# A Structured, Tagged, and Localized Visual Question Answering Dataset with Full Sentence Answers and Scene Graphs for Chest X-ray Images

### **Anonymous Author(s)**

Affiliation Address email

# **Abstract**

Visual Question Answering (VQA) enables targeted and context-dependent analysis of medical images, such as chest X-rays (CXRs). However, existing VQA datasets for CXRs are typically constrained by simplistic and brief answer formats, lacking localization annotations (e.g., bounding boxes) and little metadata (e.g., region or radiological finding/disease tags). To address these limitations, we introduce *MIMIC-Ext-CXR-QBA* (abbr. *CXR-QBA*), a large-scale CXR VQA dataset derived from MIMIC-CXR, comprising 42 million QA-pairs with multi-granular, multi-part answers, detailed bounding boxes, and structured tags. We automatically generated our VQA dataset from scene graphs (also made available), which we constructed using LLM-based information extraction from radiology reports. After automatic quality assessment, we identified 31M pre-training and 7.5M fine-tuning grade QA-pairs, providing the largest and most sophisticated VQA dataset for CXRs to date. Tools for using our dataset and the construction pipeline are available at https://anonymous.4open.science/r/mimic-ext-cxr-qba/.

# 15 1 Introduction

2

5

6

8

9

10

12

13

14

With the emergence of Large Language Models (LLMs) and Large Multimodal Models (LMMs), 16 interactive and conversational tasks have gained popularity in medical image analysis, particularly in 17 the context of *chest X-ray (CXR)* interpretation [1]–[7]. A prominent example of such interactive tasks 18 is Visual Question Answering (VQA), where a model is presented with an image and a corresponding 19 textual question, and is tasked with generating an answer. Unlike conventional medical imaging 20 approaches, which always produce the same output (such as classification labels, bounding boxes, or 21 textual reports) for a given image, VQA enables users to interactively explore and interpret images 22 in a context-dependent manner. Training robust VQA models for medical applications requires 23 high-quality, large-scale training datasets. Existing CXR VQA datasets [1], [8]-[13] suffer from several limitations: (i) they often contain only short and simplistic answers, (ii) they lack localization 25 information (such as bounding boxes), and (iii) they provide little structured metadata (e.g., region and finding/disease annotations, or uncertainty estimates). Additionally, their relatively small size 27 constrains their utility for pretraining. 28

To address these challenges, we propose a pipeline for automatic VQA dataset creation and apply it to construct a new large-scale CXR VQA dataset. Unlike prior datasets, each question-answer (QA) pair includes multi-granular, multi-part answers composed of full sentences in the style of radiology reports. Furthermore, our dataset provides detailed bounding boxes and additional structured tags (e.g., findings and regions), enhancing interpretability and facilitating the development of more advanced and transparent medical VQA models. Fig. 1 shows examples of our generated QA-pairs.



(a) Indication question.

(b) Study abnormality question.



(c) Region abnormality question.

(d) **Finding** question.

Figure 1: Examples of question-answer (QA) pairs for each of our four different types of questions. For each question (for a given chest X-ray), a detailed answer with sentences in the style of free-text radiology reports is given, supplemented by bounding boxes (for both positive and negative answers), and a set of tags (e.g. regions, findings, certainty, etc.). For more examples, we refer to Appendix A.

### 5 Our contributions are as follows:

36

37

38

39

40

41

42

43

44

47

48

49

50

51

52

53

54

56

57

58

59

- We propose an automatic scene graph construction method as an intermediate step for VQA dataset creation, utilizing LLMs, semantic entity mapping, and localization models.
- We propose a question-answer generation strategy based on the extracted scene graphs.
- Building on this approach, we introduce MIMIC-Ext-CXR-QBA (abbr. CXR-QBA), a 42M QA-pair VQA dataset derived from MIMIC-CXR [14], to be published on Physionet [15].
- We automatically evaluate the quality of the generated QA-pairs, identifying 31.2M pairs as pre-training grade and 7.5M of these as fine-tuning grade.
- We provide a detailed analysis of our dataset and demonstrate its utility on the newly proposed structured VQA task.

### 5 2 Related Work

VQA Datasets for Chest X-Rays VQA datasets (shown in Tab. 1) are scarce in the medical imaging domain, with most notable examples being VQA-RAD [8] and SLAKE [9], which are hand-labeled but limited in size (3.5K and 14K QA-pairs, respectively). On the other hand, VQA-Med at ImageCLEF 2019 [10] was automatically constructed using QA templates based on image annotations, which may limit its answer quality. To improve the quality, PMC-VQA [11] used an LLM to generated QA-pairs based on provided captions. VOA datasets for chest X-rays include MIMIC-Ext-MIMIC-CXR-VQA [12], [15] and Medical-CXR-VQA [13], [15], [16], which contain hundreds of thousands of QA-pairs (in these cases derived from MIMIC-CXR). These datasets rely on templates but use radiology reports as their original information source, where MIMIC-

Table 1: Comparison of medical VQA datasets. We present the currently largest dataset, additionally providing boxes and tags for the answers.

Dataset	#QA	Boxes	Tags	Answers
CXR-QBA (Ours)	42.2M	1	/	detailed
VQA-RAD [8]	3.5K	Х	Х	brief
SLAKE [9]	14K	Х	Х	brief
ImageCLEF [10]	15K	Х	Х	brief
PMC-VQA [11]	227K	Х	Х	brief
MIMIC-CXR-VQA [12]	377K	Х	Х	brief
Medical-CXR-VQA [13]	780K	Х	Х	brief
CheXinstruct [1]	8.5M	X	X	brief

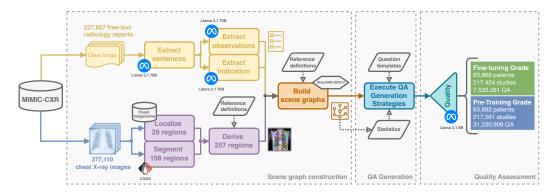


Figure 2: Overview of our dataset construction pipeline. First, we construct scene graphs based on information extracted from the radiology reports and regions localized in the images. Next, we generate question-answer pairs based on templates and the scene graphs. Finally, we automatically assess the quality of generated QA-pairs.

Ext-MIMIC-CXR-VQA leverages Chest ImaGenome's [17] scene graphs and Medical-CXR-VQA 62 employs an LLM-based extraction strategy similar to ours but without semantic entity mapping, localization, and extraction of textual descriptions. The largest chest X-ray VQA dataset to date, 64 CheXinstruct [1], contains 8.5M QA-pairs with images from multiple data sources. However, com-65 pared to our dataset, its questions and answers are less diverse, being purely template-based and 66 derived from dataset annotations instead of being directly conditioned on the reports. Additionally, 67 none of the described datasets provide the level of detail and annotation richness found in our dataset, 68 which includes bounding boxes, tags, and more detailed, multi-part answers that mirror radiology 69 report sentences. 70

Grounded Report Generation While localization is not yet common for medical VQA tasks, grounded report generation, i.e. predicting radiology reports with bounding boxes is gaining popularity. Notable examples include MAIRA-2 [18], trained on reports manually annotated with bounding boxes and MedTrinity-25M [7], a large-scale public dataset with automatically generated reports with bounding boxes. ChEX [2] is another model producing textual answers with bounding boxes. While being conditioned on textual prompts, ChEX does not support VQA tasks.

Scene Graph Construction for Chest X-Rays During our VQA dataset construction, we automatically derive scene graphs from radiology reports. A similar approach is employed by Chest ImaGenome [15], [17], [19], which uses rule-based information extraction, and RadGraph [20], which uses a relation extraction model. In contrast, our approach leverages LLM-based extraction with semantic entity mapping, enabling more comprehensive graph construction. Notably, our method defines a larger set of (localized) regions (257) and findings (221) compared to Chest ImaGenome (29 regions, 53 findings), making it a more robust foundation for VQA tasks.

# 84 3 The CXR-QBA Dataset

We present our dataset *CXR-QBA*, a large-scale *chest X-ray (CXR)* VQA dataset derived from *MIMIC-CXR* [14], [15], [21], consisting of more than 42M QA-pairs. As shown in Fig. 1, each QA sample (for a given chest X-ray) consists of a *question (Q)*, a *bounding box (B)* supplemented *answer (A)*, and additional tags (e.g. for regions, findings, certainties, and more).

To build our dataset, we propose an automatic pipeline highlighted in Fig. 2. More specifically, we

To build our dataset, we propose an automatic pipeline highlighted in Fig. 2. More specifically, we first construct (visually grounded) scene graphs based on the MIMIC-CXR radiology reports using LLM-based information extraction, semantic concept mapping, and localization models (Sec. 3.1). These scene graphs provide a structured description of the study, including sentences (derived from the report) for individual observations. They serve as a data source for our question-answer generation, where we utilize both template-based answers and answers derived from the rewritten report sentences (Sec. 3.2). Finally, we automatically assess the quality of question-answer pairs using LLM-based evaluations (Sec. 3.3). Further details are provided in Appendices D and E.

### 3.1 Scene Graph Construction

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

139

Given a MIMIC-CXR study with a radiology report and accompanying CXRs, we construct a scene graph (Appendix D.1) consisting of sentence nodes, observation nodes, region nodes, an indication node, and edges between them. Sentence nodes are directly extracted from the reports, containing 100 the raw sentences and their identified section names. Observation nodes represent individual aspects 101 described in the report's FINDINGS or IMPRESSION section, containing (i) a textual description, (ii) 102 bounding boxes for associated CXR images and (iii) additional tags, such as positivity, certainty, 103 laterality, regions, and finding classes. Region nodes are created for mentioned anatomical structures 104 and key regions. The indication node contains information from the INDICATION section, including a 105 textual description and an individual observation node, derived from the FINDINGS and IMPRESSION 106 sections, that can act as an answer to the indication. We construct these scene graphs in three steps: (a) region localization, (b) information extraction and (c) building the graphs using entity mapping. 108 We refer to Appendix E.1 for details. 109

Region Localization The bounding boxes in our scene graphs (and the derived QA-pairs) are based on fine-grained anatomical structures, allowing us to localize associated findings very precisely. We use the CXAS [22], [23] model to predict segmentation masks of 158 anatomical structures on the 377 110 CXRs from MIMIC-CXR-JPG [15], [24], [25]. Additionally, we use the bounding boxes provided by the Chest ImaGenome [15], [17], [19] dataset, which are provided for 29 anatomical structures in most frontal images of MIMIC-CXR. Next, we derive a total of 257 localized anatomical structures based on combinations (e.g. intersections, unions, super bounding boxes, etc.) of the available masks and bounding boxes. Finally, we discard any masks or boxes that are too small and derive bounding boxes from the segmentation masks. Note that we define 53 further regions/structures that are either non-localized (e.g. interstitial) or for which we do not have bounding boxes, leading to a total of 310 structures/regions.

**Information Extraction** We use the 227 827 free-text radiology reports provided by MIMIC-CXR as the main source of information for our scene graphs. Using the Llama 3.1 70B [26] model with few-shot prompting, we extract the relevant information (tags and textual descriptions) in three steps. First, we extract individual sentences from the reports and detect their sections. Next, we extract information about the INDICATION section and detect which FINDINGS or IMPRESSION sentences may provide information related to the indication. Finally, we extract individual observations described in the FINDING/IMPRESSION sentences.

Building Scene Graphs using Entity Mapping Given the extracted information from the reports 128 and the computed bounding boxes, we now construct the final scene graph. Therefore, we first map extracted tags to pre-defined sets of values, our reference definitions. This assures high quality and 130 consistency of the scene graphs and enables mapping of observations to the extracted bounding 131 boxes. The reference definitions are based on tags used in other datasets (including PadChest [27] 132 and Chest ImaGenome [19]) as well as SNOMED-CT [28]) and have been verified by clinical experts. 133 They include synonym lists, hierarchies, and relationships. For more robust mapping, we utilize the 134 BioLORD [29] model as a sentence transformer and identify the closest matching concept based on 135 their semantic embeddings. Additionally, we try to fill in missing information where possible, such 136 as inferring the region from an identified finding. Finally, we build a tree of region nodes (using the 137 reference data) and attach the indication information extracted from the report. 138

### 3.2 Question-Answer Generation

We generate question-answer pairs (Appendix D.2) using a template-based approach based on the 140 information available in the scene graphs, incorporating the textual descriptions from the observation nodes – which have been derived directly from the report – to provide diverse and fine-grained answers. Each answer may consist of multiple answer parts (as shown in Fig. 3 and Appendix A), 143 each describing an individual aspect of the answer with its own sentence, bounding boxes, and tags. 144 We categorize answer parts into three types: (i) main-answers, (ii) details, and (iii) related-information, 145 allowing for controlled answer granularity. Answer parts are generated either from templates using 146 scene graph information or directly from observation nodes (Appendix E.2). Answer parts may also 147 be structured hierarchically, where we use parent-child edges from the scene graph. 148

To generate the question-answer pairs, we employ different strategies for the four types of questions (shown in Fig. 1):

- Indication: We use the paraphrased indication as the question and create the answer based on the indication node in the scene graph, answering the indication based on information in the FINDINGS and IMPRESSION sections.
- 2. **Study abnormality**: We generate study-level questions using 13 different templates, with answer parts (Fig. 3) based on (filtered) observation nodes.
- 3. **Region abnormality**: We generate questions about individual regions using 6 different templates, considering any region mentioned and additionally randomly sampling non-mentioned regions for balancing.
- 4. Finding: We generate questions about individual findings using 7 different templates, considering any finding mentioned and additionally randomly sampling non-mentioned findings for balancing.

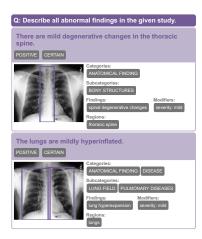


Figure 3: Answer with multiple parts for different aspects, each with a sentence, tags, and boxes.

### 3.3 Quality Assessment

The dataset construction procedure described so far allows us to automatically generate large amounts of QA-pairs. However, in each of the steps, errors may be introduced, affecting the overall quality of the datasets. For example, errors during information extraction could lead to incorrect tags, therefore leading to incorrectly filled answer templates or incorrectly selected observation nodes for answers.

In order to identify and filter such cases, we employ an automatic quality assessment strategy using an LLM as a judge. More specifically, we use Llama 3.1 8B [26] to rate question and answers by the following five criteria: *entailment* (does the answer factually align with the original report?), relevance (is the answer relevant to the question?), completeness (is the answer missing something?), as well as question and answer clarity (is the question/answer clear and grammatically correct?). Additionally, we assess the quality of the used scene graphs by identifying missing information (e.g. missing tags or localization) or issues during the construction process. Finally, we algorithmically combine these individual assessments to compute an overall quality rating as one of A++, A+, A+B, C, D, or not rated (see Appendix D.3). Based on these ratings, we propose two subsets, one for pre-training and one for fine-tuning. We exclude all non-frontal images from these datasets, as the localization quality on these images is comparatively low due to limitations in the localization models. All QA-pairs with a grade of A or better are labeled as *fine-tuning grade*, resulting in 7.5M pairs, while samples with grade B or better are considered *pre-training grade*, resulting in 31.2M pairs. 

# 4 Evaluation and Analysis

### 4.1 Evaluation of the Scene Graphs

We evaluate our scene graphs by comparing their tags and bounding boxes to hand-labeled expert annotations on MIMIC-CXR, using the scene graphs from Chest ImaGenome [17] as a baseline. First, we evaluate the plausibility of finding tags by comparing study-level labels derived from our scene graphs to two reference annotation sets: the radiologist annotations in MIMIC-CXR-JPG v.2.1.0 [24] with 13 CheXpert [33] classes and the CXR-LT 2024 [15], [30], [34] gold-standard dataset (task 2 test set) with 12 additional rare (long-tail) classes. As shown in Tab. 2a, our approach (slightly) outperforms Chest ImaGenome, with strong improvements (20%) on long-tail classes, demonstrating the value of our fine-grained finding tags (221 classes) in capturing nuanced study details. To evaluate the accuracy of finding bounding boxes, we compare them with annotations from MS-CXR [15], [31], [35] (6 classes) and REFLACX [15], [32], [36] (18 classes). We compute study-level pixel masks for each finding tag. We calculate pixel-level Intersection-over-Union (IoU), Intersection-over-Prediction (IoP), and Intersection-over-Target (IoT) for each finding class, considering only image

Table 2: Evaluation of our scene graphs, comparing finding tags (a) and associated bounding boxes (b) to expert annotations on MIMIC-CXR subsets, with 95% confidence intervals (bootstrapping, n=1000). Compared to Chest ImaGenome's [17] scene graphs, we achieve competitive or superior performance, showing that our construction process yields plausible scene graphs.

(a) Evaluation of finding tags against 13 CheXpert (CXP) classes from the MIMIC-CXR-JPG test set and 25 classes, 13 CXP and 12 long-tail (LT) classes, from the CXR-LT 2024 gold standard dataset (Appendix C.2). We report the Matthews Correlation Coefficient (MCC) macro-averaged over different finding subsets (CXP-5, CXP-7, CXP-13, LT) and micro-averaged. Compared to Chest ImaGenome, we produce slightly more accurate tags, performing especially well on long-tail classes, highlighting the importance of our fine-grained tags.

	MIM	MIMIC-CXR-JPG [24] Test [MCC]				CXR-LT 2024 [30] Gold [MCC]			
Classes	CXP-5	CXP-7	CXP-13	Micro	CXP-7	CXP-13	LT-only	CXR-LT	Micro
Ours (scene graphs)	<b>0.8</b> [0.77, 0.82]	<b>0.81</b> [0.79, 0.84]	<b>0.69</b> [0.67, 0.71]	<b>0.71</b> [0.69, 0.73]	<b>0.65</b> [0.61, 0.69]	<b>0.57</b> [0.54, 0.6]	<b>0.71</b> [0.67, 0.74]	<b>0.64</b> [0.61, 0.66]	<b>0.67</b> [0.65, 0.69]
Chest ImaGenome	0.78 [0.75, 0.81]	0.8 [0.78, 0.83]	0.66 [0.64, 0.69]	0.67 [0.65, 0.68]	<b>0.65</b> [0.61, 0.68]	0.56 [0.54, 0.59]	0.59 [0.55, 0.63]	0.58 [0.55, 0.6]	0.64 [0.62, 0.66]

(b) Evaluation of finding bounding boxes against 6 finding classes from MS-CXR and 18 classes from REFLACX (Appendix C.3). We report the pixel-level Intersection-over-Union (IoU), Intersection-over-Prediction (IoP), and Intersection-over-Target (IoT), each thresholded at 30%, and micro-averaged. Compared to Chest ImaGenome, our bounding boxes are better matching the hand-labeled boxes, especially leading to smaller and more precise boxes (larger IoP), which we assume is due to our more fine-grained region annotations.

