiyuan Cai, Zhixuan Chen, Yingxue Xu, Luyang Luo, Huiling Xiang, et al. Gsco: Towards generalizable ai in medicine via generalist-specialist collaboration. *arXiv preprint arXiv:2404.15127*, 2024.
- Wenjun Hou, Yi Cheng, Kaishuai Xu, Wenjie Li, and Jiang Liu. Recap: Towards precise radiology report generation via dynamic disease progression reasoning. *arXiv preprint arXiv:2310.13864*, 2023a.
- Wenjun Hou, Kaishuai Xu, Yi Cheng, Wenjie Li, and Jiangming Liu. Organ: Observation-guided radiology report generation via tree reasoning. *ArXiv*, abs/2306.06466, 2023b. URL <https://api.semanticscholar.org/CorpusID:259138298>.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- OpenAI Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mkadry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alexander Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alexandre Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codisoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, An drey Mishchenko, Angela Baek, Angela Jiang, An toine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, B. Ghorbani, Ben Leimberger, Ben Rossen, Benjamin Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll L. Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Chris Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mély, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Phong Duc Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Hai-Biao Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Hee woo Jun, Hendrik Kirchner, Henrique Pondé de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai

Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub W. Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Ryan Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quiñero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Joshua Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Ouyang Long, Louis Feuvrier, Lu Zhang, Lukasz Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madeleine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljube, Ma teusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Ali Yatbaz, Mengxue Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Mina Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na talie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nikolas A. Tezak, Niko Felix, Nithanth Kudige, Nitish Shirish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Phil Tillet, Prafulla Dhariwal, Qim ing Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Raphael Gontijo Lopes, Raul Puri, Reah Miyara, Reimar H. Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Ramilevich Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiye Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yuri Malkov. Gpt-4o system card. *ArXiv*, abs/2410.21276, 2024. URL <https://api.semanticscholar.org/CorpusID:273662196>.

Stephanie L. Hyland, Shruthi Bannur, Kenza Bouzid, Daniel C. Castro, Mercy Prasanna Ranjit, Anton Schwaighofer, Fernando Pérez-García, Valentina Salvatelli, Shaury Srivastav, Anja Thieme, Noel Codella, Matthew P. Lungren, Maria Teodora Wetscherek, Ozan Oktay, and Javier Alvarez-Valle. Maira-1: A specialised large multimodal model for radiology report generation. *ArXiv*, abs/2311.13668, 2023. URL <https://api.semanticscholar.org/CorpusID:265445382>.

Baoyu Jing, Pengtao Xie, and Eric Xing. On the automatic generation of medical imaging reports. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2577–2586, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1240. URL <https://aclanthology.org/P18-1240/>.

Baoyu Jing, Zeya Wang, and Eric Xing. Show, describe and conclude: On exploiting the structure information of chest x-ray reports. *arXiv preprint arXiv:2004.12274*, 2020.

-
- Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih ying Deng, Roger G. Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317, 2019.
- Charles E. Kahn, C. Langlotz, Elizabeth S. Burnside, John A. Carrino, David S. Channin, David M. Hovsepian, and D. Rubin. Toward best practices in radiology reporting. *Radiology*, 252 3:852 – 856, 2009. URL <https://api.semanticscholar.org/CorpusID:44982938>.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Ming Kong, Zhengxing Huang, Kun Kuang, Qiang Zhu, and Fei Wu. Transq: Transformer-based semantic query for medical report generation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2022. URL <https://api.semanticscholar.org/CorpusID:252369649>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Haotong Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. *Proceedings of the 29th Symposium on Operating Systems Principles*, 2023. URL <https://api.semanticscholar.org/CorpusID:261697361>.
- Yuxiang Lai, Jike Zhong, Ming Li, Shitian Zhao, and Xiaofen Yang. Med-r1: Reinforcement learning for generalizable medical reasoning in vision-language models. *ArXiv*, abs/2503.13939, 2025. URL <https://api.semanticscholar.org/CorpusID:277104458>.
- Seowoo Lee, Jiwon Youn, Hyungjin Kim, Mansu Kim, and Soon Ho Yoon. Cxr-llava: a multimodal large language model for interpreting chest x-ray images. *European Radiology*, pp. 1–13, 2025.
- Suhyeon Lee, Won Jun Kim, Jinho Chang, and Jong-Chul Ye. Llm-cxr: Instruction-finetuned llm for cxr image understanding and generation. In *International Conference on Learning Representations*, 2023. URL <https://api.semanticscholar.org/CorpusID:258823258>.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *ArXiv*, abs/2306.00890, 2023a. URL <https://api.semanticscholar.org/CorpusID:258999820>.
- Mingjie Li, Bingqian Lin, Zicong Chen, Haokun Lin, Xiaodan Liang, and Xiaojun Chang. Dynamic graph enhanced contrastive learning for chest x-ray report generation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3334–3343, 2023b. URL <https://api.semanticscholar.org/CorpusID:257631847>.
- Yaowei Li, Bang Yang, Xuxin Cheng, Zhihong Zhu, Hongxiang Li, and Yuexian Zou. Unify, align and refine: Multi-level semantic alignment for radiology report generation. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2851–2862, 2023c. URL <https://api.semanticscholar.org/CorpusID:257771453>.
- Yuan Li, Xiaodan Liang, Zhiting Hu, and Eric P Xing. Hybrid retrieval-generation reinforced agent for medical image report generation. *Advances in neural information processing systems*, 31, 2018.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Annual Meeting of the Association for Computational Linguistics*, 2004. URL <https://api.semanticscholar.org/CorpusID:964287>.
- Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. Exploring and distilling posterior and prior knowledge for radiology report generation. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13748–13757, 2021. URL <https://api.semanticscholar.org/CorpusID:235421693>.

