

Structuring Radiology Reports: Challenging LLMs with Lightweight Models

Anonymous ACL submission

Abstract

Radiology reports are critical for clinical decision-making but often lack a standardized format, limiting both human interpretability and machine learning (ML) applications. While large language models (LLMs) have shown strong capabilities in reformatting clinical text, their high computational requirements, lack of transparency, and data privacy concerns hinder practical deployment. To address these challenges, we explore lightweight encoder-decoder models (<300M parameters)—specifically T5 and BERT2BERT—for structuring radiology reports from the MIMIC-CXR and CheXpert Plus datasets. We benchmark these models against eight open-source LLMs (1B–70B parameters), adapted using prefix prompting, in-context learning (ICL), and low-rank adaptation (LoRA) finetuning. Our best-performing lightweight model outperforms all LLMs adapted using prompt-based techniques on a human-annotated test set. While some LoRA-finetuned LLMs achieve modest gains over the lightweight model on the Findings section (BLEU 6.4%, ROUGE-L 4.8%, BERTScore 3.6%, F1-RadGraph 1.1%, GREEN 3.6%, and F1-SRR-BERT 4.3%), these improvements come at the cost of substantially greater computational resources. For example, LLaMA-3-70B incurred more than 400 times the inference time, cost, and carbon emissions compared to the lightweight model. These results underscore the potential of lightweight, task-specific models as sustainable and privacy-preserving solutions for structuring clinical text in resource-constrained healthcare settings.

1 Introduction

Radiology reports play a critical role in clinical workflows by summarizing imaging findings that guide medical decisions (Kahn Jr et al., 2009). However, variations in reporting style due to individual and institutional practices as well as regional guidelines create inconsistencies that hinder

interpretability for physicians and patients (Hartung et al., 2020). Moreover, the lack of structured formats limits their usefulness as training data for machine learning (ML) applications (dos Santos et al., 2023; Steinkamp et al., 2019).

Large language models (LLMs) offer a promising solution for generating structured reports from free-form text (Adams et al., 2023; Busch et al., 2024; Hasani et al., 2024). However, deploying these models locally remains infeasible for most institutions due to the significant computational resources required (Zhang et al., 2025). Cloud-based solutions provide an alternative but introduce concerns related to data security, confidentiality, and regulatory compliance (Arshad et al., 2023; Thirunavukarasu et al., 2023). While proprietary LLMs can also be accessed via Application Programming Interface (API), this approach entails drawbacks such as dependency on a third-party vendor, potential cost increases and unpredictable changes in usage terms (Tian et al., 2024). These limitations highlight the need for smaller, open-source models that can be deployed on-device with minimal hardware requirements.

To address these challenges, we propose lightweight (<300M parameters), task-specific models for structuring free-text chest X-ray radiology reports (see Figure 1) efficiently. These models substantially reduce computational demands (Chen et al., 2024a), eliminating the need for cloud-based hosting, and enhancing data security by enabling offline deployment. We train these models on the MIMIC-CXR (Johnson et al., 2019) and CheXpert Plus (Chambon et al., 2024) datasets and structure the originally free-form reports with GPT-4 (Achiam et al., 2023) as a weak annotator, enabling large-scale supervision. We evaluate model performance on an independent test set, annotated by five radiologists (Anonymous, 2025). Our contributions include:

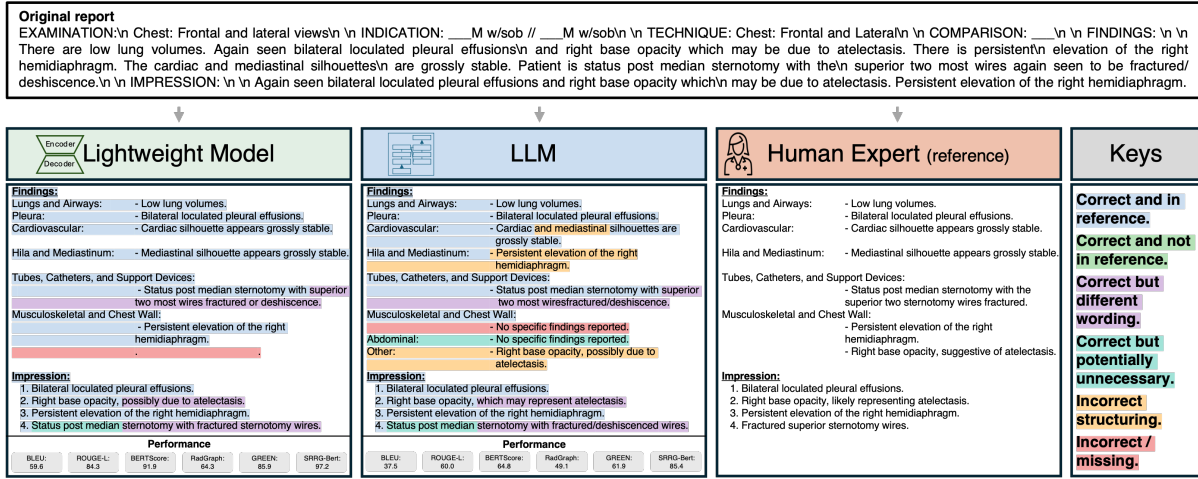


Figure 1: Overview of our study and qualitative comparison. An unstructured radiology report is structured using lightweight, task-specific models and adapted large language models (LLMs) compared to human expert annotations.

- **Lightweight Model Development and Evaluation:** We train and systematically evaluate lightweight (<300M parameters), task-specific T5 and BERT2BERT models for the task of structuring radiology reports.
- **Analysis of LLMs and Adaptation Techniques:** We assess the performance of five LLMs (3-8B parameters) under different adaptation strategies (prefix prompting, in-context learning (ICL), low-rank adaptation (LoRA)).
- **Benchmarking and Cost Analysis:** We benchmark lightweight models against LLMs of increasing size, considering model performance on the BLEU, ROUGE-L, BERTScore, F1-RadGraph, GREEN, and F1-SRRG-Bert metrics, as well as training time, inference speed and costs, and environmental impact.

2 Related Work

Beyond LLMs: Lightweight Models for Medical Text Processing

Recent studies have explored the use of LLMs, namely GPT-3.5 (OpenAI, 2022) and GPT-4, to transform free-form radiology reports into structured formats (Adams et al., 2023; Bergomi et al., 2024; Hasani et al., 2024). A recent review by Busch et al. highlights that these approaches achieve low error rates and minimal accuracy loss compared to human experts (Busch et al., 2024). However, their reliance on proprietary architectures, lack of transparency, and restrictions on patient data privacy pose significant challenges for clinical deployment (Khullar et al., 2024;

Rezaeikhonakdar, 2023). To address these limitations, similar tasks in medical NLP have adopted lightweight, task-specific models that maintain high accuracy while considerably reducing computational costs (Chen et al., 2024a; Griewing et al., 2024; Pecher et al., 2024). Existing task-specific models for radiology NLP fall into two categories: hybrid models and lightweight transformer models. Hybrid models combine rule-based methods with deep learning, enforcing domain-specific constraints but lacking flexibility (Gabud et al., 2023). In contrast, lightweight transformer models have been successfully applied to relation extraction, report coding, and summarization (Jain et al., 2021; Yan et al., 2022; Van Veen et al., 2023). While they require careful tuning to avoid hallucinations and overfitting, recent studies suggest that well-tuned lightweight models can match larger LLMs in accuracy while being far more computationally efficient (Pecher et al., 2024). Our work builds on this foundation by introducing a lightweight, task-specific model explicitly optimized for structured radiology report generation.

Model Adaptation and Finetuning

Prior work has explored a range of adaptation strategies for LLMs, from prompt-based methods to parameter-efficient finetuning (PEFT) and full finetuning, each balancing performance, data requirements, and computational cost. Prompting techniques such as prefix prompting and ICL (Brown et al., 2020; Lampinen et al., 2022) adapt models without modifying their weights. Prefix prompting typically provides instructions to guide model responses, while ICL enhances adaptation by incor-

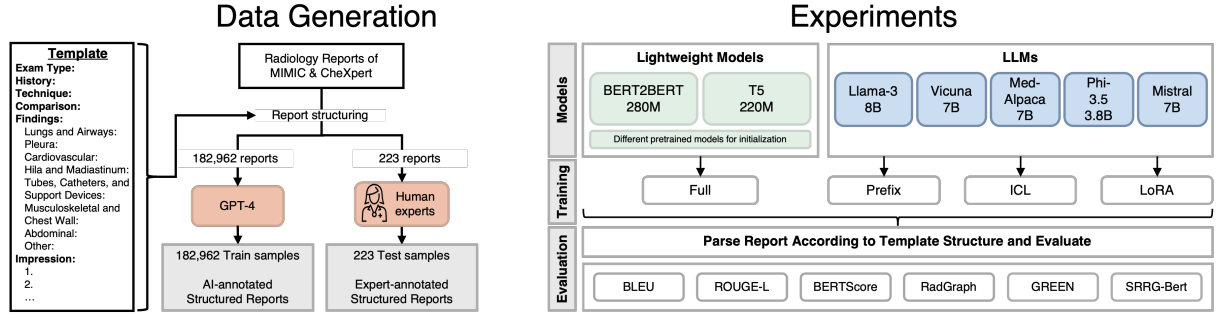


Figure 2: Left: Dataset generation from free-form radiology reports to structured radiology reports using GPT-4 (AI-based) and human experts (manual annotation). Right: Overview of our experiments including selection of lightweight models and LLMs, training/adaptation methods, and evaluation strategy and metrics.

porating task-specific examples within the prompt. However, these methods suffer from context length constraints and sensitivity to prompt phrasing (Li et al., 2023). PEFT techniques like LoRA (Hu et al., 2021), prefix-tuning (Li and Liang, 2021), and adapter layers (Houlsby et al., 2019) enable efficient adaptation with minimal computational overhead, making them well-suited for clinical NLP. While effective in low-data settings, PEFT often struggles with complex reasoning and generalization across domains (Lialin et al., 2023). In contrast, full finetuning updates all model parameters, often achieving stronger adaptation when sufficient labeled data and computational resources are available. Building on this, our approach applies full finetuning to lightweight models while leveraging GPT-4-generated structured labels to address data scarcity, enabling large-scale supervised training while preserving domain-specific accuracy.