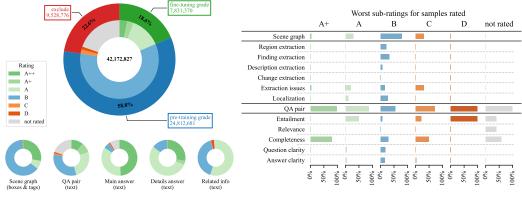
	N	AS-CXR [31	]	REFLACX [32]			
	[IoU@30]	[IoP@30]	[IoT@30]	[IoU@30]	[IoP@30]	[IoT@30]	
Ours (scene graphs)	<b>0.51</b> [0.47, 0.54]	<b>0.56</b> [0.52, 0.6]	0.94 [0.92, 0.96]	<b>0.45</b> [0.44, 0.47]	<b>0.54</b> [0.53, 0.56]	0.87	
Chest ImaGenome	0.45 [0.42, 0.49]	0.48 [0.45, 0.52]	<b>0.98</b> [0.97, 0.99]	0.42 [0.4, 0.43]	0.46 [0.44, 0.47]	<b>0.95</b> [0.94, 0.96]	

pairs with positive predictions and targets. Thresholding at 30% IoU/IoP/IoT, we micro-average the results, reporting the percentage of accurately localized finding-boxes in Tab. 2b. On the IoU metric, our scene graphs perform slightly better than the ones from Chest ImaGenome. The low IoP values indicate that bounding boxes are often too large, but high IoT values suggest that they generally cover the finding boxes well. This discrepancy arises because bounding boxes are derived from anatomical regions mentioned in reports, whereas hand-labeled annotations are more precise. Notably, our approach produces more precise boxes (higher IoP) than Chest ImaGenome, likely due to our large number of fine-grained region annotations (257 region classes).

Our analysis confirms that our scene graphs contain plausible finding tags and bounding boxes, with competitive or better quality than Chest ImaGenome. The bounding box quality, in turn, validates the plausibility of our region tags. Overall, our construction process yields high-quality scene graphs, making them a reliable foundation for generating QA samples.

### 4.2 Quality of the QA-Samples

We assessed the quality of our 42.2M QA-pairs using an LLM-as-a-judge approach (Sec. 3.3). Results are shown in Fig. 4a. We found that 18.6% were fine-tuning grade, 58.8% were pre-training grade, and 22.6% were marked for exclusion. Notably, 85% of individual main answers were rated A or higher. We also analyzed the main causes of ratings (Fig. 4b) and found that A+ samples were limited by minor incompleteness (minor details missing), A samples by minor entailment aspects (facts not explicitly mentioned in the report), while B samples were restricted by issues with region/finding/localization extraction, completeness, and text clarity. Ratings C were caused by major incompleteness or extraction issues, ratings D by contradicting entailments, while non-rated samples where due to the LLM-judge not producing parsable outputs. Using a larger LLM judge (Llama 3.1 70B), tested on a subset, reduced exclusions by 20%, but we opted for the smaller model (Llama 3.1 8B) to reduce computational requirements (we refer to Appendix C.1 for further details). Our analysis shows that even pre-training grade samples provide factually accurate answers with minor flaws, making them suitable for pre-training purposes.



(a) Overall rating (top) and sub-ratings (bottom).

(b) Reasons for ratings.

Figure 4: Results of quality assessment (Sec. 3.3). We identified a significant amount of fine-tuning grade samples, while even pre-training grade samples provide factually accurate answers, especially having high quality main answers.

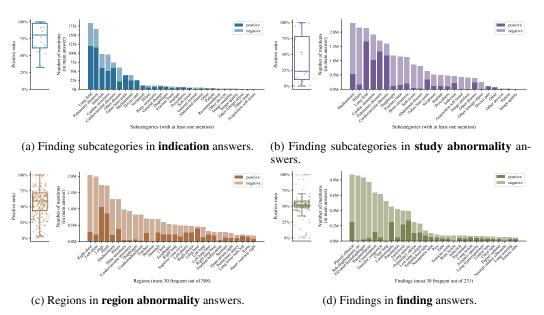


Figure 5: Distribution of tags (finding subcategories, regions, findings) mentioned in answers of different question types (indication, study abnormality, region abnormality, finding). We show their positive ratios, i.e. how often they are mentioned in positive versus in negative answers (**left**) and plot the number of positive and negative mentions of the most frequent tags (**right**). These fine-grained tags enable filtering and balancing the dataset or can be used as additional supervision.

# 4.3 Finding- and Region-Distribution in QA-Samples

Our answers include additional tags for findings (and their categories), regions, and answer positivity (positive or negative finding), enabling filtering and balancing for specific applications. For instance, undersampling negative answers can help mitigate model biases towards negative predictions. In Fig. 5 we analyze the distribution of these tags. We observe that indication questions tend to have more positive mentions (Fig. 5a) – as there is a specific indication to check for – while study abnormality questions have more negative ones (Fig. 5b) – as many samples are negative overall. In region abnormality questions, most regions are mentioned slightly more often with positive than with negative findings (Fig. 5c), while for finding questions mentions are mostly balanced (Fig. 5d). This shows the success of our balanced region/finding sampling used for these two question types.

### 4.4 Answer Characteristics

Our QA-samples provide detailed free-text answers consisting of one or even multiple sentences (i.e. answer parts). In Fig. 6, we analyze the distribution of lengths of these answers and study differences between types of answers or questions. The median answer length is 14 words, with similar lengths for most question types except for indication questions, where answers are much longer (46 words). We also observe that related information answers are much longer (22 words) than main answers (9 words) or details answers (7 words), which is expected as they can provide a lot additional context to the answers. Answers describing positive findings are typically very long (18 words), considerably longer than negative finding answers (10 words). This highlights that our dataset provides nuanced finding description in their answers, following the level of detail typically present in radiology reports.

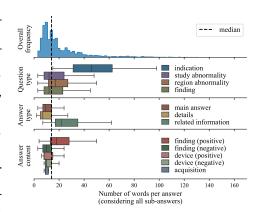


Figure 6: Distribution of answer lengths. We provide nuanced answers with detailed freetext finding descriptions.

# 5 Structured VQA Task

To demonstrate the utility of our dataset, we introduce a structured Visual Question Answering (VQA) task. This task requires models to generate free-text answers accompanied by bounding boxes and tags (e.g., findings, regions). Given a chest X-ray and a free-text question, the model must output such structured answers to respond to the query. We refer to Appendix F for further details.

Sequence Formatting for Structured VQA We implement a proof-of-concept model based on the Llava architecture [37], using Rad-DINO [38] for image encoding and the Llama 3.2 3B [26] language model connected via an MLP projection layer. Our CXR-QBA dataset provides the necessary targets, which we format into sequences using XML-style structures and special tokens to represent tags and bounding boxes (converted to relative coordinates following [18]). We then fine-tune the model for one epoch on 1M QA-pairs (MIMIC-CXR train split).

RadStrucVQA Metric For evaluation, we introduce the *RadStrucVQA* metric, which closely follows the RadFact [18] metric introduced for radiology report generation but is generalized to structured VQA. Like RadFact, we identify whether individual predicted answer parts are entailed with target answer parts and vice-versa, in our casing using Llama 3.1 8B. For entailed pairs, we compute whether they are visually grounded, i.e. whether their bounding boxes are precise enough considering their references, and whether finding and region tags are correctly reported. This is conducted bi-directionally, using either the targets as references for the predictions or vice-versa, resulting in precision or recall scores, respectively. More details can be found in Appendix F.2.

Table 3: Results on the structured VQA task, with 95% confidence intervals (bootstrapping, n=1000). **Left**: Our model trained on this task. **Right**: MAIRA-2 with adapted prompt. Our dataset enables training vision-language models to predict logically correct, visually grounded answers, supplemented by tags that facilitate thorough analysis of the model's predictions.

	Ours	MAIRA-2
Logical Prec.	<b>0.76</b> [0.75, 0.76]	0.25 [0.25, 0.26]
Logical Rec. Grounding Prec.	<b>0.75</b> [0.74, 0.76]	0.64 [0.63, 0.65]
Grounding Prec.	<b>0.87</b> [0.87, 0.88]	0.69 [0.67, 0.71]
Grounding Rec.	<b>0.89</b> [0.88, 0.89]	$\substack{0.12 \\ [0.11,0.12]}$
Finding Prec.	<b>0.68</b> [0.67, 0.69]	-
Finding Rec.	<b>0.66</b> [0.66, 0.67]	-
Finding-pos Prec.	. <b>0.41</b> [0.40, 0.43]	-
Finding-pos Rec.	<b>0.26</b> [0.25, 0.27]	-
Finding-pos Prec. Finding-pos Rec. Region Prec.	<b>0.67</b> [0.66, 0.68]	-
Region Rec.	<b>0.66</b> [0.65, 0.67]	
*Our RadStrucVO	QA implem	entation.

**Results** We evaluate our model on our fine-tuning grade dataset (MIMIC-CXR test split) and compare it to MAIRA-2 [18], a model for grounded report generation trained partially on MIMIC-

CXR, i.e. on the same images as our model. We use the frozen MAIRA-2 model, but adapt its prompt 288 to answer specific questions instead of generating full reports. Results are shown in Tab. 3. Our model 289 achieves high scores in both textual content (logical) and grounding metrics, demonstrating effective 290 training on our dataset. As expected, MAIRA-2 performed lower on all metrics, but achieved 85%291 of our models logical recall, suggesting it captured most relevant information while also including 292 extraneous details (lower precision). This aligns with its training objective of comprehensive reporting 293 294 but also indicates that our dataset's answers do not contradict with MAIRA-2's predictions, further confirming the quality of our dataset. MAIRA-2's grounding precision significantly exceeded its recall, 295 because it was trained to predict bounding boxes only for positive findings. Our model successfully 296 predicts finding and region tags in most cases. However, performance drops when focusing solely on 297 positive findings (finding-pos), indicating potential underprediction possibly due to our training 298 procedure or limitations in the pre-trained components. While further analysis would be required, this may indicate problems that could also lead to flaws in textual answers. Importantly, our datasets detailed tags enable fine-grained analysis of such issues while also enabling potential solutions like data filtering or balancing, making it well-suited for complex training scenarios. 302

# 303 6 Discussion and Conclusion

### 6.1 Use Cases and Impact

304

317

332

Our dataset is particularly well-suited for structured VQA on CXRs (Sec. 5). Additionally, its versatility also supports classical VQA tasks or grounded VQA without structure, while its large size and detailed answers make it a valuable resource for pre-training vision-language models. The accompanying tags further enable filtering and balancing of the dataset to suite specific needs.

Use cases are, however, not limited to VQA tasks. Our fine-grained scene graphs with bounding boxes, textual descriptions, and tags can serve as a versatile data source for various purposes. For instance, they can be leveraged to create customized datasets for grounded report generation or VQA, or even as a direct training source for graph generation models to predict scene graphs on unseen chest X-rays, enabling the creation of even larger datasets. Furthermore, the bounding boxes and tags provided with the scene graphs can be used for longitudinal analysis, including region-level examination. Finally, they can be used to train models for pathology localization or classification, providing fine-grained and long-tail diagnosis targets that are often lacking in existing datasets.

#### 6.2 Limitations

Our dataset was automatically constructed, relying on models and templates instead of human 318 annotations. While this enables the generation of a large number of QA-pairs, it may also introduce 319 potential errors and biases. We apply (automatic) quality assessments to mitigate these risks, but users should still be aware that the dataset may contain inaccuracies and should exercise caution, especially when using it for critical applications. Most importantly, we strictly advise against using this dataset as the sole source for fine-tuning or evaluating models used in clinical practice. 323 Furthermore, our template-based approach may limit the diversity of the dataset and may potentially 324 introduce grammatical errors. However, we partially mitigate these issues by incorporating answers 325 derived from actual report sentences and through our quality assessment measures. Additionally, our 326 approach focuses on individual chest X-ray studies, excluding longitudinal, differential questions, and 327 other (imaging) modalities. Future extensions could build upon this work, generalizing our approach 328 to broader question types and modalities. Finally, our work relies on LLMs for information extraction 329 and quality assessment. While we only use medium to small models, these still require substantial 330 computational resources for dataset creation, particularly compared to template-based methods. 331

# 6.3 Conclusion

We proposed a novel approach to constructing a large-scale CXR VQA dataset using automatic scene graph construction and question-answer generation, resulting in CXR-QBA, a dataset of 42 M QA-pairs. We hope that our dataset will serve as a valuable resource for researchers and practitioners, driving advancements in medical imaging and vision-language understanding.

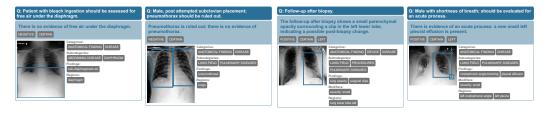
### References

- Z. Chen, M. Varma, J. Xu, et al., A Vision-Language Foundation Model to Enhance Efficiency
   of Chest X-ray Interpretation, 2024. DOI: 10.48550/arXiv.2401.12208.
- P. Müller, G. Kaissis, and D. Rueckert, "ChEX: Interactive Localization and Region Description in Chest X-Rays," in *Computer Vision ECCV 2024*, Springer Nature Switzerland, 2025, pp. 92–111. DOI: 10.1007/978-3-031-72664-4\_6.
- T. Tanida, P. Müller, G. Kaissis, *et al.*, "Interactive and Explainable Region-guided Radiology Report Generation," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 7433–7442. DOI: 10.1109/CVPR52729.2023.00718.
- C. Pellegrini, B. Busam, B. Wiestler, *et al.*, "RaDialog: Large Vision-Language Models for X-Ray Reporting and Dialog-Driven Assistance," in *Proceedings of Machine Learning Research*, 2025.
- T. Tu, S. Azizi, D. Driess, *et al.*, "Towards Generalist Biomedical AI," *NEJM AI*, vol. 1, no. 3,
   AIoa2300138, 2024. DOI: 10.1056/AIoa2300138.
- C. Li, C. Wong, S. Zhang, et al., "LLaVA-Med: Training a Large Language-and-Vision Assistant for Biomedicine in One Day," in NIPS '23: Proceedings of the 37th International Conference on Neural Information Processing Systems, 2023.
- Y. Xie, C. Zhou, L. Gao, et al., MedTrinity-25M: A Large-scale Multimodal Dataset with Multigranular Annotations for Medicine, 2025. DOI: 10.48550/arXiv.2408.02900.
- J. J. Lau, S. Gayen, A. Ben Abacha, *et al.*, "A dataset of clinically generated visual questions and answers about radiology images," *Sci Data*, vol. 5, no. 1, p. 180 251, 2018. DOI: 10.1038/sdata.2018.251.
- B. Liu, L.-M. Zhan, L. Xu, et al., "Slake: A Semantically-Labeled Knowledge-Enhanced Dataset For Medical Visual Question Answering," in 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI), 2021, pp. 1650–1654. DOI: 10.1109/ISBI48211.2021.9434010.
- A. Ben Abacha, S. A. Hasan, V. V. Datla, *et al.*, "VQA-Med: Overview of the medical visual question answering task at ImageCLEF 2019," *Proceedings of CLEF (Conference and Labs of the Evaluation Forum) 2019 Working Notes*, 2019.
- 366 [11] X. Zhang, C. Wu, Z. Zhao, *et al.*, "Development of a large-scale medical visual question-367 answering dataset," *Commun Med*, vol. 4, no. 1, pp. 1–13, 2024. DOI: 10.1038/s43856-024-368 00709-2.
- S. Bae, D. Kyung, J. Ryu, et al., MIMIC-Ext-MIMIC-CXR-VQA: A Complex, Diverse, And Large-Scale Visual Question Answering Dataset for Chest X-ray Images, 2024. DOI: 10. 13026/DEQX-D943.
- X. Hu, L. Gu, K. Kobayashi, *et al.*, "Interpretable medical image Visual Question Answering via multi-modal relationship graph learning," *Medical Image Analysis*, vol. 97, p. 103 279, 2024. DOI: 10.1016/j.media.2024.103279.
- 375 [14] A. Johnson, T. Pollard, R. Mark, *et al.*, *MIMIC-CXR Database*, 2024. DOI: 10.13026/4JQJ-376 JW95.
- 377 [15] A. L. Goldberger, L. A. N. Amaral, L. Glass, *et al.*, "PhysioBank, PhysioToolkit, and PhysioNet," *Circulation*, vol. 101, no. 23, e215–e220, 2000. DOI: 10.1161/01.CIR.101.23. e215.
- X. Hu, L. Gu, K. Kobayashi, et al., Medical-CXR-VQA dataset: A Large-Scale LLM-Enhanced
   Medical Dataset for Visual Question Answering on Chest X-Ray Images, 2025. DOI: 10.
   13026/1PM5-HY02.
- 383 [17] J. Wu, N. Agu, I. Lourentzou, et al., Chest ImaGenome Dataset, 2021. DOI: 10.13026/WV01-384 Y230.
- 385 [18] S. Bannur, K. Bouzid, D. C. Castro, et al., MAIRA-2: Grounded Radiology Report Generation, 386 2024. DOI: 10.48550/arXiv.2406.04449.
- J. T. Wu, N. N. Agu, I. Lourentzou, et al., Chest ImaGenome Dataset for Clinical Reasoning, 2021. DOI: 10.48550/arXiv.2108.00316.