-
- Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew B. A. McDermott, Willie Boag, Wei-Hung Weng, Peter Szolovits, and Marzyeh Ghassemi. Clinically accurate chest x-ray report generation. *ArXiv*, abs/1904.02633, 2019. URL <https://api.semanticscholar.org/CorpusID:102353285>.
- Kang Liu, Zhuoqi Ma, Xiaolu Kang, Zhushi Zhong, Zhicheng Jiao, Grayson Baird, Harrison X. Bai, and Qiguang Miao. Structural entities extraction and patient indications incorporation for chest x-ray report generation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2024a. URL <https://api.semanticscholar.org/CorpusID:270045672>.
- Kang Liu, Zhuoqi Ma, Mengmeng Liu, Zhicheng Jiao, Xiaolu Kang, Qiguang Miao, and Kun Xie. Factual serialization enhancement: A key innovation for chest x-ray report generation. *ArXiv*, abs/2405.09586, 2024b. URL <https://api.semanticscholar.org/CorpusID:272600508>.
- Kang Liu, Zhuoqi Ma, Zikang Fang, Yunan Li, Kun Xie, and Qiguang Miao. Priorrg: Prior-guided contrastive pre-training and coarse-to-fine decoding for chest x-ray report generation. *arXiv preprint arXiv:2508.05353*, 2025a.
- Kang Liu, Zhuoqi Ma, Xiaolu Kang, Yunan Li, Kun Xie, Zhicheng Jiao, and Qiguang Miao. Enhanced contrastive learning with multi-view longitudinal data for chest x-ray report generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 10348–10359, 2025b.
- Rui Liu, Mingjie Li, Shen Zhao, Ling Chen, Xiaojun Chang, and Lina Yao. In-context learning for zero-shot medical report generation. *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024c. URL <https://api.semanticscholar.org/CorpusID:273642593>.
- Zhizhe Liu, Zhenfeng Zhu, Shuai Zheng, Yawei Zhao, Kunlun He, and Yao Zhao. From observation to concept: A flexible multi-view paradigm for medical report generation. *IEEE Transactions on Multimedia*, 26:5987–5995, 2024d. URL <https://api.semanticscholar.org/CorpusID:266303866>.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. PEFT: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022.
- Vedran Markotić, Tina Pojužina, Dorijan Radančević, Miro Miljko, and Vladimir Pokrajčić. The radiologist workload increase; where is the limit?: mini review and case study. *Psychiatria Danubina*, 33(suppl 4):768–770, 2021.
- Qiguang Miao, Kang Liu, Zhuoqi Ma, Yunan Li, Xiaolu Kang, Ruixuan Liu, Tianyi Liu, Kun Xie, and Zhicheng Jiao. Evoke: Elevating chest x-ray report generation via multi-view contrastive learning and patient-specific knowledge. *arXiv preprint arXiv:2411.10224*, 2024.
- Yasuhide Miura, Yuhao Zhang, C. Langlotz, and Dan Jurafsky. Improving factual completeness and consistency of image-to-text radiology report generation. *ArXiv*, abs/2010.10042, 2020. URL <https://api.semanticscholar.org/CorpusID:224803298>.
- Hoang T.N. Nguyen, Dong Nie, Taivanbat Badamdorj, Yujie Liu, Yingying Zhu, Jason Truong, and Li Cheng. Automated generation of accurate & fluent medical x-ray reports. In *Conference on Empirical Methods in Natural Language Processing*, 2021. URL <https://api.semanticscholar.org/CorpusID:237347122>.
- Aaron Nicolson, Jason Dowling, Douglas Anderson, and Bevan Koopman. Longitudinal data and a semantic similarity reward for chest x-ray report generation. *Informatics in Medicine Unlocked*, 50:101585, 2024.
- Farhad Nooralahzadeh, Nicolas Perez Gonzalez, Thomas Frauenfelder, Koji Fujimoto, and Michael Krauthammer. Progressive transformer-based generation of radiology reports. *arXiv preprint arXiv:2102.09777*, 2021.

-
- Mayank Pahadia, Sonam Khurana, Hassem Geha, and S Thomas Ii Deahl. Radiology report writing skills: A linguistic and technical guide for early-career oral and maxillofacial radiologists. *Imaging Science in Dentistry*, 50:269 – 272, 2020. URL <https://api.semanticscholar.org/CorpusID:222067570>.
- Jiazhen Pan, Che Liu, Junde Wu, Fenglin Liu, Jiayuan Zhu, Hongwei Bran Li, Chen Chen, Ouyang Cheng, and Daniel Rueckert. Medvlm-r1: Incentivizing medical reasoning capability of vision-language models (vlms) via reinforcement learning. *ArXiv*, abs/2502.19634, 2025. URL <https://api.semanticscholar.org/CorpusID:276647067>.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Annual Meeting of the Association for Computational Linguistics*, 2002. URL <https://api.semanticscholar.org/CorpusID:11080756>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- Vignav Ramesh, Nathan A Chi, and Pranav Rajpurkar. Improving radiology report generation systems by removing hallucinated references to non-existent priors. In *Machine Learning for Health*, pp. 456–473. PMLR, 2022.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla P. Kiraly, Madeleine Traverse, Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, Justin Chen, Fereshteh Mahvar, Liron Yatziv, Tiffany Chen, Bram Sterling, Stefanie Anna Baby, Susanna Maria Baby, Jeremy Lai, Samuel Schmidgall, Lu Yang, Kejia Chen, Per Bjornsson, Shashir Reddy, Ryan Brush, Kenneth Philbrick, Mercy Nyamewaa Asiedu, Ines Mezerreg, Howard Hu, Howard Yang, Richa Tiwari, Sunny Jansen, Preeti Singh, Yun Liu, Shekoofeh Azizi, Aishwarya B Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ram’e, Morgane Rivi re, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-Bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Elena Buchatskaya, Jean-Baptiste Alayrac, Dmitry Lepikhin, Vladimir Feinberg, Sebastian Borgeaud, Alek Andreev, Cassidy Hardin, Robert Dadashi, L eonard Hussenot, Armand Joulin, Olivier Bachem, Yossi Matias, Katherine Chou, Avinatan Hassidim, Kavi Goel, Cl ement Farabet, Joelle K. Barral, Tris Warkentin, Jonathon Shlens, David Fleet, Victor Cotruta, Omar Sanseviero, Gus Martins, Phoebe Kirk, Anand Rao, Shravya Shetty, David Steiner, Can Kirmizibayrak, Rory Pilgrim, Daniel Golden, and Lin Yang. Medgemma technical report. *ArXiv*, abs/2507.05201, 2025. URL <https://api.semanticscholar.org/CorpusID:280150648>.
- Francesco Dalla Serra, Chaoyang Wang, Fani Deligianni, Jeffrey Dalton, and Alison Q O’Neil. Controllable chest x-ray report generation from longitudinal representations. *arXiv preprint arXiv:2310.05881*, 2023.
- Hoo-Chang Shin, Kirk Roberts, Le Lu, Dina Demner-Fushman, Jianhua Yao, and Ronald M Summers. Learning to read chest x-rays: Recurrent neural cascade model for automated image annotation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2497–2506, 2016.
- Prateek Singh and Sudhakar Singh. Chestx-transcribe: a multimodal transformer for automated radiology report generation from chest x-rays. *Frontiers in Digital Health*, 7:1535168, 2025.
- Mehreen Sirshar, Muhammad Faheem Khalil Paracha, Muhammad Usman Akram, Norah Saleh Alghamdi, Syeda Zainab Yousuf Zaidi, and Tatheer Fatima. Attention based automated radiology report generation using cnn and lstm. *Plos one*, 17(1):e0262209, 2022.
- Phillip Sloan, Philip Clatworthy, Edwin Simpson, and Majid Mirmehdi. Automated radiology report generation: A review of recent advances. *IEEE Reviews in Biomedical Engineering*, 18:368–387, 2024.