AI-Based Dataset Generation

A major challenge in developing models for structuring radiology reports is the limited availability of high-quality annotated datasets, i.e., datasets that contain both free-form and corresponding structured reports. Recent work in similar fields has explored leveraging LLMs such as GPT-4 as weak annotators to generate labels, providing a scalable alternative to manual annotation (Liyanage et al., 2024; Savelka et al., 2023). Despite their successes, studies suggest that models trained on GPT-generated data should still be rigorously evaluated against human-annotated ground truth to ensure reliability and validity (Pangakis et al., 2023).

3 Methods

In this study, we transform free-text chest X-ray radiology reports into a standardized format using

deep learning. The structured reports follow a pre-defined template based on 'RPT144' of RSNA's RadReport Template Library (Radiological Society of North America (RSNA), 2011). This template comprises the sections: Exam Type, History, Technique, Comparison, Findings, and Impression. The Findings section is further organized into organ systems: 'Lungs and Airways', 'Pleura', 'Cardiovascular', 'Tubes, Catheters, and Support Devices', 'Musculoskeletal and Chest Wall', 'Abdominal', and 'Other'. The Impression section is structured as a numbered list, prioritizing the most clinically relevant findings. As shown in Figure 2, this template is incorporated into the prompt during data annotation, and deviations from it in a structured report are penalized during evaluation. Unlike previous approaches that rely on large, general-purpose models like GPT-4, we explore the effectiveness of lightweight, task-specific models for this task.

3.1 Data

We use unstructured radiology reports from the publicly available MIMIC-CXR (Johnson et al., 2019) and CheXpert Plus (Chambon et al., 2024) datasets, preserving their original training and validation splits. To train our models in a supervised manner, we employed GPT-4 as a weak annotator, using the prompt provided in Appendix A.1 to generate structured reports that conform to our template. We obtained a total of 182,962 reports, 125,447 samples from MIMIC-CXR and 57,515 from CheXpert Plus. For evaluation and benchmarking, we conducted a human expert review of 223 reports, comprising 161 from the MIMIC-CXR test set and 72 from the CheXpert Plus validation set. Five board-certified radiologists from our institution reviewed the structured reports alongside their original free-form counterparts, assessing them for errors and

adherence to our predefined template (detailed in (Anonymous, 2025)).

3.2 Evaluation Strategies

Even though all models generate full reports, we focus our quantitative analysis on the Findings and Impression sections due to their clinical significance. Before applying our metrics, we parse these sections to assess adherence to the predefined template. In the Findings section, we identify predefined organ system headers (e.g., 'Lungs and Airways', 'Cardiovascular') and extract their corresponding observations. Metrics are computed separately for each organ system and then averaged across all identified systems. In the Impression section, we enforce a sequentially numbered format and flag any inconsistencies in ordering. To assess both linguistic quality and clinical accuracy, we use a combination of lexical and radiology-specific metrics.

Lexical Metrics To ensure comprehensive evaluation of text quality, we apply the following metrics: *BLEU* (Papineni et al., 2002) measures n-gram overlap, serving as a proxy for fluency and syntactic similarity. *ROUGE-L* (Lin, 2004) evaluates the longest common subsequence, capturing sentence-level similarity. *BERTScore* (Zhang et al., 2019) computes semantic similarity by comparing contextual embeddings from a pretrained transformer model.

Radiology-Specific Metrics To capture clinical accuracy, we apply the following metrics: *F1-RadGraph* (Delbrouck et al., 2022; Yu et al., 2023) evaluates the precision and recall of key clinical terms and relationships extracted from generated reports. *GREEN* (Ostmeier et al., 2024) assesses the factual correctness of generated radiology reports using a finetuned LLM. *F1-SRRG-Bert* (Anonymous, 2025) uses a fine-tuned BERT model to classify extracted findings into 55 disease labels, assigning each as Present, Absent, or Uncertain. It then computes the F1-score by comparing predictions from the generated report to the ground truth. Throughout this paper, our visualizations primarily focus on GREEN and F1-SRR-BERT, as GREEN correlates most strongly with expert evaluations of clinical accuracy (Ostmeier et al., 2024), while F1-SRR-BERT was specifically developed for the task of structured reporting, making their combination effective for assessing structured radiology reports.

3.3 Lightweight Models

We introduce lightweight models, which are specifically trained to structure radiology reports according to a predefined template. Our lightweight models are based on encoder-decoder architectures given their recent success in similar tasks such as radiology report generation (Aksoy et al., 2023; Chen et al., 2024b) and radiology report summarization (de Padua and Qureshi, 2024; Van Veen et al., 2023; Zhang et al., 2018). Specifically, we focused on two architectures, *T5-Base* (Raffel et al., 2020), which has 223M parameters, and *BERT2BERT* (Rothe et al., 2020), where two identical BERT models are used as the encoder and decoder, resulting in a total of 278M parameters. To investigate the influence of pretraining domains, we initialize our models with the parameters from five open-source T5 variants (Table 2) - *T5-Base* (Raffel et al., 2020)(general text), *Flan-T5-Base* (Chung et al., 2024)(instruction-tuning), *SciFive* (Phan et al., 2021)(biomedical text), *Clin-T5-Sci* (Lehman and Johnson, 2023)(biomedical text and radiology reports), and *Clin-T5-Base* (Lehman and Johnson, 2023)(radiology reports) - and four BERT variants (Table 3) - *RoBERTa-base* (Liu, 2019)(general text), *BioMed-RoBERTa* (Gururangan et al., 2020)(biomedical text), *RoBERTa-base-PM-M3-Voc-distill-align* (Lewis et al., 2020)(for simplicity named RoBERTa-PM-M3 here, biomedical text and radiology reports), and *RadBERT-RoBERTa* (Yan et al., 2022)(radiology reports). We train our lightweight models end-to-end, updating all parameters, for a maximum of ten epochs using a cosine learning rate scheduler with an initial learning rate of $1e^{-4}$, an effective batch size of 128, and the Adam optimizer. A detailed description of hyperparameters can be found in Appendix A.3. To account for variability, each configuration is trained three times with different random seeds. Following prior work (Van Veen et al., 2023), we rank pretraining datasets by relevance, assuming radiology reports to be the most relevant, followed by biomedical text (e.g., PubMed abstracts) and general-domain text (e.g., Wikipedia). However, we acknowledge that this ranking is inherently subjective and may vary depending on the specific task.

3.4 Comparison LLMs

To benchmark our lightweight models (<300M parameters), we first conduct a comprehensive comparison with instruction-tuned LLMs ranging from

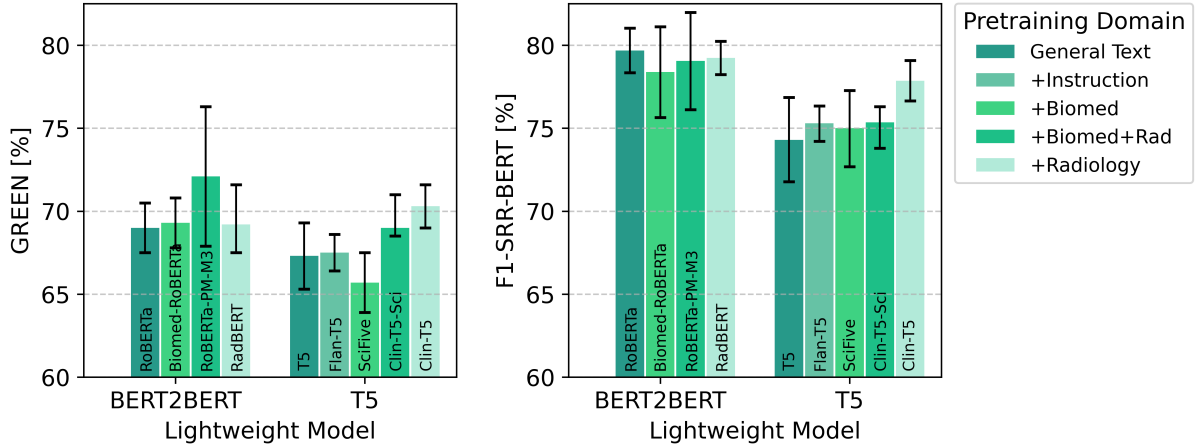


Figure 3: Performance comparison of lightweight models, initialized from pretrained models of increasing domain relevance. The plot shows the finetuned BERT2BERT and T5 models evaluated using GREEN (left) and F1-SRR-BERT (right), initialized from various pretrained models, with pretraining datasets ranging from general text (least domain-specific) to radiology (most domain-specific). Error bars denote 95% confidence intervals over the three training runs.