- S. Jain, A. Agrawal, A. Saporta, et al., "RadGraph: Extracting Clinical Entities and Relations from Radiology Reports," Proceedings of the Neural Information Processing
  Systems Track on Datasets and Benchmarks, vol. 1, 2021. [Online]. Available: https:
  //datasets benchmarks proceedings . neurips . cc / paper / 2021 / hash /
  c8ffe9a587b126f152ed3d89a146b445-Abstract-round1.html.
- A. E. W. Johnson, T. J. Pollard, S. J. Berkowitz, *et al.*, "MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports," *Sci Data*, vol. 6, no. 1, p. 317, 2019. DOI: 10.1038/s41597-019-0322-0.
- C. Seibold, S. ReiSS, S. Sarfraz, *et al.*, "Detailed Annotations of Chest X-Rays via CT Projection for Report Understanding," arXiv, 2022. DOI: 10.48550/arXiv.2210.03416.
- C. Seibold, A. Jaus, M. A. Fink, et al., Accurate Fine-Grained Segmentation of Human Anatomy in Radiographs via Volumetric Pseudo-Labeling, 2023. DOI: 10.48550/arXiv.2306.03934.
- 401 [24] A. Johnson, M. Lungren, Y. Peng, et al., MIMIC-CXR-JPG chest radiographs with structured labels, 2024. DOI: 10.13026/JSN5-T979.
- A. E. W. Johnson, T. J. Pollard, N. R. Greenbaum, et al., MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs, 2019. DOI: 10.48550/arXiv.1901.07042.
- 405 [26] A. Grattafiori, A. Dubey, A. Jauhri, et al., The Llama 3 Herd of Models, 2024. DOI: 10.48550/ 406 arXiv.2407.21783.
- 407 [27] A. Bustos, A. Pertusa, J.-M. Salinas, *et al.*, "PadChest: A large chest x-ray image dataset with multi-label annotated reports," *Medical Image Analysis*, vol. 66, p. 101797, 2020. DOI: 10.1016/j.media.2020.101797.
- 410 [28] SNOMED International, SNOMED CT, 2023. [Online]. Available: https://www.snomed.
- F. Remy, K. Demuynck, and T. Demeester, "BioLORD-2023: Semantic textual representations fusing large language models and clinical knowledge graph insights," *Journal of the American Medical Informatics Association*, vol. 31, no. 9, pp. 1844–1855, 2024. DOI: 10.1093/jamia/ocae029.
- 416 [30] G. Holste, M. Lin, S. Wang, et al., CXR-LT: Multi-Label Long-Tailed Classification on Chest X-Rays, 2025. DOI: 10.13026/RYJ9-X506.
- B. Boecking, N. Usuyama, S. Bannur, et al., MS-CXR: Making the Most of Text Semantics to Improve Biomedical Vision-Language Processing, 2024. DOI: 10.13026/9G2Z-JG61.
- 420 [32] R. Bigolin Lanfredi, M. Zhang, W. Auffermann, et al., REFLACX: Reports and eye-tracking data for localization of abnormalities in chest x-rays, 2021. DOI: 10.13026/E0DJ-8498.
- J. Irvin, P. Rajpurkar, M. Ko, et al., "CheXpert: A Large Chest Radiograph Dataset with Uncertainty Labels and Expert Comparison," Proceedings of the AAAI Conference on Artificial Intelligence, vol. 33, no. 01, pp. 590–597, 2019. DOI: 10.1609/aaai.v33i01.3301590.
- G. Holste, Y. Zhou, S. Wang, *et al.*, "Towards long-tailed, multi-label disease classification from chest X-ray: Overview of the CXR-LT challenge," *Medical Image Analysis*, vol. 97, p. 103 224, 2024. DOI: 10.1016/j.media.2024.103224.
- B. Boecking, N. Usuyama, S. Bannur, *et al.*, "Making the Most of Text Semantics to Improve Biomedical VisionLanguage Processing," in *Computer Vision ECCV 2022*, vol. 13696, Springer Nature Switzerland, 2022, pp. 1–21. DOI: 10.1007/978-3-031-20059-5\_1.
- R. Bigolin Lanfredi, M. Zhang, W. F. Auffermann, *et al.*, "REFLACX, a dataset of reports and eye-tracking data for localization of abnormalities in chest x-rays," *Sci Data*, vol. 9, no. 1, p. 350, 2022. DOI: 10.1038/s41597-022-01441-z.
- H. Liu, C. Li, Q. Wu, et al., "Visual Instruction Tuning," Advances in Neural Information Processing Systems, vol. 36, pp. 34892-34916, 2023. [Online]. Available: https://papers.nips.cc/paper\_files/paper/2023/hash/6dcf277ea32ce3288914faf369fe6de0-Abstract-Conference.html.
- F. Pérez-García, H. Sharma, S. Bond-Taylor, *et al.*, "Exploring scalable medical image encoders beyond text supervision," *Nat Mach Intell*, vol. 7, no. 1, pp. 119–130, 2025. DOI: 10.1038/s42256-024-00965-w.
- P. Qi, Y. Zhang, Y. Zhang, et al., "Stanza: A Python Natural Language Processing Toolkit for Many Human Languages," in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, 2020. [Online]. Available: https://nlp.stanford.edu/pubs/qi2020stanza.pdf.

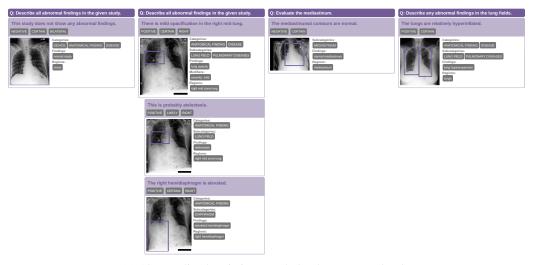
- [40] J. Lin, J. Tang, H. Tang, et al., "AWQ: Activation-aware Weight Quantization for On-Device LLM Compression and Acceleration," Proceedings of Machine Learning and Systems, vol. 6, pp. 87-100, 2024. [Online]. Available: https://proceedings.mlsys.org/paper\_files/paper/2024/hash/42a452cbafa9dd64e9ba4aa95cc1ef21-Abstract-Conference.html.
- W. Kwon, Z. Li, S. Zhuang, et al., "Efficient Memory Management for Large Language Model Serving with PagedAttention," in *Proceedings of the 29th Symposium on Operating Systems* Principles, ser. SOSP '23, Association for Computing Machinery, 2023, pp. 611–626. DOI: 10.1145/3600006.3613165.
- E. J. Hu, Y. Shen, P. Wallis, *et al.*, "LoRA: Low-Rank Adaptation of Large Language Models," 2021. [Online]. Available: https://openreview.net/forum?id=nZeVKeeFYf9.

# 456 A Example QA-Pairs from our Dataset

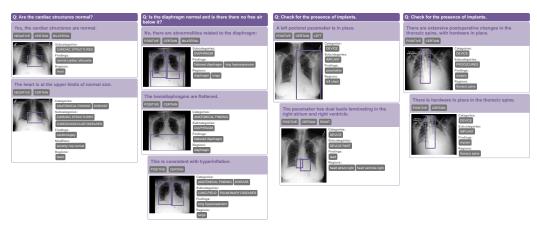


(a) Negative.

Figure 7: Examples of **indication** questions. Questions are based on the paraphrased INDICATION section while each main answer is generated based on the indication node from the scene graph (using information from the FINDINGS and IMPRESION sections).



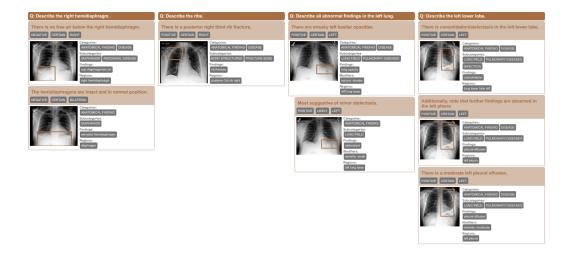
(a) Abnormality descriptions (study-level or category-level).



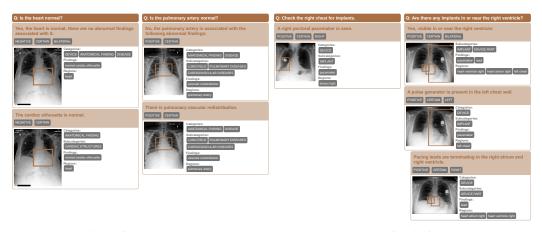
- (b) Abnormality assessment (category-level).
- (c) Device descriptions.

(b) Positive.

Figure 8: Examples of **study abnormality** questions. Questions are based on one of 13 templates. Answers may consist of several answer parts, where each describes an individual aspect (about the overall study or a finding category). Individual answer parts are constructed based on observation nodes, filtered based on finding categories relevant to the question, where individual answer parts may be organized hierarchically (indicated by indentations) based on parent-child edges in the scene graph. Additionally assessment answers (b) start with a template-based yes/no answer.



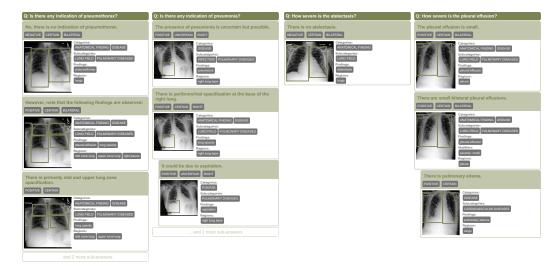
(a) Region description.



(b) Region assessment.

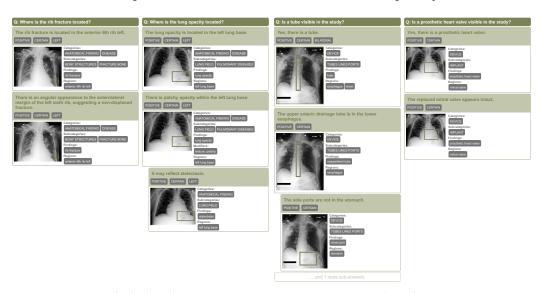
(c) Region devices.

Figure 9: Examples of **region abnormality** questions. Questions are based on one of 6 templates. Answers may consist of several answer parts, where each describes an individual aspect (about the region). Individual answer parts are constructed based on observation nodes relevant to the region, where individual answer parts may be organized hierarchically (indicated by indentations) based on parent-child edges in the scene graph. Additionally assessment answers (b) start with a template-based yes/no answer. Some templates also ask specifically about devices in the region (c).



(a) Finding assessment.

(b) Finding description.



(c) Finding location.

(d) Device.

Figure 10: Examples of **finding** questions. Questions are based on one of 7 templates. Answers start with a template-based answer part to identify the finding presence (a), provide a severity summary (b), describe the location (c), or presence of a device (d). Additional details may be provided in answer parts based on observation nodes relevant to the finding, where individual answer parts may be organized hierarchically (indicated by indentations) based on parent-child edges in the scene graph.

# 457 B Finding- and Region-Distribution in QA-Samples

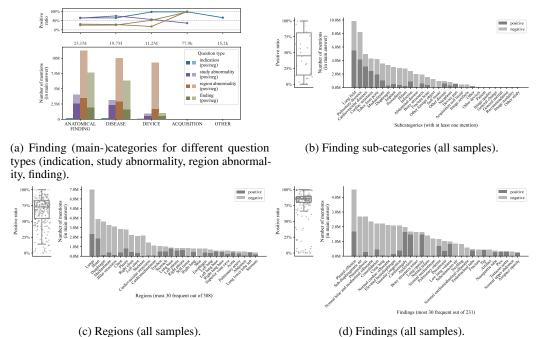


Figure 11: Tags (finding main- and sub-categories, regions, findings) mentioned in answers. We show their positive ratios (**top/left**), i.e. how often they are mentioned in positive versus in negative answer parts and plot the number of positive and negative mentions of the most frequent tags (**bottom/right**).

# 458 C Evaluation Details

# 59 C.1 QA Evaluation: Comparison of LLM-Raters

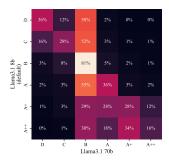


Figure 12: Confusion matrix comparing the assigned quality ratings between using Llama3.1 8b (default) and Llama3.1 70b as an LLM-judge (see Secs. 3.3 and 4.2). In most cases, ratings differ only slightly. Most importantly, low-quality samples (as rated by Llama3.1 70b) are almost never assigned to fine-tuning grades (A or higher) by Llama3.1 8b. We thus decided to use Llama3.1 8b as our default rater, as it is much more computationally efficient.

# 460 C.2 Scene Graph Evaluation: Finding Tags

Table 4: Evaluation of finding tags against the 13 CheXpert (CXP) classes from the MIMIC-CXR-JPG test set (Sec. 4.1). We show finding-level scores, macro-averages over subsets and the micro-average, with 95% confidence intervals (bootstrapping, n = 1000).

				MIMIC-CXR	-JPG [24] T	est		
	[Pre	ecision]	[R	ecall]	[	F1]	[N	ACC]
	Ours	Chest ImaG.						
Findings in CXP-5, CX	XP-7, and C	CXP-13						
Atelectasis	<b>0.82</b> [0.78, 0.87]	0.78 [0.73, 0.83]	<b>0.99</b> [0.97, 1.0]	<b>0.99</b> [0.98, 1.0]	<b>0.9</b> [0.87, 0.93]	0.88 [0.84, 0.9]	<b>0.84</b> [0.8, 0.88]	0.81 [0.77, 0.85]
Cardiomegaly	0.64 [0.58, 0.7]	<b>0.67</b> [0.61, 0.74]	<b>0.85</b> [0.78, 0.9]	0.82 [0.75, 0.87]	0.73 [0.67, 0.78]	<b>0.74</b> [0.68, 0.79]	0.61 [0.54, 0.68]	<b>0.63</b> [0.56, 0.69]
Consolidation	<b>0.83</b> [0.73, 0.91]	0.77	0.87 [0.79, 0.94]	<b>0.93</b> [0.86, 0.99]	<b>0.85</b> [0.78, 0.91]	0.84	<b>0.83</b> [0.75, 0.9]	<b>0.83</b> [0.76, 0.89]
Edema	<b>0.94</b> [0.9, 0.98]	0.9 [0.85, 0.95]	<b>0.8</b> [0.73, 0.86]	<b>0.8</b> [0.74, 0.86]	<b>0.86</b> [0.82, 0.9]	0.85 [0.8, 0.89]	<b>0.83</b> [0.77, 0.87]	0.8 [0.74, 0.85]
Pleural Effusion	<b>0.9</b> [0.86, 0.93]	0.86 [0.82, 0.9]	<b>0.98</b> [0.96, 1.0]	0.97 [0.94, 0.99]	<b>0.94</b> [0.92, 0.96]	<b>0.91</b> [0.89, 0.94]	<b>0.89</b> [0.85, 0.92]	0.85 [0.8, 0.89]
Findings in CXP-7 and	1 CXP-13							
Pneumonia	<b>0.92</b> [0.87, 0.96]	0.89 [0.84, 0.94]	0.94 [0.89, 0.97]	<b>0.95</b> [0.91, 0.98]	<b>0.93</b> [0.89, 0.96]	0.92 [0.89, 0.95]	<b>0.91</b> [0.87, 0.95]	0.9 [0.86, 0.94]
Pneumothorax	0.78 [0.64, 0.89]	<b>0.79</b> [0.66, 0.91]	0.84 [0.71, 0.95]	<b>0.89</b> [0.78, 0.98]	0.8 [0.69, 0.89]	<b>0.84</b> [0.74, 0.92]	0.79 [0.68, 0.88]	<b>0.83</b> [0.72, 0.91]
Findings in CXP-13								
Enlarged Cardiom.	0.51 [0.37, 0.65]	<b>0.61</b> [0.41, 0.8]	<b>0.39</b> [0.29, 0.51]	0.23 [0.13, 0.33]	<b>0.44</b> [0.34, 0.55]	0.33	<b>0.39</b> [0.28, 0.51]	0.33 [0.19, 0.45]
Lung Lesion	0.17 [0.12, 0.22]	<b>0.68</b> [0.56, 0.79]	0.81 [0.69, 0.91]	<b>0.87</b> [0.76, 0.95]	0.28 [0.21, 0.35]	<b>0.76</b> [0.66, 0.84]	0.25 [0.17, 0.32]	<b>0.74</b> [0.64, 0.83]
Lung Opacity	<b>0.62</b> [0.56, 0.69]	0.28 [0.24, 0.31]	0.83 [0.77, 0.89]	1.0 [1.0, 1.0]	<b>0.71</b> [0.65, 0.76]	0.43 [0.39, 0.48]	<b>0.61</b> [0.54, 0.68]	0.2 [0.17, 0.23]
Pleural Other	<b>0.54</b> [0.36, 0.71]	0.3	0.87 [0.7, 1.0]	<b>0.92</b> [0.76, 1.0]	<b>0.67</b> [0.49, 0.8]	0.45 [0.31, 0.58]	<b>0.67</b> [0.51, 0.79]	0.5 [0.37, 0.6]
Fracture	<b>0.67</b> [0.54, 0.79]	0.6 [0.47, 0.73]	<b>0.92</b> [0.82, 1.0]	0.82 [0.68, 0.93]	<b>0.77</b> [0.67, 0.86]	0.69 [0.58, 0.79]	<b>0.77</b> [0.67, 0.85]	0.68 [0.56, 0.78]
Support Devices	0.61 [0.56, 0.66]	<b>0.62</b> [0.56, 0.67]	<b>0.98</b> [0.96, 1.0]	0.83 [0.77, 0.88]	<b>0.75</b> [0.71, 0.79]	0.71 [0.66, 0.75]	<b>0.63</b> [0.59, 0.68]	0.54 [0.48, 0.61]
Macro-averages								
CheXpert-5 (CXP-5)	<b>0.83</b> [0.79, 0.85]	0.8 [0.77, 0.83]	<b>0.9</b> [0.87, 0.92]	<b>0.9</b> [0.88, 0.92]	<b>0.85</b> [0.83, 0.87]	0.84 [0.82, 0.86]	<b>0.8</b> [0.77, 0.82]	0.78 [0.75, 0.81]
CheXpert-7 (CXP-7)	<b>0.83</b> [0.8, 0.86]	0.81 [0.78, 0.84]	0.89 [0.87, 0.92]	<b>0.91</b> [0.88, 0.93]	<b>0.86</b> [0.83, 0.88]	0.85 [0.83, 0.87]	<b>0.81</b> [0.79, 0.84]	0.8
CheXpert-13 (CXP-13)		0.67 [0.65, 0.7]	<b>0.85</b> [0.83, 0.87]	<b>0.85</b> [0.82, 0.87]	<b>0.74</b> [0.72, 0.76]	0.72 [0.7, 0.74]	<b>0.69</b> [0.67, 0.71]	0.66 [0.64, 0.69]
Micro	<b>0.68</b> [0.66, 0.7]	0.63 [0.61, 0.65]	<b>0.89</b> [0.87, 0.9]	0.88 [0.86, 0.89]	<b>0.77</b> [0.75, 0.79]	0.73 [0.72, 0.75]	<b>0.71</b> [0.69, 0.73]	0.67 [0.65, 0.68]

Table 5: Evaluation of finding tags against the 13 CXP and 12 long-tail (LT) classes from the CXR-LT 2024 gold standard dataset (Sec. 4.1). We show finding-level scores, macro-averages over different subsets and the micro-average, with 95% confidence intervals (bootstrapping, n=1000).