-
- Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. Interactive and explainable region-guided radiology report generation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7433–7442, 2023. URL <https://api.semanticscholar.org/CorpusID:258179419>.
- Ryutaro Tanno, David GT Barrett, Andrew Sellergren, Sumedh Ghaisas, Sumanth Dathathri, Abigail See, Johannes Welbl, Charles Lau, Tao Tu, Shekoofeh Azizi, et al. Collaboration between clinicians and vision–language models in radiology report generation. *Nature Medicine*, 31(2): 599–608, 2025.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Lasa Team, Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Ruifeng Yuan, Jianyu Wang, Chenghao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, Yu Sun, Junao Shen, Chaojun Wang, Jie Tan, Deli Zhao, Tingyang Xu, Hao Zhang, and Yu Rong. Lingshu: A generalist foundation model for unified multimodal medical understanding and reasoning. *ArXiv*, abs/2506.07044, 2025. URL <https://api.semanticscholar.org/CorpusID:279250959>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Tony W Trinh, Atul B. Shinagare, Daniel I. Glazer, Pamela J. DiPiro, Jacob C. Mandell, Giles Boland, and Ramin Khorasani. Radiology report template optimization at an academic medical center. *AJR. American journal of roentgenology*, pp. 1 – 7, 2019. URL <https://api.semanticscholar.org/CorpusID:199662166>.
- Alison Wallis and Paul McCoubrie. The radiology report—are we getting the message across? *Clinical radiology*, 66 11:1015 – 1022, 2011. URL <https://api.semanticscholar.org/CorpusID:207014363>.
- Fuying Wang, Shenghui Du, and Lequan Yu. Hergen: Elevating radiology report generation with longitudinal data. In *European Conference on Computer Vision*, pp. 183–200. Springer, 2024.
- Jun Wang, Abhir Bhalerao, and Yulan He. Cross-modal prototype driven network for radiology report generation. In *European Conference on Computer Vision*, pp. 563–579. Springer, 2022a.
- Shansong Wang, Mingzhe Hu, Qiang Li, Mojtaba Safari, and Xiaofeng Yang. Capabilities of gpt-5 on multimodal medical reasoning. *ArXiv*, abs/2508.08224, 2025. URL <https://api.semanticscholar.org/CorpusID:280566563>.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, and Ronald M Summers. Tienet: Text-image embedding network for common thorax disease classification and reporting in chest x-rays. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9049–9058, 2018.
- Zhanyu Wang, Mingkan Tang, Lei Wang, Xiu Li, and Luping Zhou. A medical semantic-assisted transformer for radiographic report generation. *ArXiv*, abs/2208.10358, 2022b. URL <https://api.semanticscholar.org/CorpusID:251719362>.
- Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. R2gengpt: Radiology report generation with frozen llms. *ArXiv*, abs/2309.09812, 2023. URL <https://api.semanticscholar.org/CorpusID:262044687>.
- David L Weiss and Curtis P Langlotz. Structured reporting: patient care enhancement or productivity nightmare? *Radiology*, 249(3):739–747, 2008.

-
- Bin Yan and Mingtao Pei. Clinical-bert: Vision-language pre-training for radiograph diagnosis and reports generation. In *AAAI Conference on Artificial Intelligence*, 2022. URL <https://api.semanticscholar.org/CorpusID:250298613>.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Shuxin Yang, Xian Wu, Shen Ge, Xingwang Wu, S.kevin Zhou, and Li Xiao. Radiology report generation with a learned knowledge base and multi-modal alignment. *Medical image analysis*, 86: 102798, 2021. URL <https://api.semanticscholar.org/CorpusID:245634857>.
- Shuxin Yang, Xian Wu, Shen Ge, Zhuozhao Zheng, S. Kevin Zhou, and Li Xiao. Radiology report generation with a learned knowledge base and multi-modal alignment. *Medical Image Analysis*, 86:102798, 2023. ISSN 1361-8415. doi: <https://doi.org/10.1016/j.media.2023.102798>. URL <https://www.sciencedirect.com/science/article/pii/S1361841523000592>.
- Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. Automatic radiology report generation based on multi-view image fusion and medical concept enrichment. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2019. URL <https://api.semanticscholar.org/CorpusID:198148007>.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. *ArXiv*, abs/1904.09675, 2019. URL <https://api.semanticscholar.org/CorpusID:127986044>.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *ArXiv*, abs/2304.10592, 2023a. URL <https://api.semanticscholar.org/CorpusID:258291930>.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Yue Cao, Yangzhou Liu, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Han Lv, Dengnian Chen, Songze Li, Yinan He, Tan Jiang, Jiapeng Luo, Yi Wang, Cong He, Botian Shi, Xingcheng Zhang, Wenqi Shao, Junjun He, Ying Xiong, Wenwen Qu, Peng Sun, Penglong Jiao, Lijun Wu, Kai Zhang, Hui Deng, Jiaye Ge, Kaiming Chen, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *ArXiv*, abs/2504.10479, 2025. URL <https://api.semanticscholar.org/CorpusID:277780955>.
- Qingqing Zhu, Tejas Sudharshan Mathai, Pritam Mukherjee, Yifan Peng, Ronald M Summers, and Zhiyong Lu. Utilizing longitudinal chest x-rays and reports to pre-fill radiology reports. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 189–198. Springer, 2023b.