3 to 8 billion parameters: Llama-3.1-8B-Instruct (Grattafiori et al., 2024); its derivatives Vicuna-7B-v1.5 (Chiang et al., 2023), optimized for conversational tasks, and Med-Alpaca-7B (Han et al., 2023), finetuned for medical question-answering; as well as Phi-3.5-Mini-Instruct (Abdin et al., 2024) and Mistral-7B (Jiang et al., 2023). We assess three adaptation techniques: **1. Prefix Prompting.** The model is prompted using the same instructions employed during training data generation (Appendix A.1). **2. ICL.** The model is given a number of free-form reports along with their structured counterparts. These examples are manually selected from the training set to optimally represent the data distribution. **3. LoRA Finetuning.** The LLM is finetuned for five epochs on the complete training set using LoRA with a rank of eight, modifying approximately 0.1% of the model’s parameters by injecting trainable adapters into the key, query, and value projection matrices of the self-attention layers. We use a cosine learning rate scheduler with an initial learning rate of $1e^{-4}$, an effective batch size of 256 and the Adam optimizer. Detailed finetuning configurations are provided in Appendix A.4. Throughout the project, we systematically evaluated different combinations of these adaptation techniques. This included varying the number of in-context examples (1-shot, 2-shot) as well as combining *Prefix Prompting* with ICL to assess their complementary effects. We also experimented with hybrid approaches that combined LoRA finetuning with prompting-based methods.

However, these configurations did not yield consistent performance gains and introduced substantial overhead in terms of training time and memory usage, primarily due to increased input lengths.

3.5 Benchmarking Lightweight Models Against LLMs

Building on the previous experiment—which compared similarly sized LLMs under various adaptation strategies—we now turn to a scale-sensitive evaluation of our lightweight model. To this end, we benchmark its performance against LLaMA-3 models of increasing size (1B, 3B, 8B, and 70B parameters), leveraging the architectural consistency across this family to isolate the effects of model scale. Each variant is evaluated using the two most effective adaptation strategies identified in our prior experiments: *Prefix+ICL* for prompting-based approaches and *LoRA* for parameter-efficient finetuning. We then compare the computational costs associated with training and deploying the lightweight model, LLaMA-3-3B, and LLaMA-3-70B. This comparison includes the average F1-SRR-BERT score, training time per epoch, inference time per sample, inference costs per sample, and CO_2 emissions per sample. Financial costs are estimated using the Google Cloud pricing calculator¹, and CO_2 emissions are calculated with CodeCarbon (Lacoste et al., 2019). These comparisons provide insights into the trade-offs between large-scale

¹<https://cloud.google.com/products/calculator> (Assessed January 2025)

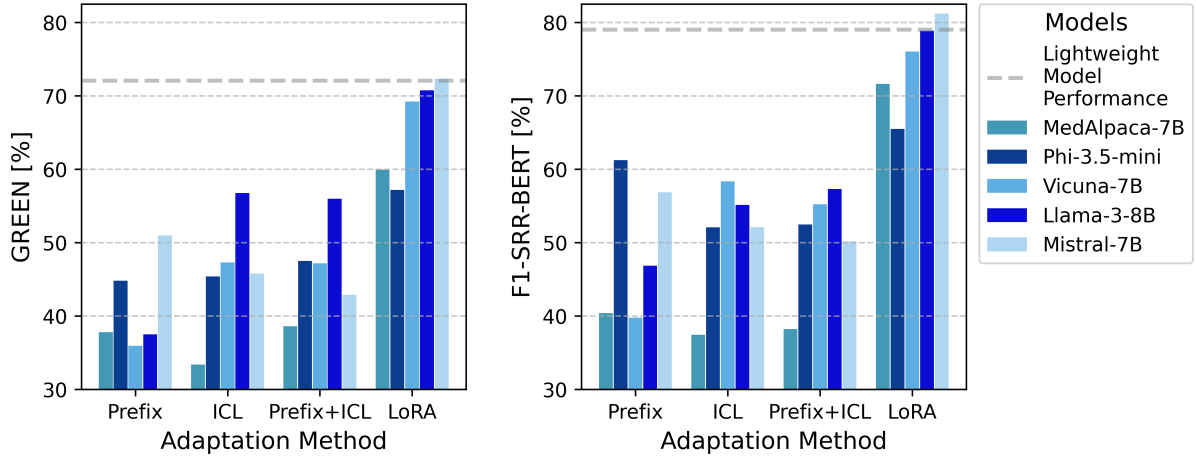


Figure 4: Comparison of LLM Adaptation Methods and the best performing lightweight model (BERT2BERT initialized from RoBERTa-PM-M3). (Left)/(Right) The figure depicts the GREEN Score/F1-SRR-BERT Score for five different LLMs across various adaptation methods, including prefix prompting, in-context learning (ICL), the combination of prefix prompting with ICL, and LoRA finetuning for five epochs.

LLMs and compact lightweight models in terms of both performance and resource efficiency.

4 Results

The models are evaluated using all metrics introduced in Section 3.2. We primarily report results using GREEN and F1-SRR-BERT Score, as they provide the most comprehensive assessments of clinical accuracy and structural consistency. However, unless stated otherwise, the observed trends hold across all metrics. A detailed comparison across all metrics is provided in Appendix A.5.

4.1 Comparison of Lightweight Models and Domain Adaptation

As introduced in Section 3.3, we initialized our lightweight models with the weights from different pretrained models. Specifically, we evaluate four different pretrained models as initializations for the BERT2BERT model and five for the T5 model (Tables 2 and 3). Each pretraining configuration was trained three times with different random seeds. Figure 3 presents the model performance for the GREEN and F1-SRR-BERT metrics, while a more comprehensive overview can be found in Table 4. For the BERT2BERT model, domain adaptation shows a clear but non-linear impact on performance. Pretraining on biomedical text improves GREEN by 0.4% over the general-text baseline, while adding radiology reports yields a more substantial 4.5% improvement. However, pretraining exclusively on radiology reports (RadBERT) provides only a marginal 0.3% increase. For the

T5 model, instruction-tuning alone leads to 0.3% improvement over the general-text baseline. Pre-training on biomedical text and radiology reports achieves a 2.5% gain, while using exclusively radiology reports leads to 4.4% increase. However, the biomedical text initialization (SciFive) underperforms the general baseline by 2.4%. Table 4 confirms that these trends persist across both datasets and sections, with scores for the Impression section being on average by $\approx 20\%$ higher. Overall, BERT2BERT models outperform T5 variants, with the best BERT2BERT model (RoBERTa-PM-M3) beating the best T5 (Clin-T5-Base) by 2.6% on GREEN and 1.5% on F1-SRR-BERT.

4.2 Adaptation of LLMs

We present the results of adapting LLMs to the structuring task as outlined in Section 3.4. Figure 4 visualizes the average test set performance on the GREEN and F1-SRR-BERT metrics across a selection of the proposed adaptation methods: prefix prompting, 2-shot in-context learning (ICL), the combination of prefix prompting and ICL, and LoRA finetuning. LoRA finetuning consistently achieves the highest performance across all models. The detailed breakdown of results across the structured Findings and Impression sections is provided in Tables 5 and 6 of the Appendix. Averaged across all five LLMs, 2-shot ICL improves performance compared to prefix prompting by 22.2%/20.6% in GREEN/F1-SRR-BERT on Findings and 9.6%/−1.0% on Impression. *Prefix+ICL* shows a 77.8%/79.2% improvement on Findings

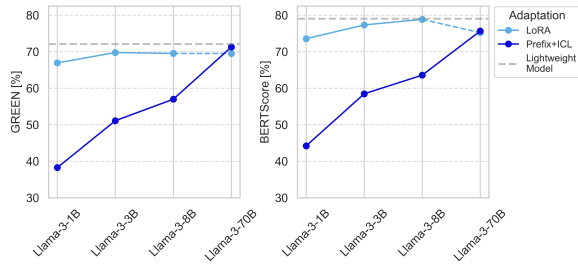


Figure 5: Model performance of LLaMA-3 models of increasing size. (Left/Right) The figure shows the GREEN and F1-SRR-BERT scores for adaptation using Prefix+ICL and LoRA finetuning, respectively. The result for the LLaMA-3-70B model with LoRA finetuning is indicated with a dashed line, as this configuration was trained for only one epoch—compared to five epochs for the other models—due to computational constraints.

but also -5.9% / -4.1% on Impression. LoRA finetuning achieves the highest scores overall, outperforming prefix prompting by 263%/237% on Findings and 8.7%/6.5% on Impression. Across LLMs, Llama-3-8B performs best in ICL methods, while Mistral-7B achieves the highest performance in LoRA finetuning. The overall best-performing configuration is Mistral-7B with LoRA finetuning.