	CXR-LT 2024 [30] Gold							
	[Pre	cision]	[Re	ecall]	[	F1]	[N	ICC]
	Ours	Chest ImaG.	Ours	Chest ImaG.	Ours	Chest ImaG.	Ours	Chest ImaG.
Findings in CXP-5, CXP-7	7, CXP-13,							
Atelectasis	0.55 [0.47, 0.62]	<b>0.56</b> [0.49, 0.61]	0.82 [0.75, 0.88]	<b>0.99</b> [0.97, 1.0]	0.66 [0.59, 0.71]	<b>0.71</b> [0.65, 0.76]	0.48 [0.4, 0.57]	<b>0.59</b> [0.54, 0.65]
Cardiomegaly	0.82 [0.76, 0.88]	<b>0.85</b> [0.79, 0.91]	<b>0.85</b> [0.79, 0.91]	0.8 [0.73, 0.86]	<b>0.84</b> [0.79, 0.88]	0.82 [0.77, 0.87]	<b>0.73</b> [0.66, 0.79]	0.72 [0.65, 0.79]
Consolidation	<b>0.82</b> [0.73, 0.91]	0.74 [0.63, 0.83]	0.86 [0.76, 0.93]	<b>0.89</b> [0.8, 0.95]	<b>0.84</b> [0.77, 0.9]	0.8 [0.73, 0.87]	<b>0.8</b> [0.72, 0.87]	0.76 [0.67, 0.83]
Edema	<b>0.73</b> [0.63, 0.83]	0.69 [0.6, 0.79]	0.64 [0.55, 0.74]	<b>0.71</b> [0.61, 0.8]	0.68 [0.6, 0.76]	<b>0.7</b> [0.62, 0.77]	0.59 [0.5, 0.69]	<b>0.6</b> [0.5, 0.69]
Pleural Effusion	<b>0.82</b> [0.76, 0.87]	0.78 [0.73, 0.84]	0.93 [0.9, 0.97]	<b>0.97</b> [0.94, 0.99]	<b>0.87</b> [0.83, 0.91]	<b>0.87</b> [0.83, 0.9]	<b>0.77</b> [0.7, 0.82]	0.76 [0.7, 0.82]
Findings in CXP-7, CXP-	13, and CX	R-LT						
Pneumonia	<b>0.38</b> [0.19, 0.58]	0.13 [0.07, 0.21]	0.45 [0.25, 0.67]	<b>0.76</b> [0.55, 0.94]	<b>0.41</b> [0.22, 0.58]	0.23 [0.13, 0.33]	<b>0.38</b> [0.19, 0.55]	0.24 [0.13, 0.34]
Pneumothorax	<b>0.85</b> [0.73, 0.95]	0.8 [0.7, 0.9]	0.85 [0.72, 0.94]	<b>0.96</b> [0.89, 1.0]	0.85 [0.75, 0.92]	<b>0.87</b> [0.8, 0.94]	0.83 [0.72, 0.91]	<b>0.86</b> [0.78, 0.93]
Findings in CXP-13 and C		[,]	2 , ,	( , , , , , , , , , , , , , , , , , , ,	,	[, ]	2,,	,
Enlarged Cardiom.	0.78 [0.57, 0.95]	<b>1.0</b> [1.0, 1.0]	<b>0.13</b> [0.07, 0.2]	0.1 [0.05, 0.17]	<b>0.22</b> [0.13, 0.32]	0.18 [0.1, 0.29]	0.24 [0.13, 0.35]	<b>0.27</b> [0.19, 0.36]
Lung Lesion	0.01	<b>0.05</b> [0.0, 0.12]	0.5 [0.0, 1.0]	<b>0.75</b> [0.0, 1.0]	0.02 [0.0, 0.06]	<b>0.1</b> [0.0, 0.22]	0.0 [-0.09, 0.1]	<b>0.18</b> [-0.02, 0.3]
Lung Opacity	<b>0.92</b> [0.87, 0.96]	0.54 [0.49, 0.59]	0.77 [0.72, 0.84]	1.0 [1.0, 1.0]	<b>0.84</b> [0.8, 0.88]	0.7 [0.66, 0.74]	<b>0.73</b> [0.66, 0.79]	0.35 [0.29, 0.4]
Pleural Other	0.11	0.23	0.19	1.0	0.14	0.38	0.1	0.45
Fracture	[0.0, 0.24] <b>0.89</b>	[0.13, 0.34] 0.84	[0.0, 0.43] <b>0.91</b>	[1.0, 1.0] 0.82	[0.0, 0.29] <b>0.9</b>	[0.23, 0.51] 0.83	[-0.05, 0.26] <b>0.89</b>	[0.34, 0.55] 0.81
Support Devices	[0.78, 0.98] <b>0.92</b> [0.88, 0.96]	[0.72, 0.94] <b>0.93</b> [0.9, 0.97]	[0.81, 0.98] <b>0.93</b> [0.89, 0.96]	[0.69, 0.93] <b>0.83</b> [0.77, 0.88]	[0.82, 0.95] <b>0.92</b> [0.9, 0.95]	[0.73, 0.9] 0.88 [0.84, 0.91]	[0.8, 0.95] <b>0.83</b> [0.77, 0.88]	[0.69, 0.89] 0.75 [0.68, 0.82]
Findings in LT-only, and G		[0.5, 0.57]	[0.09, 0.90]	[0.77, 0.00]	[0.9, 0.93]	[0.04, 0.71]	[0.77, 0.88]	[0.00, 0.02]
Calcification of the Aorta	0.95	0.95	0.43	0.93	0.6	0.94	0.61	0.94
Emphysema	[0.83, 1.0] <b>0.58</b>	[0.88, 1.0] <b>0.54</b>	[0.28, 0.58] <b>0.81</b>	[0.85, 1.0] <b>0.81</b>	[0.43, 0.73] <b>0.68</b>	[0.89, 0.99] <b>0.65</b>	[0.48, 0.73] <b>0.66</b>	[0.87, 0.99] <b>0.63</b>
Fibrosis	[0.41, 0.74] <b>0.27</b>	[0.38, 0.69] <b>0.0</b>	[0.63, 0.95] <b>0.52</b>	[0.63, 0.95] <b>0.0</b>	[0.52, 0.8] <b>0.36</b>	[0.49, 0.77] <b>0.0</b>	[0.5, 0.79] <b>0.33</b>	[0.47, 0.75]
Hernia	[0.15, 0.43] <b>1.0</b>	[0.0, 0.0] <b>0.86</b>	[0.31, 0.74] <b>0.9</b>	[0.0, 0.0] <b>0.9</b>	[0.21, 0.52] <b>0.95</b>	[0.0, 0.0] <b>0.88</b>	[0.17, 0.5] <b>0.95</b>	[0.0, 0.0] 0.87
Infiltration	[1.0, 1.0]	[0.68, 1.0] <b>0.38</b>	[0.73, 1.0] 0.33	[0.73, 1.0] <b>0.55</b>	[0.85, 1.0] 0.21	[0.74, 0.97] <b>0.44</b>	[0.85, 1.0] 0.19	[0.73, 0.97] <b>0.44</b>
Mass	[0.03, 0.3] <b>0.54</b>	[0.16, 0.67] 0.28	[0.08, 0.7]	[0.21, 0.88] <b>0.89</b>	[0.05, 0.39] <b>0.64</b>	[0.18, 0.69] <b>0.42</b>	[0.01, 0.39] <b>0.63</b>	[0.17, 0.69] 0.46
Nodule	[0.34, 0.74] <b>0.92</b>	[0.17, 0.4]	[0.56, 0.94] 0.74	[0.71, 1.0] <b>0.91</b>	[0.44, 0.78] <b>0.82</b>	[0.28, 0.56] 0.64	[0.45, 0.77] <b>0.81</b>	[0.32, 0.58]
	[0.79, 1.0]	0.5 [0.36, 0.63]	[0.57, 0.88]	[0.78, 1.0]	[0.68, 0.91]	[0.51, 0.75]	[0.68, 0.91]	[0.51, 0.74]
Pleural Thickening	<b>0.78</b> [0.62, 0.92]	0.33 [0.22, 0.45]	1.0 [1.0, 1.0]	1.0 [1.0, 1.0]	<b>0.88</b> [0.77, 0.96]	0.5 [0.36, 0.62]	<b>0.88</b> [0.78, 0.96]	0.54 [0.43, 0.64]
Pneumomediastinum	<b>0.94</b> [0.83, 1.0]	0.88 [0.73, 0.97]	<b>0.84</b> [0.7, 0.96]	<b>0.84</b> [0.7, 0.96]	<b>0.89</b> [0.78, 0.96]	0.86 [0.73, 0.94]	<b>0.88</b> [0.77, 0.95]	0.85 [0.72, 0.93]
Pneumoperitoneum	<b>0.88</b> [0.72, 1.0]	0.82 [0.65, 0.96]	0.96 [0.85, 1.0]	<b>1.0</b> [1.0, 1.0]	<b>0.92</b> [0.81, 1.0]	<b>0.9</b> [0.79, 0.98]	<b>0.91</b> [0.8, 1.0]	0.9 [0.8, 0.98]
Subcutaneous Emphysema	<b>0.97</b> [0.9, 1.0]	0.0 [0.0, 0.0]	<b>0.8</b> [0.68, 0.92]	0.0 [0.0, 0.0]	<b>0.88</b> [0.79, 0.95]	0.0 [0.0, 0.0]	<b>0.87</b> [0.78, 0.94]	$\begin{bmatrix} 0.0 \\ [0.0, 0.0] \end{bmatrix}$
Tortuous Aorta	<b>0.83</b> [0.69, 0.95]	0.79 [0.64, 0.91]	0.88 [0.75, 0.97]	<b>0.94</b> [0.84, 1.0]	0.85 [0.75, 0.93]	<b>0.86</b> [0.75, 0.94]	0.84 [0.73, 0.93]	<b>0.85</b> [0.74, 0.93]
Macro-averages								
CheXpert-5 (CXP-5)	<b>0.75</b> [0.71, 0.78]	0.72 [0.69, 0.76]	0.82 [0.79, 0.85]	<b>0.87</b> [0.84, 0.9]	<b>0.78</b> [0.75, 0.8]	<b>0.78</b> [0.75, 0.81]	0.67 [0.64, 0.71]	<b>0.69</b> [0.65, 0.72]
CheXpert-7 (CXP-7)	<b>0.71</b> [0.66, 0.75]	0.65 [0.62, 0.68]	0.77 [0.73, 0.81]	<b>0.87</b> [0.83, 0.9]	<b>0.73</b> [0.7, 0.77]	0.71 [0.69, 0.74]	<b>0.65</b> [0.61, 0.69]	<b>0.65</b> [0.61, 0.68]
CheXpert-13 (CXP-13)	<b>0.66</b> [0.63, 0.69]	0.63 [0.61, 0.65]	0.68 [0.63, 0.73]	<b>0.81</b> [0.75, 0.85]	<b>0.63</b> [0.6, 0.65]	0.62 [0.6, 0.64]	<b>0.57</b> [0.54, 0.6]	0.56 [0.54, 0.59]
LT-only	<b>0.73</b> [0.69, 0.77]	0.53 [0.48, 0.57]	<b>0.75</b> [0.7, 0.8]	0.73 [0.69, 0.77]	<b>0.72</b> [0.68, 0.75]	0.59 [0.55, 0.62]	<b>0.71</b> [0.67, 0.74]	0.59 [0.55, 0.63]
CXR-LT	<b>0.7</b> [0.67, 0.72]	0.58 [0.56, 0.6]	0.71 [0.68, 0.75]	<b>0.77</b> [0.74, 0.8]	<b>0.67</b> [0.65, 0.69]	0.61 [0.58, 0.62]	<b>0.64</b> [0.61, 0.66]	0.58 [0.55, 0.6]
Micro	<b>0.69</b> [0.67, 0.71]	0.62	0.76 [0.74, 0.78]	<b>0.8</b> [0.78, 0.82]	<b>0.72</b> [0.71, 0.74]	0.7 [0.68, 0.72]	<b>0.67</b> [0.65, 0.69]	0.64

### 461 C.3 Scene Graph Evaluation: Finding Boxes

Table 6: Evaluation of finding bounding boxes against 6 finding classes from MS-CXR (see Sec. 4.1). We show finding-level scores, macro-averages over different subsets and the micro-average, with 95% confidence intervals (bootstrapping, n=1000). We excluded 2 of the 8 finding classes, because there are no samples that have positive annotations from MS-CXR, Chest ImaGenome and our dataset.

		MS-CXR [31]							
	[IoI]	J@30]	[IoI]	P@30]	[IoT@30]				
	Ours	Chest ImaG.	Ours	Chest ImaG.	Ours	Chest ImaG.			
Atelectasis	<b>0.28</b> [0.12, 0.42]	0.1	<b>0.5</b> [0.34, 0.66]	0.14 [0.05, 0.28]	0.83 [0.71, 0.93]	<b>0.85</b> [0.74, 0.94]			
Cardiomegaly	0.96 [0.93, 0.98]	<b>0.97</b> [0.95, 0.99]	1.0 [1.0, 1.0]	1.0 [1.0, 1.0]	0.96 [0.93, 0.98]	<b>0.99</b> [0.98, 1.0]			
Consolidation	<b>0.31</b> [0.19, 0.45]	0.2	<b>0.41</b> [0.29, 0.54]	0.24 [0.12, 0.35]	0.91 [0.81, 0.98]	<b>0.98</b> [0.93, 1.0]			
Edema	<b>0.52</b> [0.32, 0.71]	<b>0.52</b> [0.32, 0.71]	<b>0.52</b> [0.32, 0.71]	<b>0.52</b> [0.32, 0.71]	1.0 [1.0, 1.0]	1.0 [1.0, 1.0]			
Pneumonia	<b>0.48</b> [0.41, 0.57]	0.28 [0.21, 0.35]	<b>0.58</b> [0.5, 0.66]	0.34	0.93 [0.88, 0.97]	1.0 [1.0, 1.0]			
Pneumothorax	0.14 [0.1, 0.18]	<b>0.15</b> [0.1, 0.2]	0.14 [0.1, 0.19]	<b>0.15</b> [0.11, 0.2]	0.96 [0.93, 0.98]	<b>0.98</b> [0.96, 1.0]			
Macro	<b>0.45</b> [0.4, 0.5]	0.37 [0.33, 0.41]	<b>0.53</b> [0.48, 0.58]	0.4 [0.36, 0.44]	0.93 [0.9, 0.95]	<b>0.97</b> [0.95, 0.98]			
Micro	<b>0.51</b> [0.47, 0.54]	0.45 [0.42, 0.49]	<b>0.56</b> [0.52, 0.6]	0.48 [0.45, 0.52]	0.94 [0.92, 0.96]	<b>0.98</b> [0.97, 0.99]			

Table 7: Evaluation of finding bounding boxes against 18 finding classes from REFLACX (see Sec. 4.1). We show finding-level scores, macro-averages over different subsets and the micro-average, with 95% confidence intervals (bootstrapping, n=1000). Note that we excluded 11 of the 29 finding classes, because there are no samples that have positive annotations from REFLACX, Chest ImaGenome and our dataset.