APPENDIX

This appendix provides comprehensive supplementary materials including limitations and future work discussions (§A), specific limitations of CheXagent-3B with long clinical contexts (§B), detailed dataset statistics with patient-level data splits (§C), prompt design examples for findings and impression generation (§D), complete instruction fine-tuning examples with clinical context integration (§E), and analysis of temporal hallucinations in radiology report generation (§F).

A LIMITATIONS AND FUTURE WORKS

A.1 SYNTHETIC DATASET AND ANNOTATIONS

Our C-SRRG dataset builds upon the SRRG dataset (Delbrouck et al., 2025), which was generated through reformulation of free-form radiology reports using large language models. This synthetic generation process introduces several potential limitations that warrant careful consideration. The use of LLM-generated content may introduce subtle hallucinations or inconsistencies that could propagate through our training pipeline. Additionally, the reformulation process may inadvertently introduce biases present in the underlying language models, potentially affecting the diversity and clinical accuracy of the generated reports.

A.2 MODEL ARCHITECTURE AND SCALE LIMITATIONS

Our experimental evaluation faces several architectural constraints that limit the full potential of our approach. First, we are constrained to backbone models with parameters up to 7B, which likely underestimates achievable performance. The computational and memory requirements of larger models present practical limitations for comprehensive evaluation across multiple architectures and clinical contexts.

Second, the multimodal large language models employed were not originally designed to handle multiple images simultaneously. In our experiments, we encountered specific limitations requiring tailored approaches. **Lingshu-7B** (Team et al., 2025) was limited to 2 images due to computational efficiency constraints, while **CheXagent-3B** (Chen et al., 2024b) was similarly restricted to 2 images due to model constraints. In contrast, **MedGemma-4B** (Sellersgren et al., 2025) demonstrated multi-image processing capabilities, requiring fewer tokens per image and enabling the use of more images in our longitudinal analysis setting.

These parameter and architectural constraints particularly impact the models’ ability to effectively integrate complex temporal relationships and multi-modal clinical information across current and prior studies. However, ongoing research trends in long-context LLM development (Dao et al., 2022; Kwon et al., 2023) suggest that future model architectures will naturally address these limitations. Scaling to larger foundation models, exploring mixture-of-experts variants, and advances in multimodal attention mechanisms represent promising directions for substantial improvements.

A.3 CLINICAL CONTEXT INTEGRATION AND SELECTION LIMITATIONS

Our approach faces constraints in both the scope and selection of clinical contexts. We impose limits on the number of images and prior studies included per clinical case, with prioritization given to the most recent studies. This limitation stems from current multimodal architectures’ context window constraints and computational overhead of processing extensive longitudinal histories. While this recency-based selection strategy captures the most clinically relevant temporal information, it may occasionally omit important historical context that could inform diagnostic reasoning. Additionally, our current implementation primarily utilizes publicly available datasets such as MIMIC-CXR (Johnson et al., 2019) and CheXpert-Plus (Chambon et al., 2024), limiting the diversity of clinical contexts.

The scope could be significantly expanded to include additional clinical modalities such as CT imaging, Electronic Health Records (EHR) (Häyrynen et al., 2008), and comprehensive patient histories. Future work should explore learned selection policies that intelligently identify the most informative clinical contexts and optimize longitudinal coverage. Retrieval-augmented generation approaches over Picture Archiving and Communication Systems (PACS) (Andriole, 2023) and EHR systems could dynamically surface the most relevant historical information for each case, unlocking the untapped potential for richer clinical context integration.

A.4 TRAINING METHODOLOGY AND DECODING LIMITATIONS

Our training approach is restricted to supervised fine-tuning with greedy decoding for reproducibility and computational efficiency. This methodology, while providing stable and consistent results, may not fully capture the nuanced decision-making processes that characterize expert radiological interpretation. The supervised learning paradigm limits the model’s ability to learn from comparative feedback and iterative refinement that occurs in clinical practice. Incorporating preference learning techniques and Reinforcement Learning (RL)-based methods with radiologists’ feedback, such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) or Direct Preference Optimization (DPO) (Rafailov et al., 2023), could enhance the fidelity and clinical appropriateness of generated reports. Furthermore, exploring retrieval-conditioned decoding strategies could improve temporal consistency and reduce hallucinations by grounding generation in verified clinical contexts.

B LIMITATION IN CHEXAGENT-3B ON C-SRRG-IMPRESSION

While most models show improvement with clinical context, CheXagent-3B exhibits a **critical failure** in following the structured report format instructions when provided with full clinical context. Instead of generating properly formatted impression sections with numbered findings, the model frequently produces single-word outputs or generic phrases. For instance, when the expected format is a multi-point structured impression such as “1. Slight decrease in size of the right apicolateral pneumothorax with chest tube in place. 2. Unchanged multifocal right-sided pulmonary opacities...”, CheXagent-3B often generates only “Pneumothorax” or “Pneumonia”. This format degradation is widespread, with the model generating non-structured outputs like “No acute cardiopulmonary process” or “Pulmonary edema” rather than detailed clinical impressions.