4.3 Benchmarking

Building on these results, we benchmark our best lightweight model against LLaMA-3 models of increasing parameter counts. Figure 5 demonstrates a general positive correlation between the LLM’s model size and performance in structuring radiology reports, with the exception of LLaMA-3-70B. Despite being the largest model, it underperforms when adapted via LoRA, likely due to insufficient training. This size-performance trend is more evident with *Prefix+ICL* adaptation. While LLaMA-3-1B achieves only 53.0%/55.9% of the lightweight model’s performance (GREEN/F1-SRR-BERT), LLaMA-3-70B reaches 98.9%/95.8%. LoRA boosts LLaMA-3-1B to 92.9%/93.0%, and enables the larger variants to slightly outperform the lightweight model on the Findings section. However, when averaged across both sections, no LLM surpasses the lightweight model. Moreover, the relative benefit of LoRA over *Prefix+ICL* diminishes as model size increases, with both methods converging in performance—and LoRA occasionally underperforming—particularly on clinically relevant metrics such as F1-RadGraph, GREEN, and F1-SRR-BERT. Given these findings, we next turn to a cost analysis. As shown in Table 1, the

lightweight model offers considerable advantages in training time, financial cost, and environmental impact—producing only 8.3% and 0.7% of the CO_2 emissions of LLaMA-3-3B and 70B, respectively. Inference efficiency follows a similar pattern: even under the least favorable deployment scenario, the lightweight model exhibits up to 91.8% lower latency and 98.4% lower emissions than LLaMA-3-70B. Under optimal conditions, these savings exceed 99.9%.

4.4 Qualitative Analysis

To complement the quantitative analysis, Figure 1 presents a qualitative comparison of BERT2BERT, Mistral-7B, and expert-reviewed reports. Both models successfully adhere to our predefined template (see Figure 2 for reference), particularly in the Findings section, where content is well-aligned with organ system categories. A full test set analysis shows that the lightweight model correctly applies the Findings and Impression section headers in all cases, while the LLM deviates in 5% of instances, occasionally using all capital letters or omitting section names in less than 1% of reports. Both models, as well as expert annotations, generally include only relevant organ systems, but occasionally report less relevant negative findings (e.g., “Pleura: - No specific findings reported”). Complete omission of relevant findings occurs in less than 1% of cases, indicating high completeness in capturing clinical details. Differences in prioritization in the Impression section are observed in fewer than 5% of reports for both models, demonstrating occasional variation but overall consistency with expert-reviewed reports.

Table 1: Trade-off between model performance and computational costs for training and inference using total training time [h], CO_2 emission during training [kg], F1-SRR-BERT Score [%], inference time [s/sample], inference cost [\$/sample], and CO_2 emissions [mg/sample] across the best-performing BERT2BERT, LLaMA-3-3B, and LLaMA-3-70B models using NVIDIA A100-80GB GPUs.

	Model	Lightweight	3B LLM	70B LLM ^o
	# Parameters	0.28B	3.21B	70.6B
	Training time [h]	2.1	15.0	44.5 ^o
	Training CO_2 eq. [kg]	0.58	7.0	82.6 ^o
Inference	SRR-BERT [%]	79.1	77.4	75.2
	Time [s]	3.1 (0.16)*	10.7	1260 (37.7) [†]
	Cost [\$]	0.0043 (2e-4)*	0.015	1.76 (0.21) [†]
	CO_2 eq. [g]	0.075 (0.0038)*	0.25	67.7 (7.9) [†]

^o Only trained for 1 epoch. Trained on four GPUs instead of one.

* For single-sample (batch-wise) processing.

[†] Executed on 1 (4) NVIDIA A100 (80GB) GPU(s).

5 Discussion

In this paper, we propose lightweight, task-specific models for structuring radiology reports into a predefined template. Despite being 10–250 times smaller than finetuned LLMs, our models achieved comparable performance while offering significant advantages in speed, cost-efficiency, and sustainability. To enable large-scale supervised training, we leveraged GPT-4 as a weak annotator to generate a training dataset, aligning chest radiology reports from MIMIC-CXR and CheXpert Plus with their corresponding structured versions as ground truth. Since GPT-generated data may contain inconsistencies and biases, we evaluated all models on a human-reviewed test set. Our study focused on two types of lightweight models, BERT2BERT and T5. Overall, our BERT2BERT model performed best when initialized from RoBERTa-PM-M3, surpassing the best T5 variant, Clin-T5-Base, by 2.6% on GREEN. Our results further indicate that pre-training on biomedical texts - particularly radiology reports - generally improved model performance. However, despite being pretrained exclusively on radiology reports, the RadBERT model did not outperform general-text variants. This suggests that pretraining factors beyond the training corpus, such as architectural choices and optimization techniques, may also influence model performance. For example, RoBERTa-PM-M3 benefited from a distillation process from RoBERTa-large-PM-M3-Voc.

To balance performance with computational feasibility, we first restricted our comparison to LLMs within the 3-8B parameter tier, evaluating different adaptation techniques within this range. We showed that LoRA finetuning consistently outperformed prefix prompting and ICL methods. As shown in Table 6, this trend was primarily driven by performance differences on the Findings section. Given that our evaluation assessed each organ system independently and assigned zero points to missing or inconsistently labeled headers (e.g., '*Lungs and Airways*' vs. '*Lungs*'), the results suggest that LoRA finetuning more effectively aligned LLM outputs with the predefined reporting template. We believe that although organ system names are provided in both the prefix prompt (see Appendix A.1) and the ICL examples, the absence of iterative feedback mechanisms in these methods made it challenging for models to internalize and consistently enforce correct structured formatting.

Among the five evaluated LLMs and four adap-

tation techniques, Mistral-7B and LLaMA-3-8B achieved the best results. Notably, MedAlpaca-7B underperformed compared to general-domain models of similar size, suggesting that current medicine-specific LLMs may not yet offer clear advantages for structured report generation. We selected LLaMA-3 models with 1B, 3B, 8B, and 70B parameters for benchmarking our lightweight model against LLMs of increasing size in Section 4.3. Under the two most effective adaptation strategies—*Prefix+ICL* and LoRA—performance generally improved with model size, with LoRA finetuning ultimately enabling larger models to surpass the lightweight model on the Findings section. This came, however, at the cost of significantly longer training times and higher inference costs.

Our qualitative analysis in Section 4.4 showed that both models (the lightweight model and Mistral-7B LLM finetuned with LoRA) followed the predefined template when tested on expert-annotated reports, omitting relevant findings in less than 1% of cases. This suggests that lightweight models (<300M parameters) can effectively learn structured formatting while maintaining clinical accuracy. Furthermore, the results indicate that our GPT-generated annotations provided a sufficient training signal, though expert review remains crucial for ensuring data reliability.

6 Conclusion

We demonstrate that lightweight, task-specific models with less than 300M parameters can effectively structure radiology reports according to a predefined template, providing a practical and scalable alternative to LLMs, while addressing concerns around computational efficiency, data privacy, and deployment feasibility. Our best-performing lightweight model, a BERT2BERT architecture initialized from two pretrained RoBERTa-PM-M3 models, achieved competitive performance while maintaining a significantly lower computational footprint. While LLaMA-3 variants with more than 3 billion parameters achieved slightly better performance on the Findings section when finetuned with LoRA, the lightweight model operated at less than 25% of their inference cost and CO_2 emissions, making it a more resource-efficient solution. These findings reinforce the lightweight model’s viability for real-world clinical applications, where infrastructure limitations, privacy regulations, and sustainability concerns play a critical role.

Limitations

First, as discussed in Section 3.1, the labels used for training our specialized models and adapting the LLMs were generated from MIMIC-CXR and CheXpert Plus reports using GPT-4 as a weak annotator. While our prompt builds on previous work, we refined it to better align with our task’s requirements (e.g., explicitly specifying organ systems for the Findings section). However, GPT-4 may introduce biases, and to mitigate this, we evaluate model performance on an independent test set annotated by five radiologists.

Second, both MIMIC-CXR and CheXpert Plus originate from hospitals in the United States - Beth Israel Deaconess Medical Center (Boston, MA) and Stanford Hospital (Stanford, CA) - and contain only chest X-rays from adult patients. As a result, these datasets may lack demographic diversity, potentially limiting generalizability to other populations.

Third, as described in Section 3, all models take full free-form reports as input and generate structured reports comprising the following sections: Exam Type, History, Technique, Comparison, Findings, and Impression. However, for quantitative evaluation, we focus exclusively on Findings and Impression, as these sections are clinically critical and exhibit the highest variability. Other sections, such as Exam Type and History, often remain unchanged and can be directly copied from the original report, making them less relevant for assessing model performance.

Fourth, 1-shot and 2-shot ICL examples were manually selected from the training set to best represent the data distribution. While we initially applied algorithmic methods to optimize alignment, manual selection proved to improve performance. This introduces a potential selection bias, which may affect the generalizability of our ICL results.