	REFLACX [32] all phases						
	Jol]	J@30]	[Iol	P@30]	[Io	Γ@30]	
	Ours	Chest ImaG.	Ours	Chest ImaG.	Ours	Chest ImaG.	
Abnormal mediastinal contour	0.08	<b>0.25</b> [0.0, 0.57]	0.08 [0.0, 0.31]	<b>0.25</b> [0.0, 0.57]	<b>1.0</b> [1.0, 1.0]	<b>1.0</b> [1.0, 1.0]	
Acute fracture	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	1.0 [1.0, 1.0]	0.0 [0.0, 0.0]	
Atelectasis	0.29	0.15	0.47	0.2	0.76	0.93	
Consolidation	[0.26, 0.33] <b>0.39</b>	[0.12, 0.17] 0.27	[0.44, 0.51] <b>0.51</b>	[0.17, 0.23] 0.34	[0.73, 0.78] 0.8	[0.91, 0.94] <b>0.95</b>	
Emphysema	[0.33, 0.45] <b>1.0</b>	[0.22, 0.32] <b>1.0</b>	[0.45, 0.57] <b>1.0</b>	[0.28, 0.4] <b>1.0</b>	[0.74, 0.85] <b>1.0</b>	[0.92, 0.97] <b>1.0</b>	
Enlarged cardiac silhouette	[1.0, 1.0] <b>0.96</b>	[1.0, 1.0] <b>0.96</b>	[1.0, 1.0] <b>1.0</b>	[1.0, 1.0] <b>0.99</b>	[1.0, 1.0] <b>0.96</b>	[1.0, 1.0] <b>0.98</b>	
Enlarged hilum	[0.94, 0.97] <b>0.5</b>	[0.95, 0.97] <b>0.5</b>	[0.99, 1.0] 0.5	[0.99, 1.0] <b>0.8</b>	[0.95, 0.97] <b>0.5</b>	[0.97, 0.99] <b>0.5</b>	
C	[0.0, 1.0]	[0.0, 1.0]	[0.0, 1.0]	[0.23, 1.0]	[0.0, 1.0]	[0.0, 1.0]	
Fracture	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	
Groundglass opacity	0.28	0.31	0.48	0.38	0.77	0.96	
Hiatal hernia	0.19	0.4	[0.4, 0.54] <b>0.27</b>	[0.32, 0.45] <b>0.4</b>	[0.71, 0.83] <b>0.94</b>	[0.93, 0.99] <b>1.0</b>	
III:-b-1	[0.0, 0.43]	[0.17, 0.67]	[0.06, 0.5]	[0.17, 0.67]	[0.78, 1.0]	[1.0, 1.0]	
High lung volume / emphysema	0.48 [0.25, 0.7]	<b>0.58</b> [0.35, 0.79]	<b>0.58</b> [0.35, 0.79]	<b>0.58</b> [0.35, 0.79]	0.9 [0.75, 1.0]	<b>1.0</b> [1.0, 1.0]	
Interstitial lung disease	0.5 [0.0, 1.0]	<b>0.8</b> [0.0, 1.0]	<b>0.8</b> [0.0, 1.0]	<b>0.8</b> [0.0, 1.0]	0.8 [0.12, 1.0]	1.0 [1.0, 1.0]	
Lung nodule or mass	0.18	0.09	0.21	0.09	0.89	0.91	
Pleural abnormality	[0.08, 0.31] <b>0.16</b>	[0.02, 0.2] 0.14	[0.1, 0.35] <b>0.2</b>	[0.02, 0.2] 0.17	[0.77, 0.97] <b>0.91</b>	[0.8, 0.98] <b>0.92</b>	
•	[0.13, 0.19]	[0.11, 0.17]	[0.17, 0.23]	[0.14, 0.2]	[0.88, 0.93]	[0.9, 0.94]	
Pleural effusion	<b>0.42</b> [0.2, 0.65]	0.37 [0.17, 0.6]	<b>0.53</b> [0.3, 0.75]	<b>0.53</b> [0.31, 0.75]	1.0 [1.0, 1.0]	0.95 [0.82, 1.0]	
Pleural thickening	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	<b>0.0</b> [0.0, 0.0]	1.0 [1.0, 1.0]	1.0 [1.0, 1.0]	
Pneumothorax	0.04	0.13	0.04	0.13	0.9	0.96	
Pulmonary edema	[0.01, 0.08] <b>0.51</b>	[0.07, 0.21] <b>0.58</b>	[0.01, 0.08] <b>0.55</b>	[0.07, 0.21] <b>0.58</b>	[0.84, 0.96] 0.95	[0.92, 0.99] <b>1.0</b>	
	[0.45, 0.56]	[0.53, 0.63]	[0.5, 0.6]	[0.53, 0.63]	[0.93, 0.97]	[1.0, 1.0]	
Macro	0.34 [0.27, 0.42]	<b>0.37</b> [0.3, 0.45]	<b>0.41</b> [0.34, 0.49]	<b>0.41</b> [0.33, 0.49]	0.84 [0.78, 0.91]	<b>0.87</b> [0.8, 0.95]	
Micro	<b>0.45</b> [0.44, 0.47]	0.42 [0.4, 0.43]	<b>0.54</b> [0.53, 0.56]	0.46 [0.44, 0.47]	0.87 [0.86, 0.88]	<b>0.95</b> [0.94, 0.96]	

### D Dataset Structure

### D.1 Scene Graph Structure

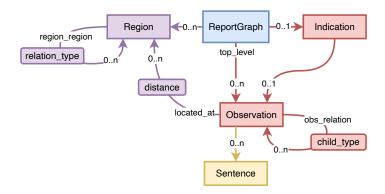


Figure 13: Scene graph structure overview.

Sentence Nodes Sentence nodes are directly associated with raw sentences in the report, i.e. there is exactly one sentence node per identified report sentence. They contain the following attributes:

- sent\_id: Identifier, unique per study. Example: S01.
- section: Name of the section that the sentence belongs to, as specified in the report. If the sentence is not part of a section, FINAL\_REPORT\_NO\_SECTION or PRE\_FINAL\_REPORT\_NO\_SECTION are used.
  Examples: FINDINGS, IMPRESSION, REASON\_FOR\_EXAM.
  - section\_type: The identified type of section used for classifying the type of content of
    the sentence. IGNORE is used for irrelevant sections.
     Examples: FINDINGS, IMPRESSION, INDICATION.
  - sentence: The raw sentence as written in the report.

**Observation Nodes** Observation nodes are created for each individually described aspect (i.e. observation) in the report's FINDINGS or IMPRESSION section. Hereby, a single sentence may be related to several observation nodes and a single observation may be derived from several sentences (if they describe related aspects). Observation nodes are structured hierarchically, i.e. they may have other observation nodes as parents. An observation node contains the following attributes:

- obs\_id: Identifier, unique per study.
   Example: 001.
   For child nodes this also contains the parent id, e.g. 001.02.
- summary\_sentence: Textual description of the observation, directly derived from the associated report sentences. In some cases, this may be an exact copy of the report sentences but it may also paraphrase parts of it.
- name: Abbreviated version of the summary\_sentence.
- child\_level: Hierarchy level, 0 for top-level, larger numbers for deeper hierarchy levels.
- child\_type: Type of parent-child relation.

  Possible options: regional\_distinction, related\_region, associated\_with, device\_part, recommendation, comparison\_only.
- regions: List of associated regions, each paired with an optional list of distance annotations. Example: [("heart", ["1 cm above"])]
  - non\_resolved\_regions: Similar to regions but with regions that could not be semantically mapped to reference definitions.

- laterality: Laterality of the region.
  Possible options: left, right, likely bilateral, bilateral, unknown.
  - default\_regions: List of regions that have been added because they are defaults for the identified findings (obs\_entities).
    - obs\_entities: List of (directly) associated findings. Example: ["pleural effusion"].

499

500

501

502

503

504

505

506

507

508

509

510

512 513

514

515

516

518

519

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

538

539

544

545

- obs\_entities\_parents: List of findings that are considered parents of findings in obs entities.
- non\_resolved\_obs\_entities: Similar to obs\_entities but with findings that could not be semantically mapped to reference definitions.
- obs\_categories: List of associated finding super-categories.
   Example: ["ANATOMICAL\_FINDING"].
  - obs\_subcategories: List of associated finding sub-categories. Example: ["LUNG\_FIELD"].
  - probability: Likelihood of the observation being positive. Short term, derived from what is mentioned in the report.
  - certainty: How certain is the observation. Derived from probability.

    Possible options: certain, likely, uncertain, comparison\_only, recommendation.
  - positiveness: Whether the observation is positive or negative. Derived from probability.
- Possible options: pos, neg, comparison\_only, recommendation.
  - modifiers: Modifiers of the finding. Dictionary with keys for each type of modifier and lists of the individual modifier values.
- Possible modifier type: severity, texture, spread, temporal.
- Example: {"severity": ["mild"], "spread": ["focal"]}.
  - change\_sentence: Optional textual description of any changes to the prior study of the same patient, if it was mentioned in the report.
  - changes: List of change types mentioned in the change\_sentence. Example: worsening.
  - from\_report: Whether this observation was explicitly mentioned in the report (true) or automatically added (false).
  - obs\_quality: Extraction quality of the observation, consisting of several individual aspects.
     See Tab. 8.
  - localization: Bounding boxes for this observation, for each associated image. Dictionary with keys equaling image ids (each study may correspond to several images). Values contain:
    - image\_id
    - **bboxes**: List of bounding boxes in the  $(x_1, y_1, x_2, y_2)$  format in original image-pixel coordinates.
    - localization\_reference\_ids: List of region names from which the bounding boxes are derived.
    - missing\_localization: List of associated region names for which no localization is available for this image.
    - is\_fallback: Whether this localization is a fallback, i.e. the original region localization was not available but a more coarse localization was used instead.
    - localization\_quality: Quality of the localization. See Tab. 8.

Region Nodes Region nodes are created for each anatomical structure mentioned in any observation and for key regions. They contain the following attributes:

- region: Name of the region and unique identifier within each study. Example: left lung.
- laterality: Laterality of the region.
- Possible options: left, right, bilateral, unknown (i.e. not clearly definable).

- localization: Bounding boxes for this region, for each associated image. Same format as for observation nodes.
  - region\_localization\_quality: Quality of the localization attribute. See Tab. 8.

**Indication Node** Each study contains an optional indication node with information extracted from the INDICATION section. If present, it contains the following attributes:

- indication\_summary: Summary of the indication, directly derived from the INDICATION section of the report, but typically paraphrased.
- patient\_info: Any information about the patient, if mentioned in the INDICATION section. A subset of the content in indication\_summary.
  - indication: Indication for the study, if mentioned in the INDICATION section. A subset of the content in indication\_summary.
  - evaluation: Any required evaluation of the patient (i.e. what should be evaluated with this study), if mentioned in the INDICATION section. A subset of the content in indication\_summary.
  - associated\_sentence\_ids: List of sent\_ids from sentence nodes that are related to the indication.
  - associated\_obs\_ids: List of obs\_ids from observation nodes that are related to the indication.
  - answer\_for\_indication: A single observation node containing the answer to the question (implicitly) asked by the provided indication. This is a special observation node with obs\_id = OIND. Its textual description is directly derived from the FINDINGS and IMPRESSION sections but conditioned on the INDICATION section.

Root Node and Relations Each study contains a single root node called the ReportGraph. It contains general metadata about the study and its scene graphs:

- patient\_id: Unique patient ID, the subject\_id from MIMIC-CXR.
- study\_id: Unique study ID, from MIMIC-CXR. Each patient may have several studies.
- study\_quality: The overall extraction quality of the scene graph for this study, consisting
  of several individual aspects. See Tab. 8.
  - study\_img\_localization\_quality: Dictionary of localization qualities for each image with keys corresponding to image IDs. See Tab. 8.

Additionally, it is connected to all other nodes and links to the top-level (root) observations. Thus, it contains the following:

sentences: List of all sentence nodes.

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

572

573 574

575

576

577

580

581

582

583

588

591

- observations: Dictionary of all observation nodes, indexed by their obs\_id.
- top\_level\_obs\_ids: List of all top-level (root) observation node IDs, i.e. their obs\_ids.
- regions: Dictionary of all region nodes, indexed by their region attribute.
- indication: The indication node, if it exists.
- Nodes can also be connected by the following relations:
- located\_at\_relations (observation region): Specifies where an observation is located with the following additional attributes:
  - distances: List of distance annotations, e.g. ["3cm above"].
  - where\_specified: How this relation was derived.
    Possible options: direct, bilateral, sub\_region.
    - obs\_relations (observation + observation): Specifies a parent-child relation between two observations, with the following additional attribute:

- child\_type: Type of parent-child relation.
  Possible options: regional\_distinction, related\_region, associated\_with,
  device\_part, recommendation, comparison\_only.
  - obs\_sent\_relations (observation \( \to \) sentence): Specifies from which sentences an observation was derived.
  - region\_region\_relations (region \rightarrow region): Specifies a relation between two regions with the following additional attribute:
    - relation\_type: Type of relation.
       Possible options: sub\_region, bilateral (the bilateral version of a region), left (the left version of a region), right (the right version of a region).

### D.2 Question-Answer Structure

596

597

599

602

603

607

608

609

610

611

612

614

616 617

619

620

621

622

623

624

625

630 631

632

633

634

635

637

638

QA-Pair Each question-answer pair consists of a free-text question (attribute question), an answer consisting of structured answer parts (attribute answers). Additionally, it contains the following metadata:

- question\_id: Identifier, unique within the associated study.
- question\_type: The QA-template used to generate this QA-pair.
- question\_strategy: The strategy used to generate QA-pair. See Sec. 3.2 and Appendix E.2.2.
- variables: Key-value pairs of variables (and their values) used during generation, e.g. to fill the template. See Appendix E.2.1.
- obs\_ids: List of obs\_idss of observation nodes (in the scene graph) from which the answer is derived.
- contains\_report\_answers: Whether any of the answer parts was derived from the report,
   i.e. from observation nodes.
  - contains\_template\_answers: Whether any of the answer parts was generated based on a template.
  - extraction\_quality: The overall extraction quality of the associated observations in the scene graph, consisting of several individual aspects. See Tab. 8.
  - question\_img\_localization\_quality: Quality of the localizations per image. See Tab. 8.
- question\_quality: The overall question-answer text quality, consisting of several individual aspects. See Tab. 9.
  - rating: The overall rating of the QA-pair. See Appendix D.3

Answers are structured hierarchically, consisting of a list of answer parts (attribute answers) and sub-answers (children) of these answers, where there can be several hierarchy levels. The hierarchy levels are derived from the parent-child structure of associated observation nodes (based on obs\_relations, Appendix D.1). Additionally, there are different types of answer parts:

- main\_answer: Required to answer the question. There is always at least one main-answer per question.
- details: Providing additional details for the main answer, which are however not mandatory to answer the question.
- related\_information: Not directly answering the question, but may be related and provides context.
- 636 Each individual answer part contains the following attributes:
  - answer\_id: Identifier, unique within each study. Contains the question\_id.
  - text: The answer text. Either generated from a template or based on summary\_sentence in the observation node (Appendix D.1).

- answer\_level: Hierarchy level, 0 for top-level answer part, larger numbers for deeper hierarchy levels (sub-answers).
  - answer\_type: Type of answer part.
  - Possible options: main\_answer, details, related\_information.
  - name\_tag: Abbreviated version of the text. Either generated from a template or based on name in the observation node (Appendix D.1).
  - laterality: Laterality of the region. See observation node (Appendix D.1). Possible options: left, right, likely bilateral, bilateral, unknown.
- regions: List of associated regions. See observation node (Appendix D.1). Distances are not provided here.
- Example: ["heart"]

643

644

645

646 647

651

652

654

655

656

657

658

659

660 661

662

663

664

665

666

667

668

669

671

672

673

674

675 676

677

679

680

684

685

687

688

689

- obs\_entities: List of (directly) associated findings. See observation node (Appendix D.1).
- Example: ["pleural effusion"].
  - obs\_entities\_parents: List of findings that are considered parents of findings in obs\_entities. See observation node (Appendix D.1).
  - obs\_categories: List of associated finding super-categories. See observation node (Appendix D.1).
  - Example: ["ANATOMICAL\_FINDING"].
    - obs\_subcategories: List of associated finding sub-categories. See observation node (Appendix D.1).
  - Example: ["LUNG\_FIELD"].
    - certainty: How certain is the observation. See observation node (Appendix D.1). Possible options: certain, likely, uncertain, comparison\_only, recommendation.
      - positiveness: Whether the observation is positive or negative. See observation node (Appendix D.1).
      - Possible options: pos, neg, comparison\_only, recommendation.
      - modifiers: Modifiers of the finding. List of pairs of modifier type and value. See observation node (Appendix D.1).
        - Possible modifier type: severity, texture, spread, temporal.
- Example: [("severity", "mild"), ("spread", "focal")].
  - localization: Bounding boxes for this answer part, for each associated image. Dictionary
    with keys equaling image ids (each study may correspond to several images). See observation
    node (Appendix D.1).
    - sub\_answers: List of child answers (deeper in the hierarchy). Each sub-answer is another
      answer-part with all attributes and potentially further sub-answers.
      - from\_report: Whether this answer part is derived from the report, i.e. an observation node (true), or from a template (false).
    - extraction\_quality: The overall extraction quality of the associated observations in the scene graph, consisting of several individual aspects. See Tab. 8.
    - answer\_quality: The overall answer text quality, consisting of several individual aspects. See Tab. 9.

### 682 D.3 Quality

- Ratings We distinguish between the following overall ratings for each QA-pair:
  - A++: Perfect and complete content; all information in the answer is explicitly mentioned in the report.
    - A+: Perfect and mostly complete content; all information in the answer is explicitly mentioned in the report, but some minor details may be missing or irrelevant.
    - A: Very good content with minor issues not affecting the overall quality; some tags or boxes
      may be inferred or minor issues (e.g. grammatical) may be present in the text.

- **B**: Good content; factually correct answers, which may however be not fully complete or slightly unclear.
  - C: Poor content; answers may be misleading or contain completely unclear information.
  - **D**: Incorrect content; answers may be contradicting the report and are not usable.
  - not rated: Quality could not be assessed, e.g. due to invalid LLM-rater outputs.
- These ratings are derived based on individual aspects that will be described in the following paragraphs.
  Possible quality levels for each aspects and the resulting rating are presented in Tabs. 8 and 9. The
  final rating is computed as the minimum (worst) rating over all individual aspects.
- Scene Graph Extraction Quality For each scene graph, we provide a quality rating based on how well it could be constructed/extracted. Tab. 8 shows the considered aspects with their potential quality levels and resulting ratings.

Table 8: Quality levels for the 6 scene graph quality aspects, with their resulting ratings.

	Quality level	Value	Resulting rating
How we	ell are region tags identified? (attribute reg	ions)	
	NO_REGIONS	0	В
n ion	DEFAULT_REGIONS_ONLY	1	В
Region	CONTAINS_DEFAULT_REGIONS	2	A
ex x	CONTAINS_NON_RESOLVED_REGIONS	3	A
	RESOLVED_REGIONS_ONLY	4	A++
How we	ell are finding tags identified? (attribute obs	_entitie	es)
ion	NO_ENTITIES	0	В
Finding	CONTAINS_NON_RESOLVED_ENTITIES	1	A
EX EX	RESOLVED_ENTITIES_ONLY	2	A++
How we	ell are textual descriptions extracted? (attri	butes sum	nary_sentence and name)
tion	CHANGE_IN_SENTENCE_OR_NAME	0	В
Description	UNDERSCORES_IN_SENTENCE_OR_NAME	1	A
Des	NO_ISSUES	2	A++
How we	ell are mentions of change extracted? (attri	butes char	age_sentence and change)
	CHANGE_SENTENCE_REMOVED	0	В
Change	UNDERSCORES_IN_CHANGE_SENTENCE	1	A
Cha extra	CONTAINS_NON_RESOLVED_CHANGES	2	A
	NO_ISSUES	3	A++
Have th	nere been any issues in the extraction and s	cene grap	h construction pipeline?
	DISCARDED	-1	D
п	NON_INTERPRETABLE	0	C
extraction	MOSTLY_INTERPRETABLE	1	В
Extra	IGNORABLE	2	A
_	FIXABLE	3	A+
	NO_ISSUES	4	A++
How we	ell could observations/regions be localized?	(attribute	localization))
	NO_LOCALIZATION	0	В
Localization	FALLBACK_LOCALIZATION	1	В
aliza	INCOMPLETE_LOCALIZATION	2	A
Loc	BBOX_LOCALIZATION	3	A++
	BBOX_AND_MASK_LOCALIZATION	4	A++

Finding extraction is also referred to as entity extraction, description extraction as sentence/name quality.