Table 11: **Performance degradation of CheXagent-3B on C-SRRG-Impression with full clinical context.** The model shows dramatic drops across all metrics when provided with complete clinical context.

Model	Full Clinical Context	Split	Traditional Metrics				F1-SRR-BERT		
			BLEU	ROUGE-L	BERT Score	F1-RadGraph	Precision	Recall	F1-Score
CheXagent-3B	✗	Valid	9.44	34.03	61.82	19.30	63.80	63.48	59.10
		Test	7.83	29.40	59.82	16.13	57.18	59.18	54.27
		Test-reviewed	7.42	28.60	58.35	13.71	51.32	56.34	49.74
	✓	Valid	2.57	21.76	40.10	13.10	74.48	49.40	54.05
		Test	2.40	17.54	33.79	9.78	66.56	41.04	45.99
		Test-reviewed	2.89	19.61	37.88	11.87	64.18	41.27	46.44

The performance metrics in Tab. 11 reveal the severity of this issue: when provided with full clinical context, traditional metrics plummet dramatically (BLEU: 9.44→2.57, ROUGE-L: 34.03→21.76, BERTScore: 61.82→40.10 on validation set). This catastrophic degradation suggests that CheXagent-3B, likely trained primarily on shorter sequence lengths, struggles to process and integrate the extensive clinical context while maintaining adherence to the structured output format. The model’s inability to handle long input sequences effectively undermines its utility for clinical applications requiring comprehensive context integration.

C DETAILED DATASET STATISTICS

We provide detailed statistics of our clinical context chest X-ray dataset, focusing on patient distribution across splits.

Patient Distribution Across Splits. Our dataset maintains strict patient-level separation across training, validation, and test splits to prevent data leakage. As shown in the patient overlap heatmaps, the training set contains 83,147 unique patients for findings and 125,947 unique patients for impression tasks. The validation sets include 434 patients for findings and 477 patients for impression, while the test sets contain 274 patients for findings and 423 patients for impression. The test-reviewed splits comprise 173 patients for findings and 172 patients for impression, with 106 and 108 patients respectively shared with the test split. This patient-level split ensures that clinical studies from the same patient do not appear across different evaluation splits, with zero patient overlap between training and evaluation sets. The distribution maintains clinical diversity while preserving the integrity of comprehensive clinical contexts within patient histories.

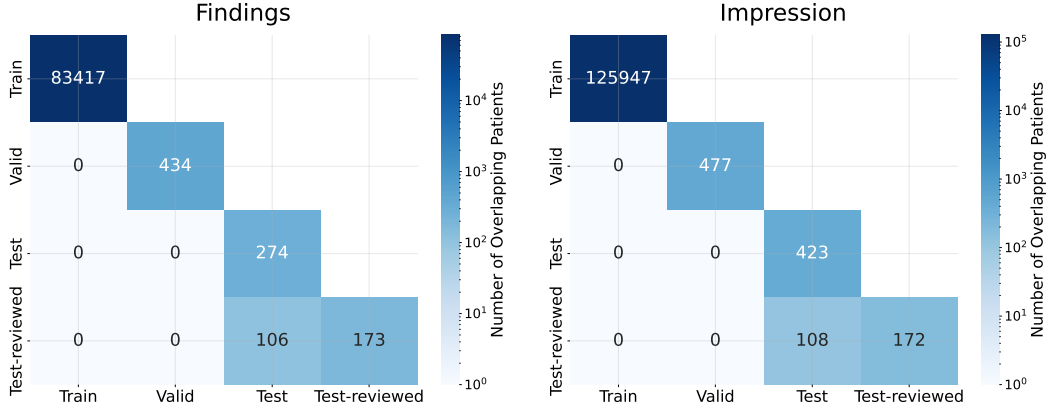


Figure 10: **Patient overlap heatmaps** across train, valid, test, and test-reviewed splits.

D PROMPT DESIGNS

In this section, beyond formats for impression with prior studies (Fig. 7), current study (Fig. 8), we provide examples of our design choices for prompts, *e.g.*, formats for findings with/without prior studies (Figs. 11 and 12), impression without prior studies (Fig. 13), prior studies (Fig. 14), and response format for findings and impression (Figs. 9 and 15), used for training and evaluation.

SYSTEM PROMPT:

You are an expert radiologist.

USER PROMPT:

Analyze the current chest X-ray images and compare them with the previous studies to write the FINDINGS section of a radiology report. Use standard medical terminology and note any changes from the prior studies, focusing on the most recent comparisons. Consider the available clinical contexts when formulating your findings.

Figure 11: **Prompt format for FINDINGS** for C-SRRG.

SYSTEM PROMPT:

You are an expert radiologist.

USER PROMPT:

Analyze the chest X-ray images and write the FINDINGS section of a radiology report. Use standard medical terminology and organize findings by anatomical regions. Consider the available clinical contexts when formulating your findings.

Figure 12: **Prompt format for FINDINGS without previous history** for C-SRRG.

SYSTEM PROMPT:

You are an expert radiologist.

USER PROMPT:

Analyze the chest X-ray images and write the IMPRESSION section of a radiology report. Provide a concise clinical summary and diagnosis based on the imaging findings. Consider the available clinical contexts when formulating your impression.

Figure 13: **Prompt format for IMPRESSION without previous history** for C-SRRG.

=== PREVIOUS STUDY N ===

Previous Clinical History/Indication:

Status post coronary artery bypass grafting (CABG), post chest tube removal.

Previous Technique:

Standard posteroanterior and lateral chest radiographs.

Previous Findings:

Lungs and Airways:

- No pneumothorax visualized
- Low lung volumes
- Increased infiltrate in the left upper lung

Pleura:

- Retrocardiac opacity
- Layering left pleural effusion

Cardiovascular:

- Vascular plethora

Other:

- Patchy areas of alveolar edema on the right
- Endotracheal (ET) tube positioned 4.3 cm above the carina
- Nasogastric (NG) tube within the stomach

Previous Impression:

1. Worsened infiltrate in the left upper lung suggestive of an acute process.
2. Presence of retrocardiac opacity and layering left pleural effusion.
3. Vascular plethora and patchy alveolar edema on the right, which may indicate fluid overload or cardiac etiology.
4. Proper positioning of the ET and NG tubes.