Fifth, while we initially experimented with full-parameter finetuning for select LLMs, we found that it did not yield substantial performance improvements over LoRA. Given the significantly higher computational and time demands of full finetuning, we opted to use LoRA as an efficient adaptation strategy for all LLMs within our resource constraints.

Sixth, we initially also evaluated GPT-4 using prefix prompting and ICL. However, since it was used for data annotation and provided as a reference for radiologist, its results may be biased in its favor.

To account for this, we excluded GPT-4 from the discussion to avoid misleading comparisons.

Seventh, while we expected the LLMs—particularly the larger models—to outperform the lightweight model given their scale, this was not consistently observed under our current finetuning setup. Although we performed basic hyperparameter tuning and employed established adaptation techniques, the finetuning process may not have been sufficiently extensive or optimized to fully leverage the capabilities of these models. This is especially true for LLaMA-3-70B, which was limited to a single epoch of training due to computational constraints.

Eighth, while our selection of LLMs aims to represent both the current state of the art and a range of model sizes, one could argue for the inclusion of more domain-specific models tailored to the medical field. We include MedAlpaca-7B as a representative example, but find that it underperforms compared to general-domain models of similar scale, suggesting that current medicine-specific LLMs may not yet offer a clear advantage for the structuring task evaluated here.

References

- Marah Abidin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Lisa C Adams, Daniel Truhn, Felix Busch, Avan Kader, Stefan M Niehues, Marcus R Makowski, and Keno K Bresslem. 2023. Leveraging gpt-4 for post hoc transformation of free-text radiology reports into structured reporting: a multilingual feasibility study. *Radiology*, 307(4):e230725.
- Nurbanu Aksoy, Nishant Ravikumar, and Alejandro F Frangi. 2023. Radiology report generation using transformers conditioned with non-imaging data. In *Medical Imaging 2023: Imaging Informatics for Healthcare, Research, and Applications*, volume 12469, pages 146–153. SPIE.
- Anonymous. 2025. Automatic structured radiology report generation. Under review.
- Hassan B Arshad, Sara A Butt, Safi U Khan, Zulqarnain Javed, and Khurram Nasir. 2023. Chatgpt and

718	artificial intelligence in hospital level research: poten-	Langlotz. 2022. Improving the factual correctness of	774
719	tial, precautions, and prospects. <i>Methodist DeBakey</i>	radiology report generation with semantic rewards.	775
720	<i>cardiovascular journal</i> , 19(5):77.	<i>arXiv preprint arXiv:2210.12186</i> .	776
721	Laura Bergomi, Tommaso M Buonocore, Paolo Anton-	Daniel Pinto dos Santos, Elmar Kotter, Peter Milden-	777
722	azzo, Lorenzo Alberghi, Riccardo Bellazzi, Lorenzo	berger, and Luis Martí-Bonmatí. 2023. Esr paper	778
723	Preda, Chandra Bortolotto, and Enea Parimbelli.	on structured reporting in radiology—update 2023.	779
724	2024. Reshaping free-text radiology notes into struc-	<i>Insights into Imaging</i> , 14(1):199.	780
725	tured reports with generative question answering		
726	transformers. <i>Artificial Intelligence in Medicine</i> ,	Roselyn Gabud, Portia Lapitan, Vladimir Mariano, Ed-	781
727	154:102924.	uardo Mendoza, Nelson Pampolina, Maria Art An-	782
728	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	tonette Clariño, and Riza Theresa Batista-Navarro.	783
729	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	2023. A hybrid of rule-based and transformer-based	784
730	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	approaches for relation extraction in biodiversity lit-	785
731	Askell, et al. 2020. Language models are few-shot	erature. In <i>Proceedings of the 2nd Workshop on</i>	786
732	learners. <i>Advances in neural information processing</i>	<i>Pattern-based Approaches to NLP in the Age of Deep</i>	787
733	<i>systems</i> , 33:1877–1901.	<i>Learning</i> , pages 103–113.	788
734	Felix Busch, Lena Hoffmann, Daniel Pinto Dos Santos,	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	789
735	Marcus R Makowski, Luca Saba, Philipp Prucker,	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	790
736	Martin Hadamitzky, Nassir Navab, Jakob Nikolas	Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,	791
737	Kather, Daniel Truhn, et al. 2024. Large language	Alex Vaughan, et al. 2024. The llama 3 herd of mod-	792
738	models for structured reporting in radiology: past,	els. <i>arXiv e-prints</i> , pages arXiv–2407.	793
739	present, and future. <i>European Radiology</i> , pages 1–		
740	14.	Sebastian Griewing, Fabian Lechner, Niklas Gremke,	794
741	Pierre Chambon, Jean-Benoit Delbrouck, Thomas	Stefan Lukac, Wolfgang Janni, Markus Wallwiener,	795
742	Sounack, Shih-Cheng Huang, Zhihong Chen, Maya	Uwe Wagner, Martin Hirsch, and Sebastian Kuhn.	796
743	Varma, Steven QH Truong, Curtis P Langlotz, et al.	2024. Proof-of-concept study of a small language	797
744	2024. Chexpert plus: Hundreds of thousands of	model chatbot for breast cancer decision support—a	798
745	aligned radiology texts, images and patients. <i>arXiv</i>	transparent, source-controlled, explainable and data-	799
746	<i>e-prints</i> , pages arXiv–2405.	secure approach. <i>Journal of Cancer Research and</i>	800
747	Dong Chen, Shuo Zhang, Yueting Zhuang, Siliang	<i>Clinical Oncology</i> , 150(10):1–12.	801
748	Tang, Qidong Liu, Hua Wang, and Mingliang Xu.		
749	2024a. Improving large models with small models:	Suchin Gururangan, Ana Marasović, Swabha	802
750	Lower costs and better performance. <i>arXiv preprint</i>	Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	803
751	<i>arXiv:2406.15471</i> .	and Noah A Smith. 2020. Don’t stop pretraining:	804
752	Qi Chen, Yutong Xie, Biao Wu, Xiaomin Chen, James	Adapt language models to domains and tasks. <i>arXiv</i>	805
753	Ang, Minh-Son To, Xiaojun Chang, and Qi Wu.	<i>preprint arXiv:2004.10964</i> .	806
754	2024b. Act like a radiologist: Radiology report gen-		
755	eration across anatomical regions. In <i>Proceedings</i>	Tianyu Han, Lisa C Adams, Jens-Michalis Papaioan-	807
756	<i>of the Asian Conference on Computer Vision</i> , pages	nou, Paul Grundmann, Tom Oberhauser, Alexander	808
757	1–17.	Löser, Daniel Truhn, and Keno K Bresssem. 2023.	809
758	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	Medalpaca—an open-source collection of medical	810
759	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	conversational ai models and training data. <i>arXiv</i>	811
760	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion	<i>preprint arXiv:2304.08247</i> .	812
761	Stoica, and Eric P. Xing. 2023. <i>Vicuna: An open-</i>		
762	<i>source chatbot impressing gpt-4 with 90%* chatgpt</i>	Michael P Hartung, Ian C Bickle, Frank Gaillard, and	813
763	<i>quality</i> .	Jeffrey P Kanne. 2020. How to create a great radiol-	814
764	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	ogy report. <i>Radiographics</i> , 40(6):1658–1670.	815
765	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi		
766	Wang, Mostafa Dehghani, Siddhartha Brahma, et al.	Amir M Hasani, Shiva Singh, Aryan Zahergivar, Beth	816
767	2024. Scaling instruction-finetuned language models.	Ryan, Daniel Nethala, Gabriela Bravomontenegro,	817
768	<i>Journal of Machine Learning Research</i> , 25(70):1–53.	Neil Mendhiratta, Mark Ball, Faraz Farhadi, and	818
769	Raul Salles de Padua and Imran Qureshi. 2024. Leverag-	Ashkan Malayeri. 2024. Evaluating the performance	819
770	ing summary of radiology reports with transformers.	of generative pre-trained transformer-4 (gpt-4) in	820
771	<i>Artificial Intelligence in Health</i> , 1(4):85–96.	standardizing radiology reports. <i>European Radiol-</i>	821
772	Jean-Benoit Delbrouck, Pierre Chambon, Christian	<i>ogy</i> , 34(6):3566–3574.	822
773	Bluethgen, Emily Tsai, Omar Almusa, and Curtis P		
774		Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	823
775		Bruna Morrone, Quentin De Laroussilhe, Andrea	824
776		Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	825
777		Parameter-efficient transfer learning for nlp. In <i>In-</i>	826
778		<i>ternational conference on machine learning</i> , pages	827
779		2790–2799. PMLR.	828