702

701

692

693

**QA Text Quality** For each QA-pair, we provide quality rating for its text, i.e. the question text and the textual descriptions in its answer parts. Tab. 9 shows the considered aspects with their potential quality levels and resulting ratings.

Table 9: Quality levels for the 5 QA-pair text quality aspects, with their resulting ratings.

	Quality level	Value	Resulting rating
Does	s the answer factually align with the origin	ıal report	?
(rate	d per answer-part, given the question and the	e report)	
	NON_ALIGNED_CONTRADICTING	-3	D
	NON_ALIGNED_MISLEADING	-2	C
ent	NON_ALIGNED_NON_INFERABLE	-1	В
Entailmen	ALIGNED_GENERAL_STATEMENT	0	A
En	ALIGNED_NEGATIVE_NOT_MENTIONED	1	A+
	ALIGNED_INFERABLE	2	A++
	ALIGNED_MENTIONED	3	A++
Is th	e answer relevant for the given question?		
(rate	d per answer-part, given the question but ind	lependent	of the report)
e	IRRELEVANT_INFO	-2	A
Relevance	REDUNDANT_INFO	-1	A
Relev	RELATED_INFO	0	$A + (A +\!$
	RELEVANT_MAIN_ANSWER	1	$A \!+\!\!+\! (A \ for \ \texttt{related\_information} \ answer)$
Does	s the answer cover all aspects in the report	t that are	relevant to the question?
(rate	d for the full answer, given the question and	the report	)
só.	INCOMPLETE_MISLEADING	-2	C
Completeness	INCOMPLETE_NON_MISLEADING	-1	В
ıplet	NOT_ANSWERED	0	В
Con	DETAILS_MISSING	1	A+
	FULLY_COMPLETE	2	A++
Is th	e generated question clear and grammatic	cally corr	ect?
(rate	d for the question, given nothing else)		
	UNANSWERABLE	-3	C
arity	UNRELATED_TO_CHEST_XRAY	-2	В
Question clarity	UNCLEAR_QUESTION	-1	В
estic	GRAMMATICAL_ERRORS	0	A
õ	UNUSUAL_SENTENCE_STRUCTURE	1	A
	OPTIMAL	2	A++
Is th	e answer clear and grammatically correct	:?	
(rate	d per answer-part, given nothing else)		
>-	NOT_UNDERSTANDABLE	-2	C
larit	UNCLEAR_ANSWER	-1	В
Answer clarity	GRAMMATICAL_ERRORS	0	A
Ansv	UNUSUAL_SENTENCE_STRUCTURE	1	A
,	OPTIMAL	2	A++

### E Dataset Construction Details

### 708 E.1 Scene Graph Construction

### 709 E.1.1 Region Localization

717

718

719

720

721

722

723

724

725

726

727

728

729

730

735

We use the CXAS [22], [23] model to predict segmentation masks of 158 anatomical structures on the 377,110 CXRs from MIMIC-CXR-JPG [15], [24], [25]. Additionally, we use the bounding boxes provided by the Chest ImaGenome [15], [17], [19] dataset, which are provided for 29 anatomical structures in most frontal images of MIMIC-CXR. The masks predicted by CXAS are post-processed with morphological operations to filter out outlier pixels.

We specify 257 localized regions in our reference definitions. For each of these regions, we define how the bounding boxes are derived. We consider the following options:

- CXAS masks: Some regions are directly associated with one of the 158 anatomical structures for which the CXAS model predicts segmentation masks. In these cases, we compute the bounding box around the predicted segmentation mask.
- Chest ImaGenome boxes: Some regions are directly associated with one of the 29 anatomical structures for which Chest ImaGenome provides bounding boxes. In such cases, we use these provided bounding boxes if no CXAS masks are associated.
- **Bilateral regions**: Some regions refer to a pair of bilateral regions (e.g. *lungs* refers to *left lung* and *right lung*). In these cases, we simply use the two bounding boxes of the left and right versions, but do not fuse them.
- Parent regions: For some regions we do not have exact correspondences to available masks
  or boxes but we have available sub-regions. In these cases, we compute the super bounding
  box, a single box, around all specified child regions.
- Fusions: In some rare cases, we combine multiple individual masks or bounding boxes. We compute intersections or unions of boxes or masks, before inferring the final bounding box.

After computing all regions, we filter out regions with a too small bounding box area. For images where a specific region is not available, we try to use alternative regions as fallbacks instead, e.g. using a more coarse parent regions as an alternative. Note that this is often the case for lateral images as there no Chest ImaGenome boxes are available.

### **E.1.2** Information Extraction

**Extracting the Sentences** First, we extract individual sentences from the reports, detect their 736 sections (e.g. FINDINGS, IMPRESSION, INDICATION, ...), discard sentences without relevant information, and merge sentences containing similar information (e.g. if findings are described in 738 both the FINDINGS and IMPRESSION section). Therefore, each full report is passed in a single 739 step to the LLM, which predicts the individually separated sentences as well as their sections and 740 related sentences. We use the prompt shown in Listing 1 (with few-shot examples similar to Listing 741 2) and apply it to the full radiology report. After parsing the LLM outputs, we apply the Stanza [39] 742 tokenizer to each identified sentence and try to further split it. The LLM also identified potentially 743 related sentences. We use this information to identify sentence clusters containing related information. Such sentence clusters are the basis for the next step, i.e. observation extraction. We successfully extracted sentence from 227 626 studies (reports) while having parse errors for 209 studies.

Listing 1: LLM prompt used for sentence extraction.

```
Extract all sentences from the given textual report.
748
    You will be given a (free-text) medical radiology report describing
749
750
        \hookrightarrow one or more chest X-rays of a single patient.
751
    # Rules:
752
    - Split the report into sentences and extract all sentences in the
753
        \hookrightarrow report.
754
    - Do not rewrite the sentences!
755
    - For each sentences, identify the its section name (written in the
756
        \hookrightarrow report).
757
    If a sentence is not part of a section but is part of the "FINAL
758
        \hookrightarrow REPORT", then use "FINAL\_REPORT\_NO\_SECTION". If a sentence
759
        \hookrightarrow is not part of a section but the sentence is before the "FINAL
760
        \hookrightarrow REPORT", use the section name "PRE\_FINAL\_REPORT\_NO\_SECTION"
761
762
      For each sentence, classify the content written therein into one of \hookrightarrow the following types: <code>[EXAM\_TECHNIQUE</code>, <code>INDICATION</code>, <code>FINDINGS</code>,
763
764
        \hookrightarrow IMPRESSION, PRE\_FINAL\_REPORT, IGNORE]. This is typically
765
        \hookrightarrow inferred from the section name but may also be influenced by
766
        \hookrightarrow the content of the sentence. Some example sections names for
767
        \hookrightarrow each type are given below:
768
         EXAM\_TECHNIQUE: EXAMINATION, EXAM, TECHNIQUE
769
         INDICATION: INDICATION, INDICATIONS, HISTORY, CLINICAL HISTORY,
770
             \hookrightarrow CLINICAL, REASON, REASON FOR EXAM
771
772
         FINDINGS: FINDING, FINDINGS
         IMPRESSION: IMPRESSION, IMPRESSIONS, RECOMMENDATION
773
         PRE\_FINAL\_REPORT: WET\_READ, WET\_READ\_VERSION\_#1, PRE\_FINAL\
774
             \hookrightarrow _REPORT\_NO\_SECTION
775
         IGNORE: COMPARISON, COMPARISONS, REFERENCE EXAM, NOTIFICATION
776
    - Split the report into individual sentences and report each sentence
777
        \hookrightarrow in its own line, removing any newlines present in the sentence.
778
      For enumerations: each point is considered an independent sentence!
779
        \hookrightarrow Remove the numbering.
780
781
      Specify sentence IDs of similar, previous sentences that each
782

ightarrow sentence could be merged with. A sentence should be merged with
        \hookrightarrow all previous sentences that either describe the same aspect or
783
        \hookrightarrow that refer to each other (e.g. if a sentence provides further
784
        \hookrightarrow details to a previous one). A bullet point may also be
785
        \hookrightarrow associated with a sentence, even if the other sentence has a
786
        \hookrightarrow different bullet number or none at all.
787
    - Follow the examples given below!
788
789
    # Examples:
790
     <FEWSHOT>
791
792
    # Input Report (extract data from this report):
793
    --- START OF REPORT ---
794
     <REPORT>
    --- END OF REPORT ---
796
797
    # Hints:
798
     Infer the output format from the examples!
799
    - Do not add any explanations or text BEFORE or AFTER the extracted
800
        \hookrightarrow sentences, i.e. start with the first sentence!
801
802
    # Proceed:
883
```

Listing 2: Few-shot example for sentence extraction.

```
***Example: Report***
806
       - START OF REPORT ---
807
                                        FINAL REPORT
808
    PORTABLE CHEST OF \_\_\_
809
810
    COMPARISON: \_\_\ radiograph.
811
812
    FINDINGS: No pleural effusion or pneumothorax.
813
    --- END OF REPORT --
814
815
    ***Example 5: Output ***
816
    [S01] FINAL\_REPORT\_NO\_SECTION(EXAM\_TECHNIQUE) - merge with []:
817
        \hookrightarrow PORTABLE CHEST OF \_\_\_
    [SO2] COMPARISON(IGNORE) - merge with []: \_\_\ radiograph.
819
    [SO3] FINDINGS(FINDINGS) - merge with []: No pleural effusion or
820
        \hookrightarrow pneumothorax.
821
```

Extracting the Observations In this step, we consider each sentence cluster (as identified during sentence extraction), in the FINDINGS and IMPRESSION sections. A sentence cluster contains one or more sentences that describe related aspects and may stretch over one of both of these sections. From each of these clusters, we now extract mentioned observations using the prompt shown in Listing 3 with few-shot examples similar to Listing 4. We apply this prompt to each sentence cluster individually and extract zero, one, or multiple observations each. The output is provided in the json-format and follows a similar structure as the final observation node, but we optimized it to be easy to fill by the LLM. The LLM is allowed to freely assign values to each of the json-fields. For name and summary\_sentence, we prompt the model to stay close to the original sentence, but it must remove any mentions of change and only keep the part relevant to the individual observation (if several observations are mentioned in one sentence). We successfully extracted observations from 227 266 studies (reports) while having parse errors for 360 studies.

823

824

825

826

829

830

831

832

833

Listing 3: LLM prompt used for observation extraction.

```
835
    Extract structured information from the given textual report.
836
    You will be given sentences from a (free-text) medical radiology
837
        \hookrightarrow report describing one or more chest X-rays of a single patient.
838
839
    # Guidelines:
840
    <GUIDE>
841
842
    # Rules:
843
844
    - Follow the examples given below!
845
    # Examples:
846
     <FEWSHOT>
847
848
    # Hints:
849
      Check for any "change" modifiers (see guidelines).
850
      If there is a "change" modifier, rewrite the "summary\_sentence"
851
        \hookrightarrow such that it describes only what is visible in the current
852
853
        \hookrightarrow image, without any mentions of change or comparisons! Describe
        \hookrightarrow the change in the "change\_sentence". Do this for all top-level
854
           AND child observations.
855
      Make sure to include all children of observations, even if they
856
        \hookrightarrow repeat information from the parent!
857
    # Proceed with the Input Sentence:
859
    Sentence(s): <SENT>
860
    Output JSON-List:
862
```

Listing 4: Few-shot example for observation extraction.

```
Sentence(s): Left more than right basilar atelectasis.
864
    Output JSON-List:
865
866
         {
867
              "name": "bibasilar atelectasis", "entity": "atelectasis",
868
              "probability": "positive", "change": null,
869
              "summary\_sentence": "Bibasilar atelectasis.",
870
              "change\_sentence": null,
871
              "regions": ["bibasilar"],
872
              "children": [
873
                   {
874
                        "child\_type": "regional\_distinction", "name": "left
875
                        \hookrightarrow basilar atelectasis", "entity": "atelectasis", "probability": "positive", "change": null,
876
877
                        "summary\_sentence": "Left more than right basilar
878
                            \hookrightarrow atelectasis.",
879
                        "change\_sentence": null,
880
                        "regions": ["left basilar"]
881
                   }
882
              1
883
         }
884
    ]
885
```

**Extracting the Indication** Next, we extract information about the INDICATION section and detect which FINDINGS or IMPRESSION sentences may provide information related to the indication. Therefore, the extracted INDICATION sentences and a list of all FINDINGS and IMPRESSION sentences are passed to the LLM using the prompt shown in Listing 5 with few-shot examples similar to Listing 6. The LLM predicts a json-structure containing several text fields for summaries of aspects in the indication, an answer\_for\_indication derived from the FINDINGS and IMPRESSION section, as well as relevant sentence IDs. We successfully extracted indictions from 227 596 studies (reports) while having parse errors for 30 studies.

887

888

891

892

893

Listing 5: LLM prompt used for indication extraction.

```
895
896
    Extract structured information from the given (free-text) medical
897
        \hookrightarrow report.
    You will be given the indication sentence from a report and
898
        \hookrightarrow additionally the sentences from the findings section.
899
900
901
    # Rules:
      Extract / summarize the given indication information. Use only the
902
        \hookrightarrow provided indication sentence.
903
      Additionally, identify the finding sentences associated with the
904
        \hookrightarrow indication, i.e. the sentence that answer the quesiton of the
905
        \hookrightarrow indication or are highly relevant to it. Based on these finding
906
            sentences, provide an answer to the question asked in the
907
        \hookrightarrow indication.
908
    - Follow the examples given below!
909
910
    # Examples:
911
     <FEWSHOT>
912
913
    # Hints:
914
      For each attribute, write full sentences instead of single terms or
915
        \hookrightarrow bullet points.
916
      In the "answer\_for\_indication", describe in YOUR OWN WORDS how the
917
        \hookrightarrow question asked in the evaluation can be answered based on the
918
        \hookrightarrow findings. Only include the key information.
919
    - Use the JSON structure from the examples!
920
921
922
    # Proceed with the Input:
    **Input:**
923
```

```
924 INDICATION: <IND>
925 FINDINGS:
926 <FIND>
927

928 **Output JSON:**
```

Listing 6: Few-shot example for indication extraction.

```
931
932
    **Input:**
    INDICATION: \_\_\F with new onset ascites
                                                      // eval for infection
    FINDINGS:
934
    [S01] There is no focal consolidation, pleural effusion or
935
        \hookrightarrow pneumothorax.
936
937
    [S02] No acute cardiopulmonary process.
938
    **Output JSON: **
939
    {
940
        "patient\_info": "female",
941
        "indication": "New onset ascites.",
942
        "evaluation": "Evaluate for infection.";
943
        "indication\_summary": "Female with new onset ascites; should be
944
             \hookrightarrow evaluated for infection.",
945
        "associated\_findings": ["S02"],
946
        "answer\_for\_indication": "Evaluation for infection is negative:
947
            \hookrightarrow There is no acute cardiopulmonary process."
948
949
    }
```

### E.1.3 Building Scene Graphs

**Entity Mapping** We apply semantic entity mapping to modifiers (used to fill the attributes probability, certainty, positiveness, and modifiers), regions (attribute regions), finding entities (attribute obs entities), and changes (attribute changes).

For each of these we consider the associated tags extracted by the LLM during observation extraction and encode them into text embeddings using the BioLORD [29] model. We also encode all potential tags and their synonyms, defined for each type of tag in our reference definitions. Then we compute the cosine similarities of each tag with all reference tags of the same type. We pick the reference tag with the highest cosine similarity but threshold it at 0.5. If no reference tag was identified with cosine similarity  $\geq 0.5$ , then we mark the tag as non-resolved. For finding entities, we follow a slightly more complicated matching approach. Instead of only considering the finding entity tags extracted by the LLM, we also consider pairs of these entities and extracted region tags as well as the extracted summary sentences and names for matching. We then try to match each of those with the reference finding tags and pick the ones with the highest cosine similarities.

The matched reference finding tags are stored in the obs\_entities attribute (non-resolved ones are kept in non\_resolved\_obs\_entities), matched reference regions are stored in the regions attribute, where we also store the distance as identified by the LLM (non-resolved regions are kept in non\_resolved\_regions). The matched changes are store in the changes attribute (non-resolved changes are discarded). For all modifiers, we use the modifier type defined for the matched reference tag. We matched all modifiers against all types of modifiers, which means that the modifier type identified by the LLM can be overwritten during matching. Finally, we extract the probability from the modifiers (this is a special modifier type), store it in the probability attribute and infer the certainty and positiveness attributes from it (using the reference definitions). The remaining modifiers are stored in the modifiers attribute (non-matched ones are discarded).

We additionally try to identify the laterality of the observations. Here, we do not use semantic entity mapping but rely on keywords instead. We consider the raw finding entities, regions, as well as the summary sentences, and search for any laterality-related mentions such as *left*, *right*, *bilateral*, and related terms. From this we infer the laterality and store it into the laterality attribute.

979 **Reference Data and Standardization** Using the reference definitions, we infer all 980 obs\_entities\_parents, obs\_categories, obs\_subcategories, and default\_regions from 981 the matched obs\_entities.