Previous Images:

[Images: Previous chest X-ray images]

Figure 14: **An example of prior study** for C-SRRG.

ASSISTANT RESPONSE:

FINDINGS:

Lungs and Airways:

- No pneumothorax or focal consolidation.

Pleura:

- Enlarged small right pleural effusion.

Cardiovascular:

- Mildly enlarged heart.

Hila and Mediastinum:

- Normal hilar and mediastinal contours.

Tubes, Catheters, and Support Devices:

- Right internal jugular (IJ) catheter terminates at the superior cavoatrial junction.

Figure 15: **An example of ground-truth assistant responses** in the C-SRRG-Findings dataset.

Findings Example - Part 1: Current Study Context

USER PROMPT:

Analyze the current chest X-ray images and compare them with the previous studies to write the FINDINGS section of a radiology report. Use standard medical terminology and note any changes from the prior studies, focusing on the most recent comparisons. Consider the available clinical contexts when formulating your findings.

=== CURRENT CLINICAL HISTORY/INDICATION ===

Evaluation for fluid overload.

=== CURRENT TECHNIQUE ===

Standard frontal chest radiography protocol.

=== CURRENT COMPARISON ===

Prior radiographs and CT scans.

=== CURRENT IMAGES ===

[Images: Current chest X-ray images]

Figure 16: **Findings generation example (Part 1)** in the C-SRRG-Findings dataset.

E INSTRUCTION TUNING DATASET PROMPT EXAMPLE

We provide detailed instruction fine-tuning examples that showcase the comprehensive clinical context utilized in our approach. These examples demonstrate how all available clinical information is systematically integrated into our instruction tuning dataset, including patient medical history, imaging techniques, previous study findings, and temporal comparisons. The following multi-part examples illustrate the complete structure of our training data, highlighting how comprehensive clinical contexts including temporal, multi-view, and metadata information are preserved and leveraged for clinical reasoning in radiology report generation.

E.1 FINDINGS GENERATION EXAMPLE

The first example demonstrates the generation of the FINDINGS section, which requires detailed anatomical observation and temporal comparison across multiple studies (Figs. 16 to 18):

E.2 IMPRESSION GENERATION EXAMPLE

The second example demonstrates the generation of the IMPRESSION section, which requires clinical synthesis and diagnostic reasoning (Figs. 19 to 21):

F HALLUCINATION ANALYSIS

In this section, we examine a critical limitation of radiology report generation models trained without clinical context, specifically their tendency to hallucinate temporal comparisons when referencing non-existent prior studies. We first demonstrate that dataset ground truth reports contain temporal statements that become hallucinations when clinical context is absent, as radiologists naturally write these comparisons when they have access to prior studies. We then analyze how models trained without such clinical context systematically produce these hallucinations, even for patients with no imaging history. Finally, we quantify these hallucinations by detecting the frequency of temporal statements on the generated reports on the evaluation set without clinical context.

Dataset Hallucination. Ground truth radiology reports in clinical datasets frequently contain temporal statements such as “new from prior exam,” “unchanged,” or “stable compared to previous study.” These temporal references are clinically appropriate when radiologists have access to prior imaging studies for comparison. However, when language models are trained on these reports without access to the corresponding clinical context and prior studies, they learn to replicate these

Findings Example - Part 2: Previous Study 1

=== PREVIOUS STUDY 1 (Most Recent) ===

Previous Clinical History/Indication:

Patient with a history of multifocal after CABG, currently presenting with symptoms suggestive of CHF or pneumonia.

Previous Technique:

A single frontal chest radiograph was obtained.

Previous Comparison:

Multiple prior radiographs

Previous Findings:

Lungs and Airways:

- No definitive consolidation observed on this examination; however, subsequent CT confirms presence at the right base
- Mild pulmonary edema

Pleura:

- Moderate right pleural effusion, unchanged
- No pneumothorax

Cardiovascular:

- Moderate cardiomegaly noted
- Aortic tortuosity present

Tubes, Catheters, and Support Devices:

- Status post median sternotomy with CABG and valve replacements

Previous Impression:

1. Mild pulmonary edema with right pleural effusion.
2. Right base consolidation not definitive on radiograph but confirmed on subsequent CT.

Previous Images:

[Images: Previous chest X-ray images]

Figure 17: **Findings generation example (Part 2)** in the C-SRRG-Findings dataset.

temporal language patterns indiscriminately. This training paradigm creates a systematic hallucination problem where models generate temporal comparison statements even for patients with no prior imaging history. The following examples demonstrate these temporal hallucinations present in ground truth radiology reports from the dataset, showing how temporal comparison statements appear without proper clinical context (Figs. 22 to 24):

Trained Model Hallucination. Our trained models also exhibit hallucination patterns, particularly when generating temporal comparisons despite having access only to the current study. The following examples from MedGemma-4B model trained without clinical history demonstrate how models can inappropriately generate comparative language even when no prior studies are available for comparison (Figs. 25 to 30):

Findings Example - Part 3: Previous Study 2 & Response

=== PREVIOUS STUDY 2 ===

Previous Clinical History/Indication:

The patient presents with tachycardia.

Previous Technique:

Posteroanterior (PA) and lateral chest radiographs.

Previous Findings:

Cardiovascular:

- Status post sternotomy with mitral and aortic valve replacements
- Moderate cardiomegaly, unchanged

Hila and Mediastinum:

- Stable mediastinal and hilar contours

Lungs and Airways:

- New mild-to-moderate elevation of the right hemidiaphragm suggesting volume loss
- Increased patchy opacity in the right lower lobe, concerning for pneumonic consolidation
- Probable new opacification in the right middle lobe
- New vague opacity in the right suprahilar region of the right upper lobe, potentially indicating early pneumonia

Pleura:

- No definite pleural effusion

Previous Impression:

1. Findings suggest multifocal pneumonia involving the right lung.

ASSISTANT RESPONSE:

FINDINGS:

Lungs and Airways:

- No pneumothorax or focal consolidation.