829	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoy-	883
830	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	anov. 2020. Pretrained language models for biomed-	884
831	and Weizhu Chen. 2021. Lora: Low-rank adap-	ical and clinical tasks: understanding and extending	885
832	tation of large language models. <i>arXiv preprint</i>	the state-of-the-art. In <i>Proceedings of the 3rd clin-</i>	886
833	<i>arXiv:2106.09685</i> .	<i>ical natural language processing workshop</i> , pages	887
		146–157.	888
834	Saahil Jain, Ashwin Agrawal, Adriel Saporta,	Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lian-	889
835	Steven QH Truong, Du Nguyen Duong, Tan Bui,	min Zheng, Joseph Gonzalez, Ion Stoica, Xuezhe Ma,	890
836	Pierre Chambon, Yuhao Zhang, Matthew P Lungren,	and Hao Zhang. 2023. How long can context length	891
837	Andrew Y Ng, et al. 2021. Radgraph: Extracting	of open-source llms truly promise? In <i>NeurIPS 2023</i>	892
838	clinical entities and relations from radiology reports.	<i>Workshop on Instruction Tuning and Instruction Fol-</i>	893
839	<i>arXiv preprint arXiv:2106.14463</i> .	<i>lowing</i> .	894
840	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	895
841	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Optimizing continuous prompts for generation. <i>arXiv</i>	896
842	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	<i>preprint arXiv:2101.00190</i> .	897
843	laume Lample, Lucile Saulnier, et al. 2023. Mistral		
844	7b. <i>arXiv preprint arXiv:2310.06825</i> .		
845	Alistair Johnson, Lucas Bulgarelli, Tom Pollard,	Vladislav Lialin, Vijeta Deshpande, and Anna	898
846	Steven Horng, Leo Anthony Celi, and Roger Mark.	Rumshisky. 2023. Scaling down to scale up: A guide	899
847	2020. Mimic-iv. <i>PhysioNet</i> . Available online at:	to parameter-efficient fine-tuning. <i>arXiv preprint</i>	900
848	https://physionet.org/content/mimiciv/1.0/ (accessed	<i>arXiv:2303.15647</i> .	901
849	August 23, 2021), pages 49–55.	Chin-Yew Lin. 2004. Rouge: A package for automatic	902
		evaluation of summaries. In <i>Text summarization</i>	903
850	Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz,	<i>branches out</i> , pages 74–81.	904
851	Nathaniel R Greenbaum, Matthew P Lungren, Chih-		
852	ying Deng, Roger G Mark, and Steven Horng.	Yinhan Liu. 2019. Roberta: A robustly opti-	905
853	2019. Mimic-cxr, a de-identified publicly available	mized bert pretraining approach. <i>arXiv preprint</i>	906
854	database of chest radiographs with free-text reports.	<i>arXiv:1907.11692</i> , 364.	907
855	<i>Scientific data</i> , 6(1):317.		
856	Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H	Chandreen R Liyanage, Ravi Gokani, and Vijay Mago.	908
857	Lehman, Mengling Feng, Mohammad Ghassemi,	2024. Gpt-4 as an x data annotator: Unraveling its	909
858	Benjamin Moody, Peter Szolovits, Leo Anthony Celi,	performance on a stance classification task. <i>PloS one</i> ,	910
859	and Roger G Mark. 2016. Mimic-iii, a freely accessi-	19(8):e0307741.	911
860	ble critical care database. <i>Scientific data</i> , 3(1):1–9.	NCBI. 1996. PubMed.	912
861	Charles E Kahn Jr, Curtis P Langlotz, Elizabeth S Burn-	NCBI. 2000. PubMed Central (pmc).	913
862	side, John A Carrino, David S Channin, David M		
863	Hovsepian, and Daniel L Rubin. 2009. Toward	OpenAI. 2022. Gpt-3.5. https://openai.com/ .	914
864	best practices in radiology reporting. <i>Radiology</i> ,		
865	252(3):852–856.	Sophie Ostmeier, Justin Xu, Zhihong Chen, Maya	915
866	Dhruv Khullar, Xingbo Wang, and Fei Wang. 2024.	Varma, Louis Blankemeier, Christian Bluethgen,	916
867	Large language models in health care: Charting a	Arne Edward Michalson, Michael Moseley, Curtis	917
868	path toward accurate, explainable, and secure ai.	Langlotz, Akshay S Chaudhari, et al. 2024. Green:	918
869	<i>Journal of General Internal Medicine</i> , pages 1–3.	Generative radiology report evaluation and error no-	919
		tation. <i>arXiv preprint arXiv:2405.03595</i> .	920
870	Alexandre Lacoste, Alexandra Luccioni, Victor	Nicholas Pangakis, Samuel Wolken, and Neil Fasching.	921
871	Schmidt, and Thomas Dandres. 2019. Quantifying	2023. Automated annotation with generative ai re-	922
872	the carbon emissions of machine learning. <i>arXiv</i>	quires validation. <i>arXiv preprint arXiv:2306.00176</i> .	923
873	<i>preprint arXiv:1910.09700</i> .		
874	Andrew K Lampinen, Ishita Dasgupta, Stephanie CY	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	924
875	Chan, Kory Matthewson, Michael Henry Tessler,	Jing Zhu. 2002. Bleu: a method for automatic eval-	925
876	Antonia Creswell, James L McClelland, Jane X	uation of machine translation. In <i>Proceedings of the</i>	926
877	Wang, and Felix Hill. 2022. Can language models	<i>40th annual meeting of the Association for Computa-</i>	927
878	learn from explanations in context? <i>arXiv preprint</i>	<i>tional Linguistics</i> , pages 311–318.	928
879	<i>arXiv:2204.02329</i> .	Branislav Pecher, Ivan Srba, and Maria Bielikova. 2024.	929
880	Eric Lehman and Alistair Johnson. 2023. Clinical-t5:	Comparing specialised small and general large lan-	930
881	Large language models built using mimic clinical	guage models on text classification: 100 labelled	931
882	text. <i>PhysioNet</i> .	samples to achieve break-even performance. <i>arXiv</i>	932
		<i>preprint arXiv:2402.12819</i> .	933

934	Long N Phan, James T Anibal, Hieu Tran, Shaurya	Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan,	989
935	Chanana, Erol Bahadroglu, Alec Peltekian, and Gré-	Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser	990
936	goire Altan-Bonnet. 2021. Scifive: a text-to-text	Ururahy Nunes Fonseca, Henrique Min Ho Lee,	991
937	transformer model for biomedical literature. <i>arXiv</i>	Zahra Shakeri Hossein Abad, Andrew Y Ng, et al.	992
938	<i>preprint arXiv:2106.03598</i> .	2023. Evaluating progress in automatic chest x-ray	993
		radiology report generation. <i>Patterns</i> , 4(9).	994
939	Radiological Society of North America (RSNA). 2011.	Kuo Zhang, Xiangbin Meng, Xiangyu Yan, Jiaming Ji,	995
940	Radreport: Radiology reporting templates. template	Jingqian Liu, Hua Xu, Heng Zhang, Da Liu, Jingjia	996
941	rpt144 . Accessed: 2024-02-07.	Wang, Xuliang Wang, et al. 2025. Revolutionizing	997
942	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	health care: The transformative impact of large lan-	998
943	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	guage models in medicine. <i>Journal of Medical Inter-</i>	999
944	Wei Li, and Peter J Liu. 2020. Exploring the lim-	<i>net Research</i> , 27:e59069.	1000
945	its of transfer learning with a unified text-to-text		
946	transformer. <i>Journal of machine learning research</i> ,	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q	1001
947	21(140):1–67.	Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-	1002
		uating text generation with bert. <i>arXiv preprint</i>	1003
948	Delaram Rezaeikhonakdar. 2023. Ai chatbots and chal-	<i>arXiv:1904.09675</i> .	1004
949	lenges of hipaa compliance for ai developers and ven-		
950	dors. <i>Journal of Law, Medicine & Ethics</i> , 51(4):988–	Yuhao Zhang, Daisy Yi Ding, Tianpei Qian, Christo-	1005
951	995.	pher D Manning, and Curtis P Langlotz. 2018. Learn-	1006
		ing to summarize radiology findings. <i>arXiv preprint</i>	1007
952	Sascha Rothe, Shashi Narayan, and Aliaksei Severyn.	<i>arXiv:1809.04698</i> .	1008
953	2020. Leveraging pre-trained checkpoints for se-		
954	quence generation tasks. <i>Transactions of the Associ-</i>		
955	<i>ation for Computational Linguistics</i> , 8:264–280.		
956	Jaromir Savelka, Kevin D Ashley, Morgan A Gray,		
957	Hannes Westermann, and Huihui Xu. 2023. Can gpt-		
958	4 support analysis of textual data in tasks requiring		
959	highly specialized domain expertise? <i>arXiv preprint</i>		
960	<i>arXiv:2306.13906</i> .		
961	Jackson M Steinkamp, Charles Chambers, Darco Lale-		
962	vic, Hanna M Zafar, and Tessa S Cook. 2019. Toward		
963	complete structured information extraction from ra-		
964	diology reports using machine learning. <i>Journal of</i>		
965	<i>digital imaging</i> , 32:554–564.		
966	Arun James Thirunavukarasu, Darren Shu Jeng Ting,		
967	Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan,		
968	and Daniel Shu Wei Ting. 2023. Large language		
969	models in medicine. <i>Nature medicine</i> , 29(8):1930–		
970	1940.		
971	Shubo Tian, Qiao Jin, Lana Yeganova, Po-Ting Lai,		
972	Qingqing Zhu, Xiuying Chen, Yifan Yang, Qingyu		
973	Chen, Won Kim, Donald C Comeau, et al. 2024. Op-		
974	portunities and challenges for chatgpt and large lan-		
975	guage models in biomedicine and health. <i>Briefings</i>		
976	<i>in Bioinformatics</i> , 25(1):bbad493.		
977	Dave Van Veen, Cara Van Uden, Maayane Attias,		
978	Anuj Pareek, Christian Bluethgen, Malgorzata Po-		
979	lacin, Wah Chiu, Jean-Benoit Delbrouck, Juan		
980	Manuel Zambrano Chaves, Curtis P Langlotz, et al.		
981	2023. Radadapt: Radiology report summarization		
982	via lightweight domain adaptation of large language		
983	models. <i>arXiv preprint arXiv:2305.01146</i> .		
984	An Yan, Julian McAuley, Xing Lu, Jiang Du, Eric Y		
985	Chang, Amilcare Gentili, and Chun-Nan Hsu. 2022.		
986	Radbert: adapting transformer-based language mod-		
987	els to radiology. <i>Radiology: Artificial Intelligence</i> ,		
988	4(4):e210258.		