- Next, we inspect the summary\_sentence and name attributes (extracted by the LLM) for underscores or mentions of changes. We track such issues (which are used for quality assessment) but do not apply any cleanup. Similarly, we check the change\_sentence for underscores and assert that it contains mentions of changes.
- We further inspect the structure of observations and their children. If an observation mentions multiple different findings and has one child for each of these findings, then we lift these children to the top-level and discard the parent. Similarly, we merge multiple duplicate observations into one.
- Finally, we try to resolve missing regions or improve their precision. If no regions could be extracted, we rely on the default\_regions derived from the obs\_entities instead, but consider the identified laterality. We also check whether these default\_regions are more precise than extracted ones. Then we check whether any identified region contradicts the identified laterality and remove them. We then either split or merge bilateral versions of the same region.
- Graph Construction Based on the matched regions, we associate bounding boxes with the observations if available. Additionally, we build a tree of all mentioned regions and fill missing intermediate regions based on the reference data. This allows us to build a graph of region nodes relevant to the study.
- We construct region\_region\_relations based on the reference data alone.
  located\_at\_relations are constructed based on the regions attribute of observations
  (direct specified). Additionally, we infer located\_at\_relations relations for sub regions
  (sub\_region) and bilateral versions of regions (bilateral). obs\_relations are constructed
  based on the parent-child structure of observations and their child type, as predicted by the LLM.
  obs\_sent\_relations are constructed based on the sentences each observation was derived from.
- Finally, we attach the indication information extracted from the report. Therefore, we build an additional observation node based on the LLM-extracted answer\_for\_indication and the LLM-extracted associated associated sentences, from which we can infer the associated observations and can infer all relevant tags.

### **E.2** Question-Answer Generation

### E.2.1 Template Engine

1008

1009

1013

1014

1015

1016 1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

To construct QA-pairs, we develop a template-engine that considers the information in the scene graphs to construct the answers. The template engine generates a QA-pair by running the following steps:

- 1. Filter observations and studies based on the template configuration, e.g. only keeping observations of specific sub-categories.
- 2. Run a QA-strategy (indication, study abnormality, region abnormality, or finding) on the remaining scene graph. The strategy provides multiple named subsets of observations, variables to fill the template, as well as an overall state consisting of multiple tags (e.g. is the study positive, are there any devices, ...).
- 3. Construct the template-based main answer by selecting and filling the answer-template based on the state returned by the QA-strategy and the returned variables. Tags and bounding boxes can be inferred from defined observation subsets. (Not all templates provide such main answers)
- 4. Pick observation subsets identified by the QA-strategy and convert the observations into answer parts. The template configuration defines which subsets are picked and how they are ordered. Additionally, template-based prefix- or fallback-answers can be defined for each subset. Some subsets can also be excluded based on the QA-strategy state. These answers can be main-answers, details, or related information as defined in the configuration.

Additionally, the template engine supports variables, i.e. each template can be used to generate multiple QA-pairs. Variables can either be defined as lists (configured in the template) or can be

provided by the QA-strategy (which might infer variables from the current scene graph, e.g. all mentioned regions). The question may then also contain such template variables.

### **E.2.2** Strategies and Templates

**Indication** In this strategy, we use the extracted indication (if available) as the question. More precisely, we use the indication\_summary attribute from the indication node as the question text. The main-answer is constructed from the indication observation (i.e. the answer to the indication based on the finding sentences), while detail answer parts are constructed based on all associated finding observations. We include this question, if an indication observation is present in the scene graph.

**Study abnormality** In this strategy, we generate questions about abnormalities. This includes descriptions of the full study or specific categories of observations (e.g. devices), description of only abnormal findings, and yes/no questions of whether there are positive findings (overall or of specific categories) present in the study. We use the templates defined in Tab. 10.

The strategy identifies five types of observations: (i) finding (positive), (ii) finding (negative), (iii) device (positive), (iv) device (negative), (v) acquisition. Based on the specific template, these are selected as main-answers, details, or related information. Additionally, a template answer can be included, which is selected based on whether the study is abnormal or not. Some templates use different subcategories as variables, i.e. one question is generated for each of the defined subcategories, where observations are filtered based on this subcategory.

Table 10: Study abnormality templates.

Template (ID)	Question Example	Variables	Main answer	Details	Related Inf.
describe_all B01_describe_all	Describe the given study.	-	finding (positive) device (positive) device (negative) finding (negative) acquisition	-	-
describe_abnormal B02_describe_abnormal	Describe all abnormal findings in the given study.	=	finding (positive)	-	device (positive)
is_abnormal B03_is_abnormal	Are there any abnormal findings?	=	template finding (positive)	finding (negative)	device (positive) finding (negative)
is_normal BO4_is_normal	Is the study normal?	-	template	finding (positive) finding (negative)	finding (negative) device (positive)
describe_subcat B08_describe_subcat	Evaluate the cardiac structures.	subcategory	finding (positive) finding (negative)	=	=
describe_abnormal_subcat B09_describe_abnormal_subcat	Describe any pulmonary diseases and disorders suggested by the study.	subcategory	finding (positive)	_	_
is_abnormal_subcat B10_is_abnormal_subcat	Are there any fractures or bone diseases apparent from the study?	subcategory	template finding (positive)	finding (negative)	finding (negative)
is_normal_subcat B11_is_normal_subcat	Are the mediastinal and hilar contours normal?	subcategory	template finding (positive)	finding (negative)	finding (negative)
describe_device B12_describe_device	Check the presence and posi- tion of devices, tubes, lines, and other foreign objects.	subcategory	device (positive) device (negative)	_	_
has_devices B13_has_devices	Are there any signs of prior surgical procedures?	subcategory	template device (positive)	device (negative)	device (negative)
describe_acquisition B14_describe_acquisition	Assess the image quality and describe aspects related to image acquisition.	-	acquisition	-	-
describe_imaging_artifacts B15_describe_imaging_artifacts	Describe any apparent imag- ing artifacts and imaging- related shadows.	-	acquisition	-	-
has_imaging_artifacts B16_has_imaging_artifacts	Are there any imaging artifacts or imaging-related shadows?	-	template acquisition	-	-

**Region abnormality** In this strategy, we generate question about anatomical regions. This includes describing regions, answering yes/no questions about the abnormality of regions, or describing specific aspects of regions (e.g. devices). We use the templates defined in Tab. 11.

For a given region, the strategy first identifies observations associated with that region and classifies them into the five types defined in the study abnormality strategy. Additionally, it identifies observations in related regions. This includes positive findings in parent regions or the opposite laterality. Additionally, a template answer can be included, which is selected based on whether the region is abnormal or not.

Before generating QA-pairs, the strategy first identifies a set of regions. For each of these regions an individual QA-pair is generated. The set of regions is computed as follows: We always include a set of pre-defined default regions (the lungs, the heart, ...) and include all regions explicitly mentioned in observations, as well as their parent regions. Additionally, we randomly sample regions. Their sampling probabilities are computed based on how often they are associated with positive vs. negative findings, i.e. the more often a region is associated with positive findings and the less often it is associated with negative findings, the more often we sample it as a question. This assures that we generate additional negative questions for regions that are only/mostly mentioned with positive findings.

Table 11: Region templates.

Template (ID)	Question Example	Variables	Main answer	Details	Related Inf.
describe_region C01_describe_region	Describe the left lung.	region	finding (positive) device (positive) finding (negative) device (negative)	-	related regions
describe_abnormal_region CO2_describe_abnormal_region	Describe all abnormal findings in the lung bases.	region	finding (positive)	-	device (positive) related regions
is_abnormal_region CO3_is_abnormal_region	Are there any abnormal findings in the mediastinum?	region	template finding (positive)	device (positive) finding (negative)	related regions
is_normal_region CO4_is_normal_region	Is the heart normal?	region	template finding (positive)	region (positive)	finding (negative) related regions
describe_region_device CO7_describe_region_device	Check the right chest for implants.	region subcategory	device (positive) device (negative)	_	related regions
has_region_device CO8_has_region_device	Are there any tubes, lines, or ports in or near the left lung?	region subcategory	template device (positive)	device (negative)	device (negative) related regions

**Finding** In this strategy, we generate question about specific findings (radiological findings, diseases, devices, ...). This includes descriptions of findings, yes/no questions about the presence of findings, location of findings, and severity of findings. We use the templates defined in Tab. 12.

For a given finding/device entity, the strategy first identifies observations associated with it and classifies them into the five types defined in the study abnormality strategy. Additionally, it identifies observations that contain related finding/device entities. This includes parent findings (i.e. findings that are parents of the current one), same subcat findings (i.e. findings having the same sub-category), correlated findings (based on statistics computed over the whole scene graph dataset), indications of the current finding, and findings that are indicative of the current finding. The observation subset can be selected based on the template configuration. Additionally, a template answer can be included, which is selected based on whether the finding is present or not and based on severity levels. This template may also be filled with information about the localization of the finding.

Before generating QA-pairs, the strategy first identifies a set of finding/device entities. For each of these entities an individual QA-pair is generated. The set of entities is computed as follows: We always include a set of pre-defined default entities and include all entities explicitly mentioned in observations, as well as their parent entities. Additionally, we randomly sample entities. Their sampling probabilities are computed based on how often they are mentioned positively vs. negatively (over all scene graphs), i.e. the more often a finding is mentioned positively and the less often it is mentioned negatively, the more often we sample it as a question. This assures that we generate additional negative questions for findings that are only/mostly mentioned positively.

Table 12: Finding templates.

ion Example	Variables	Main answer	Details	Related Inf.
ibe the pleural effusion.	finding	finding (positive) finding (negative)	-	parent findings indications indicative of same subcat correlated
e any indication of pneu-?	finding	template	finding (positive) finding (negative)	parent findings indications same subcat correlated
e is the lung nodule lo-	finding	template	finding (positive) finding (negative)	
severe is the car- galy?	finding	template	finding (positive) finding (negative)	
ibe the endotracheal	device	device (positive) device (negative)	-	parent findings same subcat
acemaker visible in the	device	template	device (positive) device (negative)	same subcat
e are the surgical clips d?	device	template	device (positive) device (negative)	same subcat
	d?			

### 1089 E.3 Quality Assessment

1088

Scene Graph Quality The scene graph quality aspects are computed by simply inspecting the observations nodes and checking which fields are set or empty. Additionally, we track issues during the graph construction procedure and derive quality aspects from them.

QA Quality We automatically assess the quality of the textual content of QA-pairs using Llama 3.1 8B [26] as a judge for the five criteria presented in Tab. 9.

For rating *entailment* (Listing 7), we condition the model on the report, the question, as well as the answer parts and we rate each answer part individually.

Listing 7: LLM prompt used for entailment evaluation of generated answers.

```
You will be given a Report (medical report of a chest X-ray study), a
1098
         \hookrightarrow Question (about the study), and an Answer (to the question)
1099
         \hookrightarrow consiting of several (numbered) sentences.
1100
1101
     Your task is to assess/rate whether each of the answer sentences is
1102
         \hookrightarrow true or not, given a the reference report about the chest X-ray
1103
1104
         \hookrightarrow . This task is known as entailment verification.
     Assess the quality of each answer sentence independently and use one
1105
1106
         \hookrightarrow of the rating options provided below to assess how well the
         \hookrightarrow facts in each sentences align with the report.
1107
1108
     # Guidelines:
1109
      Rate each sentence in the Answer individually; do NOT use any prior
1110
         \hookrightarrow answer sentences as context or source
1111
      Provide the rating for each answer sentences in its own line
1112
1113
         \hookrightarrow starting with the sentence number followed by your rating
      For each sentence, use ONE of the rating options provided below, do
1114
        \hookrightarrow NOT use any other options
1115
     - An example format will be provided
1116
1117
     - DO NOT REPEAT the question or answer sentences in your response!
1118
     ## Rating Options -- ONLY USE ONE OF THE FOLLOWING OPTIONS
1119
     - ALIGNED_MENTIONED: Answer aligns with the report (is factually
1120
         \hookrightarrow correct) and all facts are explicitly stated in the report.
1121
1122
         Example: The same finding is described in the answer and the
             \hookrightarrow report
1123
```

```
- ALIGNED INFERABLE: Answer aligns with the report (is factually
1124
         \hookrightarrow correct) but some facts are NOT explicitly stated in the report
1125
         \hookrightarrow , can however be derived from what is written there.
1126
         Example: The answer provides a more general description of what is
1127
1128
                 written in the report.
     - ALIGNED_NEGATIVE_NOT_MENTIONED: Answer does NOT contradict the
1129
1130
         \hookrightarrow report (may factually correct) but some facts (negative
         \hookrightarrow findings) cannot be derived from the report, are however likely
1131
         \,\hookrightarrow\, correct because they are negative findings and nothing
1132
1133
         \hookrightarrow contradictory is mentioned in the report.
1134
         Example: The answer mentions that a finding is not present but
             \hookrightarrow this is not explicitly mentioned and does not contradict
1135
             \hookrightarrow anything in the report.
1136
     - ALIGNED_GENERAL_STATEMENT: Answer is a more general statement or
1137
         \hookrightarrow summary that is not explicitly mentioned but aligns roughly
1138
1139
         \hookrightarrow with the overall report.
         Example: Summaries of whether the study is positive or negative.
1140
     - NON_ALIGNED_NON_INFERABLE: Answer does NOT contradict the report but
1141
         \hookrightarrow the correctneaa of some facts cannot be validated using the
1142
1143
         \hookrightarrow report.
         Example: The answer mentions that a finding is present but this is
1144
             \hookrightarrow never mentioned in the report and cannot be concluded from
1145
1146
     - NON_ALIGNED_MISLEADING: Answer does NOT directly contradict the
1147
        \hookrightarrow report but the description is highly misleading considering the
1148
         \hookrightarrow report.
1149
1150
         Example: The answer mentions that a finding is not present, which
             \hookrightarrow is never mentioned in the report but could likely be
1151
             \hookrightarrow present considering the report.
1152
     - NON_ALIGNED_CONTRADICTING: Answer contradict with the report.
1153
         Example: The answer describes that a finding is not present, which
1154
             \hookrightarrow is however mentioned as present in the report or vice
1155
             \hookrightarrow versa.
1156
1157
     # Example Format:
1158
1159
     Report:
1160
     --- START OF REPORT ---
1161
     --- END OF REPORT ---
1162
1163
     Question: ...
1164
1165
     Answer (2 sentences to rate):
1166
     [01] First answer sentence.
1167
     [02] Second answer (last sentence in this example).
1168
1169
     Rating (provide 2 ratings):
1170
     [01] ALIGNED_MENTIONED
1171
1172
     [02] NON_ALIGNED_NON_INFERABLE
1173
1174
     # Proceed with the following Report, Question, and Answer sentences:
1175
1176
     Report:
1177
     --- START OF REPORT ---
     <REPORT>
1178
     --- END OF REPORT ---
1179
1180
1181
     Question: <QUEST>
1182
     Answer ( <NUMANS> sentences to rate):
1183
     <ANSWERS>
1184
1185
     Rating (provide <NUMANS> ratings):
1189
```

For rating *relevance* (Listing 8), we condition the model on the question as well as the answer parts (but not on the report) and we rate each answer part individually.

Listing 8: LLM prompt used for relevance evaluation of generated answers.

```
1190
1191
     You will be given a question (about a chest X-ray study), and an
        \hookrightarrow answer (to the question) consiting of several (numbered)
        \hookrightarrow sentences.
1193
1194
     Your task is to assess/rate whether each of the answer sentences
1195
        \hookrightarrow relevant to answer the question or is redundant.
     Assess the quality of each answer sentence and use one of the rating
1197
        \hookrightarrow options provided below.
1198
1199
    # Guidelines:
1200
     - Rate each sentence in the answer individually; but check for
1201
        \hookrightarrow redundancy with previous sentences.
1202
     - Provide the rating for each answer sentences with the sentence
1203
      \hookrightarrow number followed by your rating. For each sentence, use ONE of the rating options provided below, do
1204
1205
        \hookrightarrow NOT use any other options.
1206
     - An example format will be provided.
1207
1208
     - DO NOT REPEAT the question or answer sentences in your response!
1209
     ## Rating Options -- ONLY USE ONE OF THE FOLLOWING OPTIONS
1210
     - RELEVANT_MAIN_ANSWER: Fullfill ALL of the following
1211
         a) Are relevant to the question
1212
         b) Are needed to answer the question or provide details
1213
1214
         c) Are not redundant to previous RELEVANT_MAIN_ANSWER sentences
      RELATED_INFO: Fullfill ALL of the following
1215
         a) Are NOT needed to answer the question
1216
1217
         b) Provide additional context related to the question or other
1218
             \hookrightarrow answer sentences
1219
         c) Are not redundant to any previous sentences
     - REDUNDANT_INFO(...): Fullfill ALL of the following
1220
         a) Would fullfill criteria a-b) for RELEVANT_MAIN_ANSWER, or
1221
1222

→ RELATED_INFO

1223
         b) Contains exactly the same information that was already provided
             \hookrightarrow in a previous sentence of the same type (ONLY consider
1224
             \hookrightarrow previous sentences here!)
1225
1226
         c) Does not provide any additional details or related information
         d) Could be removed without changing the content of the answer
1227
1228
         Note: replace ... with the sentences IDs OF PREVIOUS SENTENCE with
             \hookrightarrow which the current sentence is redundant
1229
      IRRELEVANT_INFO: Fullfill ALL of the following
1230
1231
         a) Does not classify as any of the above
         b) No information in the sentence is relevant or related to the
1232
1233
             \hookrightarrow question
1234
     # Example Format:
1235
     Question: ...
1237
     Answer (4 sentences to rate):
1238
     [01] First answer sentence.
1239
1240
     [02] Second answer.
     [03] Third answer, containint no additional information, everything
1241
         \hookrightarrow was already mentioned in 01 and 02.
1242
1243
     [04] Fourth sentence.
1244
     Rating (provide 4 ratings):
1246
     [01] RELEVANT_MAIN_ANSWER
     [02] IRRELEVANT_INFO
1247
     [03] REDUNDANT_INFO(01,02)
1248
1249
     [04] RELATED_INFO
1250
```

```
1251
     # Proceed with the following Question and Answer sentences:
1252
     Question: <QUEST>
1253
1254
1255
     Answer ( <NUMANS> sentences to rate):
     <ANSWERS>
1256
1257
    Rating (provide <NUMANS> ratings):
1358
```

1260 For rating *completeness* (Listing 9), we condition the model on the report, the question, as well as the full answer and we rate the full answer as a whole.