Pleura:

- Enlarged small right pleural effusion.

Cardiovascular:

- Mildly enlarged heart.

Hila and Mediastinum:

- Normal hilar and mediastinal contours.

Tubes, Catheters, and Support Devices:

- Right internal jugular (IJ) catheter terminates at the superior cavoatrial junction.

Figure 18: **Findings generation example (Part 3)** in the C-SRRG-Findings dataset.

Example Case with Clinical Context - Part 1: Current Study

USER PROMPT:

Analyze the current chest X-ray images and compare them with the previous studies to write the IMPRESSION section of a radiology report. Provide a concise clinical summary and diagnosis, noting any changes from the prior studies, focusing on the most recent comparisons. Consider the available clinical contexts when formulating your impression.

=== CURRENT CLINICAL HISTORY/INDICATION ===

Male with end-stage renal disease on hemodialysis, multiple orthopedic hardware, bioprosthetic aortic valve replacement, coronary artery disease status post coronary artery bypass grafting, permanent pacemaker for sick sinus syndrome, admitted for MRSA bacteremia, transferred to the critical care unit for hypotension during anesthesia induction. Patient intubated for respiratory status assessment.

=== CURRENT TECHNIQUE ===

Chest single view

=== CURRENT COMPARISON ===

Prior imaging at an unspecified time.

=== CURRENT IMAGES ===

[Images: Current chest X-ray images]

Figure 19: **Impression generation example (Part 1)** in the C-SRRG-Impression dataset.

Example Case with Clinical Context - Part 2: Previous Study 1

=== PREVIOUS STUDY 1 (Most Recent) ===

Previous Clinical History/Indication:

Status post coronary artery bypass grafting (CABG), post chest tube removal.

Previous Technique:

Standard posteroanterior and lateral chest radiographs.

Previous Findings:

Lungs and Airways:

- No pneumothorax visualized
- Low lung volumes
- Increased infiltrate in the left upper lung

Pleura:

- Retrocardiac opacity
- Layering left pleural effusion

Cardiovascular:

- Vascular plethora

Other:

- Patchy areas of alveolar edema on the right
- Endotracheal (ET) tube positioned 4.3 cm above the carina
- Nasogastric (NG) tube within the stomach

Previous Impression:

1. Worsened infiltrate in the left upper lung suggestive of an acute process.
2. Presence of retrocardiac opacity and layering left pleural effusion.
3. Vascular plethora and patchy alveolar edema on the right, which may indicate fluid overload or cardiac etiology.
4. Proper positioning of the ET and NG tubes.

Previous Images:

[Images: Previous chest X-ray images]

Figure 20: **Impression generation example (Part 2)** in the C-SRRG-Impression dataset.

Example Case with Clinical Context - Part 3: Previous Study 2 & Response

=== PREVIOUS STUDY 2 ===

Previous Clinical History/Indication:

Patient presents with portions and exertion x.

Previous Technique:

The examination was performed using frontal and lateral views of the chest.

Previous Findings:

Tubes, Catheters, and Support Devices:

- Dual lead left-sided pacemaker extending to the expected positions of the right atrium and right ventricle.

Pleura:

- Slight blunting of the posterior costophrenic angles, which may suggest very trace pleural effusions.

Lungs and Airways:

- No focal consolidation observed.

Cardiovascular:

- The cardiac silhouette appears stable.

Hila and Mediastinum:

- The mediastinal silhouette is stable.

Musculoskeletal and Chest Wall:

- Right proximal humerus hardware is present but not well evaluated due to the imaging technique.

Other:

- Slight prominence of the interstitium, which may indicate minimal interstitial edema.

Previous Impression:

1. Possible very trace pleural effusions.
2. Minimal interstitial edema.
3. Stable cardiac and mediastinal silhouettes.
4. Presence of dual lead left-sided pacemaker in expected positions.
5. Right proximal humerus hardware present but not adequately assessed.

ASSISTANT RESPONSE:

IMPRESSION:

1. Worsened left basilar opacity, which may represent atelectasis or pneumonitis in the appropriate clinical context.
2. Mild improvement in pulmonary vascularity.

Figure 21: **Impression generation example (Part 3)** in the C-SRRG-Impression dataset.

Dataset Hallucination Example 1: Temporal Information Fabrication

Structured Report:

Exam Type: Chest radiograph.

Technique: Portable anteroposterior (AP) chest radiography was performed.

History: A male patient with hep C cirrhosis and large right pleural effusion status post thoracocentesis. Evaluate for resolution of pleural effusion.

Comparison: Prior portable AP chest radiograph

Findings:

Lungs and Airways:

- Mild inflation of the right upper lobe
- Collapsed right lower lobe
- No consolidation in the left lung

Pleura:

- Moderate pleural effusion within the right pleural space.
- Moderate right pneumothorax, new from prior exam.
- No left pleural effusion or pneumothorax.

Cardiovascular:

- No significant mediastinal shift observed.

Hila and Mediastinum:

- Mediastinum appears unremarkable

Impression:

1. Moderate right-sided pneumothorax.
2. Moderate right pleural effusion.
3. Inflation of the right upper lobe with collapse of the right lower lobe.
4. No mediastinal shift.

Hallucination: The phrase “new from prior exam” represents temporal information that cannot be verified from the current study alone, if not with previous history.

Figure 22: Dataset hallucination example 1 in SRRG dataset.

Dataset Hallucination Example 2: Stability Assumption Without Comparison

Structured Report:

Exam Type: Chest radiograph

Technique: Standard frontal and lateral chest radiographic views were performed.

History: Atrial fibrillation (AF), coronary artery disease (CAD), congestive heart failure (CHF).