A Appendix

1009

A.1 GPT-4 prompt template for structuring of radiology reports

1010

The following prompt was executed with GPT-4 "Turbo 1106 preview" via Azure services to structure free-text radiology reports according to our template. The account was explicitly opted out of human review.

1011

1012

1013

Your task is to improve the formatting of a radiology report to a clear and concise radiology report with section headings.

Guidelines:

1. Section Headers: Each section should start with the section header followed by a colon. Provide the relevant information as specified for each section.
2. Identifiers: Remove sentences where identifiers have been replaced with consecutive underscores ('___').
3. Findings and Impression Sections: Focus solely on the current examination results. Do not reference previous studies or historical data.
4. Content Restrictions: Strictly include only the content that is relevant to the structured sections provided. Do not add or extrapolate information beyond what is found in the original report. If the original report doesn't contain the information necessary to generate a section, write the section header and then leave the section empty. Do not make up any findings.!

Sections to include (if applicable):

1. Exam Type: Provide the specific type of examination conducted.
2. History: Provide a brief clinical history and state the clinical question or suspicion that prompted the imaging.
3. Technique: Describe the examination technique and any specific protocols used.
4. Comparison: Note any prior imaging studies reviewed for comparison with the current exam.
5. Findings:

Describe all positive observations and any relevant negative observations for each organ or organ system under distinct headers. Start with the organ system name followed by a colon, then list observations.

Here is the corresponding template:

Organ 1:

 - Observation 1

Organ 2:

 - Observation 1
 - Observation 2

Use only the following headers for organ systems:

- Lungs and Airways
- Pleura
- Cardiovascular
- Hila and Mediastinum
- Tubes, Catheters, and Support Devices
- Musculoskeletal and Chest Wall
- Abdominal
- Other

6. Impression: Summarize the key findings with a numbered list from the most to the least clinically relevant. Ensure all findings are numbered.

The radiology report to improve is the following: \{report\}

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

A.2 Overview of model checkpoints and pre-training data

Table 2: Pretrained T5 models used for initialization along with details of their pretraining corpus.

Model	Description
T5-BASE (Raffel et al., 2020)	Original model, pre-trained on C4.
FLAN-T5-BASE (Chung et al., 2024)	Additional instruction-prompt tuning.
SCIFIVE (Phan et al., 2021)	Fine-tuned on PubMed Abstract (NCBI, 1996), and PubMed Central (NCBI, 2000).
CLIN-T5-SCI (Lehman and Johnson, 2023)	Fine-tuned on PubMed, MIMIC-III (Johnson et al., 2016), and MIMIC-IV (Johnson et al., 2020).
CLIN-T5-BASE (Lehman and Johnson, 2023)	Fine-tuned on MIMIC-III and MIMIC-IV.

Table 3: Pretrained RoBERTa models used for initialization of the BERT2BERT model along with details of their pretraining corpus.

Model	Description
RoBERTa-base (Liu, 2019)	Baseline version, pretrained on Books and Wikipedia.
BioMed-RoBERTa (Gururangan et al., 2020)	Pretrained on PubMed abstracts and PubMed Central.
RoBERTa-base-PM-M3-Voc-distill-align (Lewis et al., 2020)	Pretrained on PubMed abstracts, PubMed Central full-text articles, and MIMIC-III.
RadBERT-RoBERTa (Yan et al., 2022)	Fine-tuned on radiology reports from the Veterans Affairs health care system.

A.3 Considerations and hyperparameters for end-to-end training

We train all expert models (BERT2BERT and T5 instances) with the following set of hyperparameters:

- Cosine learning rate scheduler, starting at $1e^{-4}$, with 5% warm-up ratio before decay.
- Maximum of 10 epochs, with early stopping enabled by loading the best model at the end based on validation performance.
- Batch size of 32 per device for training and 16 for evaluation, with four gradient accumulation steps, resulting in an effective batch size of 128 for training.
- Adam optimizer with $\beta_2 = 0.95$ and weight decay of 0.1.
- Sequence lengths: Model processes a maximum input length of 370 tokens, with generated outputs constrained between 120 and 286 tokens.

We experimented with different learning rate schedulers and initial learning rates but found the here presented set to give better performance in the validation loss.

A.4 Considerations and hyperparameters for parameter-efficient fine-tuning

As discussed in Section 3.4, we initially finetune all LLMs using the same hyperparameters. We apply LoRA and adjust the target modules to align with each LLM’s architecture. We find that, due to their comparable size, using the same LoRA rank and scaling factor leads to a similar proportion of updated parameters across all models ($\sim 0.1\%$). We use the following set of hyperparameters:

- Cosine learning rate scheduler, starting at $1e^{-4}$, with 5% warm-up ratio before decay.
- Maximum of 5 epochs, with early stopping enabled by loading the best model at the end based on validation performance.
- LoRA adaptation with rank $r = 8$ and scaling factor $\alpha = 8$ to enable parameter-efficient fine-tuning.
- Batch size of 16 per device for training and 1 for evaluation, with 16 gradient accumulation

steps, resulting in an effective training batch size of 256.

- Adam optimizer with $\beta_2 = 0.95$ and weight decay of 0.1.

We use similar settings as in expert model fine-tuning but reduce the maximum number of epochs due to computational constraints. The results in Section 4.3 later confirm our initial estimate for the optimal LoRA rank.

A.5 Detailed Evaluations of Model Performance

Table 4: Detailed comparison of expert models. This table presents test set evaluations of our finetuned expert models initialized from different pre-trained checkpoints. Each model was trained three times with different random seeds and evaluated on the Findings sections of the MIMIC (F_M) and CheXpert (F_C) test sets, as well as their corresponding Impression sections (I_M and I_C).

Model	Section	BLEU	ROUGE-L	BERTScore	RadGraph	GREEN	SRR-BERT
BERT2BERT							
roberta-base	F_M	31.3	62.2	67.4	54.8	66.1	73.0
	F_C	30.6	59.0	64.7	50.1	63.0	69.4
	I_M	41.1	65.4	79.7	57.5	65.6	81.8
	I_C	51.1	74.9	86.3	66.1	82.0	94.5
roberta-biomed	F_M	31.6	60.4	65.4	53.1	62.8	70.4
	F_C	29.4	57.8	63.8	48.2	62.1	70.0
	I_M	34.0	65.5	79.9	58.0	69.1	81.8
	I_C	48.3	74.1	86.1	65.3	82.0	91.3
roberta-PM	F_M	33.3	62.6	67.4	54.3	67.0	71.9
	F_C	32.8	62.5	67.3	53.8	64.2	72.8
	I_M	42.0	66.1	79.8	56.5	71.8	81.4
	I_C	53.4	77.6	87.5	67.7	86.4	90.1
roberta-rad	F_M	32.6	62.1	66.8	54.9	64.8	71.8
	F_C	29.4	59.2	64.2	50.7	61.0	69.1
	I_M	42.3	67.5	80.6	58.9	69.7	81.7
	I_C	52.4	76.6	87.2	65.7	86.7	94.3
T5							
T5-Base	F_M	26.4	52.8	58.8	64.9	58.6	63.6
	F_C	26.0	57.2	61.9	49.1	59.7	66.5
	I_M	35.8	61.7	77.7	56.2	69.8	80.1
	I_C	48.5	73.2	85.8	67.9	81.2	87.1
Flan-T5-Base	F_M	27.9	55.9	61.0	48.0	59.3	65.4
	F_C	30.3	59.2	63.5	51.1	62.2	66.2
	I_M	37.3	62.0	77.6	55.5	66.2	77.8
	I_C	51.6	76.1	87.1	68.6	82.3	91.7
SciFive	F_M	24.1	49.3	55.6	43.4	56.4	62.0
	F_C	24.6	54.1	60.5	47.2	56.7	65.7
	I_M	38.6	63.2	78.8	59.5	71.8	82.9
	I_C	46.8	71.4	85.1	68.1	77.8	89.4
Clin-T5-Sci	F_M	28.7	59.0	64.4	50.7	62.4	68.9
	F_C	23.4	52.5	57.1	44.0	56.1	62.0
	I_M	33.6	59.4	76.2	51.4	63.8	76.3
	I_C	46.7	71.8	84.6	62.8	84.0	93.0
Clin-T5-Base	F_M	29.8	58.3	64.0	50.9	62.7	68.6
	F_C	27.1	57.3	62.0	49.0	60.9	68.1
	I_M	37.6	63.3	78.9	55.7	68.7	80.2
	I_C	48.4	74.8	85.5	67.9	88.8	94.6

Table 5: Comparison of LLM performance across different adaptation and finetuning methods. Results are averaged over all samples in the expert-reviewed MIMIC and CheXpert test sets and reported separately for the Findings and Impression sections. The highest score for each model across adaptation techniques is highlighted.

Model	Method	BLEU	ROUGE-L	BERTScore	Radgraph	GREEN	F1-Score
Findings Section							
Medalpaca-7B	Prefix	0.0	0.0	0.0	0.0	0.0	0.0
	1-shot ICL	0.0	0.2	1.4	0.1	0.1	0.9
	2-shot ICL ICL	0.0	0.0	0.0	0.0	0.0	0.0
	Prefix+ICL	0.0	2.3	7.6	0.7	11.4	5.4
	LoRA	19.7	45.4	50.5	41.3	51.0	57.1
Phi-3.5-mini	Prefix	11.0	34.6	38.9	26.7	38.1	46.5
	1-shot ICL	8.6	21.5	24.8	20.1	25.6	26.4
	2-shot ICL	6.8	20.1	24.1	18.5	23.2	25.8
	Prefix+ICL	14.3	35.3	40.7	28.8	38.3	43.6
	LoRA	17.8	43.8	49.5	39.0	46.7	52.9
Vicuna-7B	Prefix	0.0	0.0	0.0	0.0	0.0	0.0
	1-shot ICL	5.9	21.5	29.2	17.5	22.8	32.4
	2-shot ICL	7.1	19.8	24.6	17.0	22.6	28.2
	Prefix+ICL	7.4	23.7	30.9	19.0	26.3	32.2
	LoRA	32.7	62.1	66.8	54.2	66.1	70.6
LLaMA-3-8B	Prefix	2.4	10.9	12.8	8.6	13.1	12.7
	1-shot ICL	13.1	35.6	42.1	30.6	40.1	46.4
	2-shot ICL	13.7	36.4	42.1	31.1	38.0	46.4
	Prefix+ICL	18.7	44.7	51.1	37.6	48.6	56.6
	LoRA	35.0	62.9	68.4	54.4	68.1	74.0
Mistral-7B	Prefix	8.2	26.8	30.3	6.9	32.5	35.8
	1-shot ICL	6.5	15.2	18.4	14.7	16.9	19.4
	2-shot ICL	5.9	14.9	18.1	12.5	18.5	18.4
	Prefix+ICL	14.3	30.6	35.6	24.8	34.1	38.9
	LoRA	37.5	69.3	73.6	61.2	72.4	77.7
Impression Section							
Medalpaca-7B	Prefix	23.6	55.1	63.9	52.0	75.6	80.8
	1-shot ICL	23.3	54.0	60.7	50.3	66.8	74.1
	2-shot ICL	25.8	56.5	66.7	57.4	77.2	76.5
	Prefix+ICL	18.4	46.7	60.8	39.8	65.2	63.8
	LoRA	17.4	53.5	63.4	38.4	68.9	86.2
Phi-3.5-mini	Prefix	19.2	45.7	63.7	43.7	51.5	76.0
	1-shot ICL	24.4	48.6	66.8	47.7	65.3	77.8
	2-shot ICL	32.6	48.5	66.8	51.9	71.8	79.2
	Prefix+ICL	27.1	52.5	69.7	46.7	64.2	74.2
	LoRA	39.3	64.4	77.3	56.2	67.5	78.1
Vicuna-7B	Prefix	34.0	64.8	73.7	57.8	71.9	79.6
	1-shot ICL	38.8	64.7	77.5	61.5	71.9	84.3
	2-shot ICL	36.8	62.9	76.8	59.5	71.8	82.3
	Prefix+ICL	37.7	64.9	77.0	56.6	70.1	81.4
	LoRA	38.0	63.7	70.9	54.3	72.4	81.5
LLaMA-3-8B	Prefix	25.5	55.4	70.7	51.3	61.9	77.5
	1-shot ICL	9.7	27.5	45.6	33.1	73.5	63.9
	2-shot ICL	10.6	30.1	49.3	32.6	74.0	68.2
	Prefix+ICL	15.9	45.3	62.5	41.9	65.4	70.6
	LoRA	35.3	65.3	72.0	54.7	74.2	83.7
Mistral-7B	Prefix	33.6	63.4	78.4	56.0	69.5	78.0
	1-shot ICL	38.3	65.6	76.2	62.9	67.4	82.0
	2-shot ICL	39.2	66.0	77.2	62.9	67.4	82.0
	Prefix+ICL	42.6	70.7	80.2	61.9	55.0	86.1
	LoRA	42.3	67.6	74.8	57.0	76.1	84.8

Table 6: Detailed comparison of LLM adaptation methods for the Findings and Impression sections. The table shows average values across all five LLMs (excluding GPT-4), along with percentage changes relative to performance under prefix prompting.

Method	BLEU	ROUGE-L	BERTScore	Radgraph	GREEN	F1-SRR-BERT
Findings Section						
Prefix	4.31	14.4	16.4	8.43	16.7	19.7
1-shot ICL	6.79	18.8	23.2	16.6	21.1	25.1
	↑57.5%	↑30.2%	↑41.4%	↑96.8%	↑26.1%	↑27.3%
2-shot ICL	6.67	18.2	21.8	15.8	20.4	23.8
	↑54.8%	↑26.3%	↑33.0%	↑87.4%	↑22.2%	↑20.6%
Prefix+ICL	11.0	27.3	33.2	22.2	29.7	35.3
	↑155%	↑89.6%	↑102%	↑163%	↑77.8%	↑79.2%
LoRA	28.5	56.7	61.7	50.0	60.7	66.5
	↑562%	↑293%	↑277%	↑493%	↑263%	↑237%
Impression Section						
Prefix	27.2	56.9	70.1	52.1	66.1	78.4
1-shot ICL	26.9	52.0	65.3	50.7	70.4	77.0
	↓-1.1%	↓-8.5%	↓-6.8%	↓-2.7%	↑6.5%	↓-11.8%
2-shot ICL	26.8	52.8	67.3	52.8	72.4	77.6
	↓-1.5%	↓-7.2%	↓-3.9%	↑1.3%	↑9.6%	↓-1.0%
Prefix+ICL	28.4	56.0	70.0	49.4	62.2	75.2
	↑4.4%	↓-1.6%	+0.0%	↓-5.2%	↓-5.9%	↓-4.1%
LoRA	34.4	62.9	71.6	52.1	71.8	83.5
	↑26.8%	↑10.6%	↑2.2%	+0.0%	↑8.7%	↑6.5%

Table 7: Comparison of lightweight and LLM model performance. Results are averaged over all samples in the expert-reviewed MIMIC and CheXpert test sets and reported separately for the Findings and Impression sections. The highest score for each model across adaptation techniques is highlighted.

Model	Method	BLEU	ROUGE-L	BERTScore	Radgraph	GREEN	F1-Score
Findings Section							
BERT2BERT	Full Training	32.9	62.6	67.4	54.0	66.4	72.3
LLaMA-3-1B	Prefix+ICL	3.7	11.6	17.3	12.2	11.9	16.5
	LoRA	29.8	58.8	64.0	50.5	62.3	67.9
LLaMA-3-3B	Prefix+ICL	10.9	29.6	36.4	24.7	33.3	40.8
	LoRA	33.4	65.6	69.8	54.6	68.8	75.4
LLaMA-3-8B	Prefix+ICL	18.7	44.7	51.1	37.6	48.6	56.6
	LoRA	35.0	62.9	68.4	54.4	68.1	74.0
LLaMA-3-70B	Prefix+ICL	25.4	53.3	60.2	41.3	53.4	63.1
	LoRA	30.2	59.1	64.2	51.2	63.3	68.9
Impression Section							
BERT2BERT	Full Training	47.7	71.9	83.7	62.1	77.8	85.8
LLaMA-3-1B	Prefix+ICL	21.7	51.6	65.8	44.6	64.6	71.9
	LoRA	39.3	64.5	78.9	55.4	71.6	79.2
LLaMA-3-3B	Prefix+ICL	21.2	48.9	66.0	46.0	68.9	76.2
	LoRA	42.1	64.9	78.3	58.7	70.7	79.3
LLaMA-3-8B	Prefix+ICL	15.9	45.3	62.5	41.9	65.4	70.6
	LoRA	35.3	65.3	72.0	54.7	74.2	83.7
LLaMA-3-70B	Prefix+ICL	21.4	57.5	68.5	69.0	89.2	88.3
	LoRA	32.3	64.8	77.9	57.6	75.8	81.5

Table 8: Template adherence errors across the three best-performing models on 233 test samples.

Evaluation Category	BERT2BERT	LLaMA-3-8B	LLaMA-3-70B
Missing or misspelled headers	0	0	0
Different organ system names	0	14	35
Inconsistencies in bullet/enumeration formatting	0	80	61
Mismatch of mentioned organ systems	130	136	141
of which potentially irrelevant	100	113	111
of which potentially relevant	30	23	30