Comparison: Prior chest radiographs

Findings:

Cardiovascular:

- Mild to moderate cardiomegaly, unchanged.
- Tortuous but stable aorta.

Tubes, Catheters, and Support Devices:

- Transvenous pacemaker/AICD with leads terminating in the right atrium and right ventricle.
- Median sternotomy wires are aligned and intact.

Lungs and Airways:

- Lungs are clear with no evidence of consolidation, pleural effusion, pneumothorax, or overt pulmonary edema.

Impression:

1. No radiographic evidence for acute cardiopulmonary process.

Hallucination: The term “**unchanged**” implies comparison with prior studies, which is problematic without previous history.

Figure 23: **Dataset hallucination example 2** in SRRG dataset.

Dataset Hallucination Example 3: Generic Temporal Statement

Structured Report:

Exam Type: Chest radiograph

Technique: Single AP upright portable chest radiograph.

History: Shortness of breath.

Comparison: Prior chest radiograph

Findings:

Lungs and Airways:

- Low lung volumes
- Minimal bibasilar atelectasis
- Subcentimeter left lower lung rounded calcification, stable, possibly representing a calcified granuloma
- No focal consolidation
- No overt pulmonary edema

Pleura:

- No pleural effusion
- No evidence of pneumothorax

Cardiovascular:

- Stable cardiac and mediastinal silhouettes

Hila and Mediastinum:

- Ovoid calcification over the left mediastinum, stable

Tubes, Catheters, and Support Devices:

- Right-sided Port-A-Cath in place, terminating at the cavoatrial junction, without significant change in position

Musculoskeletal and Chest Wall:

- Status post median sternotomy

Impression:

1. No significant interval change.

Hallucination: The phrase “No significant interval change” is a generic temporal statement that, if not with prior studies, represents a non-specific conclusion that could mask important clinical details.

Figure 24: Dataset hallucination example 3 in SRRG dataset.

Trained Model Hallucination Example 1: Temporal Comparisons in Findings

Sample ID: 12

Generated Report:

FINDINGS:

Tubes, Catheters, and Support Devices:

- Endotracheal tube terminates 4.5 cm above the carina.
- Right internal jugular central venous catheter terminates in the right atrium.
- Left pectoral pacemaker with leads in the right atrium and right ventricle.
- Median sternotomy wires are intact.

Cardiovascular:

- Moderate cardiomegaly, unchanged.

Lungs and Airways:

- Low lung volumes.
- No evidence of pneumothorax.

Pleura:

- Small left pleural effusion, unchanged.

Other:

- Retrocardiac opacity, likely representing atelectasis, unchanged.

Hallucination: Multiple instances of “unchanged” indicate temporal comparisons when only analyzing a single study.

Figure 25: Trained model hallucination example 1 on test-reviewed split using MedGemma-4B.

Trained Model Hallucination Example 2: Interval and Persistent Language

Sample ID: 60

Generated Report:

FINDINGS:

Pleura:

- Interval placement of a right-sided pigtail catheter.
- Decreased size of the right pleural effusion.
- No evidence of pneumothorax.

Lungs and Airways:

- Persistent right basilar opacity, likely representing atelectasis.
- Left lung is clear.

Cardiovascular:

- Stable cardiomeastinal silhouette.

Hallucination: Terms like “Interval placement”, “Decreased size”, and “Persistent” all imply temporal comparisons.

Figure 26: Trained model hallucination example 2 on test-reviewed split using MedGemma-4B.

Trained Model Hallucination Example 3: Explicit Temporal Language

Sample ID: 37

Generated Report:

FINDINGS:

Lungs and Airways:

- No significant interval change in the appearance of the lungs.

Pleura:

- Small right pleural effusion, unchanged.

Cardiovascular:

- Stable cardiomeastinal silhouette.

Tubes, Catheters, and Support Devices:

- Tracheostomy tube in place.
- Right internal jugular (IJ) dialysis catheter in place.
- Vascular stent in the right subclavian region.

Other:

- No evidence of pneumothorax.

Hallucination: Explicit temporal language (“No significant interval change”, “unchanged”) when no prior studies are available.

Figure 27: Trained model hallucination example 3 on test-reviewed split using MedGemma-4B.

Trained Model Hallucination Example 4: Unchanged Position in Impression

Sample ID: 0

Generated Report:

IMPRESSION:

1. No evidence of pneumothorax.
2. Stable small bilateral pleural effusions.
3. Unchanged position of the right chest tube.
4. Normal cardiomeastinal silhouette.
5. Presence of a vagal nerve stimulator.

Hallucination: The term “Unchanged position” suggests comparison with prior imaging that is not available.

Figure 28: Trained model hallucination example 4 on test-reviewed split using MedGemma-4B.

Trained Model Hallucination Example 5: Interval and Persistent Terms

Sample ID: 36

Generated Report:

IMPRESSION:

1. Interval increase in bilateral pleural effusions.
2. Persistent left lower lobe atelectasis.
3. Mild pulmonary edema.
4. Stable cardiomegaly.

Hallucination: Terms “Interval increase” and “Persistent” indicate temporal comparisons without access to prior studies.

Figure 29: Trained model hallucination example 5 on test-reviewed split using MedGemma-4B.

Trained Model Hallucination Example 6: New vs Unchanged Terminology

Sample ID: 40

Generated Report:

IMPRESSION:

1. New right lower lobe opacity, which may represent pneumonia or aspiration.
2. Small right pleural effusion.
3. Stable mild cardiomegaly.
4. Unchanged position of the right-sided central venous catheter.

Hallucination: The term “Unchanged position” implies knowledge of prior imaging studies that the model does not have access to. **New** also suggests temporal comparison but may be acceptable if referring to acute findings.

Figure 30: Trained model hallucination example 6 on test-reviewed split using MedGemma-4B.

G THE USE OF LLMS

We used LLMS solely for light editing such as correcting grammatical errors and polishing some words. They did not contribute to research ideation, experiments, analysis, or substantive writing.