1261

Listing 9: LLM prompt used for completeness evaluation of generated answers.

```
1262
     You will be given a Report (medical report of a chest X-ray study), a
1263
1264
         \hookrightarrow Question (about the study), and an Answer (to the question)
1265
     Your task is to assess/rate whether the provided Answer contains all
1266
1267
         \hookrightarrow the necessary information to answer the Question, considering
1268
         \hookrightarrow the Report as the source of truth.
1269
     # Guidelines:
1270
1271
     - Do not assess whether the answer is correct but whether it is
1272
         \hookrightarrow contains all relevant information from the Report to answer the
         \hookrightarrow Question.
1273
1274
     - Use ONE of the rating options provided below, do NOT use any other
1275
         \hookrightarrow options.
      Answer with a short explanation (a few words) followed by "->" and
1276
1277
         \hookrightarrow the rating option.
1278
     - An example format will be provided.
     - DO NOT REPEAT the report, question, or answer sentences in your
1279
1280
         \hookrightarrow response!
1281
     ## Rating Options -- ONLY USE ONE OF THE FOLLOWING OPTIONS
1282
     - FULLY_COMPLETE: All facts from the report that are relevant to the
1283
         \,\hookrightarrow\, question are included and the question is answered.
1284
1285
     - DETAILS_MISSING: The main facts from the report that are relevant to
1286
         \hookrightarrow the question are included BUT some details are missing.
     - NOT_ANSWERED: While facts from the report may be contained, the
1287
         \ensuremath{\hookrightarrow} answer does not relate to the question at all.
1288
     - INCOMPLETE_NON_MISLEADING: Main facts are missing, but these should
1289
         \hookrightarrow not lead to a misrepresentation of the facts (e.g. only some
1290
         \hookrightarrow negative findings are not mentioned).
1291
     - INCOMPLETE_MISLEADING: Important facts are missing, such that the
1292
         \hookrightarrow answer may mislead the reader.
1293
1294
     # Example Format:
1295
1296
     Report:
     --- START OF REPORT ---
1297
1298
     --- END OF REPORT ---
1300
     Question: ...
1301
1302
1303
     Answer (to rate):
1304
1305
1306
     Rating (your task):
1307
     severity is missing -> DETAILS_MISSING
1308
1309
     # Proceed with the following Report, Question, and Answer:
1310
1311
     Report:
1312
     --- START OF REPORT ---
     <REPORT>
1313
```

```
--- END OF REPORT ---
1314
1315
     Question: <QUEST>
1316
1317
1318
     Answer (to rate):
      <ANSWERS>
1319
1320
     Rating (your task):
1321
```

For rating *question clarity* (Listing 10), we condition the model on the question only and rate it. 1323

Listing 10: LLM prompt used for question clarity evaluation of generated questions.

```
1324
     You will be given a medical Question about a radiological chest X-ray
1325
         \hookrightarrow study (which is not provided).
1326
     Your task is to assess/rate the clarity of the Question, i.e. whether
1327
1328
         \hookrightarrow its wording is clear and unambiguous, and whether it is easy to
         \hookrightarrow understand and answer.
1329
1330
1331
     # Guidelines:
     - Use ONE of the rating options provided below, do NOT use any other
1332
        \hookrightarrow options
1333
      Answer with a short explanation (a few words) followed by "->" and
1334
         \hookrightarrow the rating option
1335
     - An example format will be provided
1336
1337
     - DO NOT REPEAT any part of the question or answer sentences in your
1338
         \hookrightarrow response!
1339
1340
     ## Rating Options -- ONLY USE ONE OF THE FOLLOWING OPTIONS
     - OPTIMAL: The question is mostly clear, unambiguous, and can be
1341
         \hookrightarrow answered. It is well-structured and concise without grammatical
1342
1343
            errors.
     - UNUSUAL_SENTENCE_STRUCTURE: The question is mostly clear,
1344
1345
         \hookrightarrow unambiguous, and can be answered. However, the sentence
1346
         \hookrightarrow structure is unusual or complex. There are no grammatical
         \hookrightarrow errors.
1347
1348
       GRAMMATICAL_ERRORS: The question is mostly clear, unambiguous, and
1349
         \hookrightarrow can be answered. However, there are grammatical errors that may
            affect the clarity. The sentence may or may not be well-
1350
        \hookrightarrow structured.
1351
     - UNRELATED_TO_CHEST_XRAY: The question is mostly clear and
1352
         \hookrightarrow unambiguous. However, it does not make sense to ask this
1353
         \hookrightarrow question about a chest X-ray study, because it does not relate
1354
         \hookrightarrow to the content that can be observed in a chest X-ray. There may
1355
         \hookrightarrow or may not be grammatical errors or unusual sentence structure
1356
1357
       UNCLEAR_QUESTION: The question may be misunderstood, is ambiguous,
1358
1359
         \hookrightarrow or otherwise unclear. Any answer could be misleading or
         \hookrightarrow incorrect, even with proper medical knowledge and context.
1360
        \hookrightarrow There may or may not be grammatical errors or unusual sentence
1361
        \hookrightarrow structure.
1362
1363
     Note that simply stating the indication/history motivating the study
1364
         \hookrightarrow is considered a valid question (and should not be rated as
1365
1366

→ UNCLEAR_QUESTION solely for not being an explicit question)!

1367
     # Example Format:
1368
     Question (to rate): ...
1369
1370
     Rating (your task):
1371
     The question is unrelated to chest X-rays -> UNRELATED_TO_CHEST_XRAY
1372
1373
1374
1375
     # Proceed with the following Question:
     Question (to rate): <QUEST>
1376
```

```
1377 Rating (your task):
```

For rating *answer clarity* (Listing 11), we condition the answer parts only (but not on the report or question) and we rate each answer part individually.

Listing 11: LLM prompt used for answer clarity evaluation of generated answers.

```
1382
1383
     You will be given a medical Answer to an unknown question about a
        \hookrightarrow radiological chest X-ray study (which is not provided).
1384
     Your task is to assess/rate the clarity of each sentence of the Answer
1385
         \hookrightarrow , i.e. whether its wording is clear and unambiguous, and
1386
         \hookrightarrow whether it is easy to understand.
1387
1388
     # Guidelines:
1389
      Rate each sentence in the Answer individually; do NOT use any prior
1390
         \hookrightarrow answer sentences as context or source
1391
1392
      Provide the rating for each answer sentences in its own line
        \hookrightarrow starting with the sentence number followed by your rating
1393
     - For each sentence, use ONE of the rating options provided below, do
1394
        \hookrightarrow NOT use any other options
1395
     - An example format will be provided
1396
     - DO NOT REPEAT the question or answer sentences in your response!
1397
1398
     ## Rating Options -- ONLY USE ONE OF THE FOLLOWING OPTIONS
1399
1400
     - OPTIMAL: The answer sentence is mostly clear and unambiguous. It is
1401
         \hookrightarrow well-structured and concise without grammatical errors.
      UNUSUAL_SENTENCE_STRUCTURE: The answer sentence is mostly clear and
1402
         \hookrightarrow unambiguous. However, the sentence structure is unusual or
1403
1404
         \hookrightarrow complex. There are no grammatical errors.
     - GRAMMATICAL_ERRORS: The answer sentence is mostly clear and
1405
1406
         \hookrightarrow unambiguous. However, there are (severe) grammatical errors
         \hookrightarrow that affect the clarity. The sentence may or may not be well-
1407
         \hookrightarrow structured.
1408
1409
      UNCLEAR ANSWER: The answer sentence may be misunderstood, is
         \hookrightarrow ambiguous, or otherwise unclear. There may or may not be
1410
         \hookrightarrow grammatical errors or unusual sentence structure.
1411
     - NOT_UNDERSTANDABLE: The answer sentence cannot be understood at all.
1412
         \hookrightarrow It is completely unclear, nonsensical, gibberish, or
1413
        \hookrightarrow contradictory in itself. There may or may not be grammatical
1414
1415
         \hookrightarrow errors or unusual sentence structure.
1416
     # Example Format:
1417
1418
     Answer (4 sentences to rate):
     [01] This first sentence.
1419
     [02] This is the second answer sentence.
1420
1421
     [03] Some text where it is unclear what is meant.
1422
     [04] This is the last answer sentence.
1424
     Rating (provide 4 ratings):
     [01] GRAMMATICAL_ERRORS
1425
     [02] OPTIMAL
1426
1427
     [03] UNCLEAR_ANSWER
     [04] OPTIMAL
1428
1429
1430
     # Proceed with the following Answer sentences:
1431
     Answer ( <NUMANS> sentences to rate):
1432
     <ANSWERS>
1433
1434
    Rating (provide <NUMANS> ratings):
1435
```

#### 437 E.4 Resources for Dataset Construction and Evaluation

#### 1438 E.4.1 Source Datasets

- MIMIC-CXR [14], [15], [21] We use the MIMIC-CXR dataset version 2.1.0 (https://physionet.org/content/mimic-cxr/2.1.0/ as the source of radiology reports from which we extract the scene graphs. It contains 227 835 radiographic (chest X-ray) studies performed at the Beth Israel Deaconess Medical Center in Boston, MA, USA. It is licensed under the PhysioNet Credentialed Health Data License 1.5.0.
- MIMIC-CXR-JPG [15], [24], [25] We use the MIMIC-CXR-JPG dataset version 2.1.0 (https: //physionet.org/content/mimic-cxr-jpg/2.1.0/) as the source of images for localization (the CXAS segmntation model is applied on these images). Additionally, we use the provided radiologist annotations (mimic-cxr-2.1.0-test-set-labeled.csv) as targets to evaluate the quality of extracted finding tags (Tabs. 2a and 4). The dataset is derived from MIMIC-CXR and is licensed under the PhysioNet Credentialed Health Data License 1.5.0.
- Chest ImaGenome [15], [17], [19] We use the Chest ImaGenome Dataset version 1.0.0 (https://physionet.org/content/chest-imagenome/1.0.0/) as a source of anatomical region bounding boxes for localization. Additionally, we use their provided scene graphs as a baseline for the evaluations of our scene graphs (Tabs. 2 and 4 to 7). It contains scene graphs for 242 072 frontal images from MIMIC-CXR that have been created using rule-based natural language processing and CXR atlas-based bounding box detection. The dataset is derived from MIMIC-CXR and is licensed under the PhysioNet Credentialed Health Data License 1.5.0.
- CXR-LT 2024 [15], [30], [34] We use the CXR-LT 2024 dataset version 2.0.0 (https://physionet.org/content/cxr-lt-iccv-workshop-cvamd/2.0.0/) as targets to evaluate the quality of extracted finding tags (Tabs. 2a and 5). More precisely, we use the gold standard dataset provided for Task 2 in the CXR-LT 2024 challenge tasks (406 reports, 26 classes). The dataset is derived from a small subset of MIMIC-CXR and was hand-labeled by radiologists. It is licensed under the PhysioNet Credentialed Health Data License 1.5.0.
- MS-CXR [15], [31], [35] We use the MS-CXR dataset version 1.1.0 (https://physionet.org/content/ms-cxr/1.1.0/) as targets to evaluate the quality of extracted finding boxes (Tabs. 2b and 6). The dataset contains 1162 image-sentence pairs of bounding boxes and corresponding phrases (and their finding classes) for 8 different findings. It is derived from a small subset of MIMIC-CXR and was hand-labeled by radiologists. It is licensed under the PhysioNet Credentialed Health Data License 1.5.0.
- REFLACX [15], [32], [36] We use the REFLACX dataset version 1.0.0 (https://physionet.org/content/reflacx-xray-localization/1.0.0/) as targets to evaluate the quality of extracted finding boxes (Tabs. 2b and 7). The dataset provides eye-tracking data collected for 3032 frontal chest x-rays from the MIMIC-CXR dataset. Additionally, it provides hand-labeled ellipses localizing for several anomalies present in the images. We only use the ellipses but do not use the eye-tracking data. It is licensed under the PhysioNet Credentialed Health Data License 1.5.0.

# 1475 E.4.2 Models

- Llama 3.1 70B [26] We use the AWQ-INT4 [40] quantized version of Llama 3.1 70B Instruct provided by the Huggingface hub at https://huggingface.co/hugging-quants/Meta-Llama-3. 1-70B-Instruct-AWQ-INT4. The model is derived from the https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct and is licensed under the LLAMA 3.1 COMMUNITY LICENSE AGREEMENT. We limit the maximum number of tokens to 6144.
- Llama 3.1 8B [26] We use the AWQ-INT4 [40] quantized version of Llama 3.1 70B provided by the Huggingface hub at https://huggingface.co/hugging-quants/Meta-Llama-3.1-8B-Instruct-AWQ-INT4. The model is derived from the https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct and is licensed under the LLAMA 3.1 COMMUNITY LICENSE AGREEMENT. We limit the maximum number of tokens to 8192.

CXAS [22], [23] We use the model provided by the CXAS Python library https: //pypi.org/project/cxas/. See also https://github.com/ConstantinSeibold/
ChestXRayAnatomySegmentation. It is licensed under the Attribution-NonCommercialShareAlike 4.0 International license. We run segmentation of all anatomical structures on half the original image resolution (half original image width and height).

BioLORD [29] We use the BioLORD-2023-C variant provided by the Huggingface model hub at https://huggingface.co/FremyCompany/BioLORD-2023-C and licensed under the MIT license. To apply the model, we use the Sentence Transformers library (https://github.com/UKPLab/sentence-transformers), which is licensed under the Apache-2.0 license.

Model Inference Details For all LLM-based information extraction steps, we rely on the vLLM library [41] (https://github.com/vllm-project/vllm, Apache-2.0 license) for inference. We run all models with temperature = 0.0. All json-outputs are parsed using the Pydantic libary (https://docs.pydantic.dev).

# 1499 E.4.3 Computational Costs

Each dataset construction step can run on an individual Nvidia A100 GPU, but we use multiple GPUs 1500 in parallel, with each GPU responsible for a different subset of the dataset. Semantic segmentation of 1502 all 158 anatomical structures using the CXAS models takes about 6 seconds per image, leading to a total of about 628 GPU hours. Sentence extraction takes about 1 second per study (report), leading to 1503 a total of about 65 GPU hours (for 227 835 studies). Observation extraction takes about 1.7 seconds 1504 per study, leading to a total of about 108 GPU hours. Indication extraction takes about 0.3 seconds per 1505 study, leading to a total of about 19 GPU hours. Scene graph construction (including entity matching) 1506 takes about 0.6 seconds per study, leading to a tool of about 38 hours. Question-answer generation 1507 1508 does not require a GPU but takes about 9 seconds per study (including all question templates and strategies), leading to a total of about 24 days. However, multiple processes can be run in parallel 1509 on a single machine, leading to an effective time of only about a day for all 42M QA-pairs. Quality 1510 assessment of OA texts again requires a GPU and consists of 5 individual steps that can be run in 1511 parallel. Overall the assessment takes about 6 GPU days for all 42M QA-pairs. 1512

# E.5 Dataset Release

1513

1523

1524

1525

1526

1527

1528 1529

1530

1531 1532

1533

1534

1535

We release the dataset as a credentialed dataset on the Physionet [15] platform (https://physionet. 1514 org/) under the PhysioNet Credentialed Health Data License 1.5.0 (https://physionet.org/ 1515 about/licenses/physionet-credentialed-health-data-license-150/). This makes the 1516 dataset openly accessible to all researches credentialed by PhysioNet, which requires a short online training. This type of hosting is required because we derived our dataset from the MIMIC-CXR [14] dataset. Additionally, this is also a responsible safeguard to protect the data that is (indirectly) derived 1519 from patient health data. While enabling researchers access to the dataset, it limits the access for 1520 other purposes. Additionally, it requires researchers to complete a privacy and ethics course. Code to 1521 construct the dataset and to train on it is made openly available. 1522

Societal Impact As a large vision-language dataset for medical imaging, this dataset has significant potential for societal impact. However, its use as a training source for models employed in clinical or medical applications also poses a substantial risk of misdiagnosis, highlighting the need for caution. Therefore, we strongly advise against relying solely on this dataset for fine-tuning or evaluating such models. On the other hand, this dataset can facilitate the development of large and interactive VQA models, which can provide supplemental information for patients, serve as a training tool for healthcare professionals, or optimize clinical workflows. The provided annotations, including bounding boxes and tags, further enhance its utility by providing a level of transparency and explainability in model predictions, allowing for more informed interpretation and analysis. By sparking research in this direction, this dataset can contribute to the advancement of the field and ultimately lead to positive long-term societal impacts. Nevertheless, it is essential to approach this dataset with caution, recognizing its limitations and potential risks if used improperly. As such, we consider this dataset a valuable research asset, but not yet suitable as a (sole) training source for real-world medical applications, emphasizing the need for careful evaluation and validation.

# F Structured VQA Task

1538

# F.1 Further Structured VQA Results

Table 13: Further results of our structured VQA task (Sec. 5). We show all the metrics from Tab. 3 with additional sub-metrics of our RadStrucVQA metric. Besides our default VQA model and the MAIRA-2 baseline, we also show alternative settings of our VQA model, namely training without bounding boxes and/or tags and predicting bounding boxes and tags before or after the text. Apart from these adaptions, the experimental setup was the same as in Sec. 5. We found that none of these adaptions has major influences on the results (apart from being capable of predicting boxes/tags), indicating that text, boxes, and tags in our dataset do not contradict each other. However, we observed minor improvements in text quality by adding bounding boxes

	Model		Ours (ablations)		Ours (default)	MAIRA-2 [18]	
	Boxes		/	/	<b>✓</b>		
	Tags	Х	X	×	✓	✓	✓
	Text	only text	after text	before text	after text	before text	after text
	Logical Prec	0.75	0.76	0.76	0.76	0.76	0.25
	Logical Rec	0.74	0.75	0.75	0.75	0.75	0.64
	Logical F1.	0.74	0.74	0.75	0.75	0.75	0.27
RadFact*	Grounding Prec	-	0.88	0.87	0.88	0.87	0.69
dFа	Grounding Rec	-	0.88	0.90	0.88	0.89	0.12
Ra	Grounding F1.	-	0.83	0.83	0.83	0.83	0.32
	Spatial Prec	-	0.68	0.67	0.68	0.67	0.12
	Spatial Rec	-	0.67	0.68	0.68	0.68	0.07
	Spatial F1.	-	0.63	0.63	0.64	0.63	0.06
	Finding Prec.	_	_	_	0.68	0.68	_
	Finding Rec.	-	-	_	0.67	0.66	_
	Finding F1	_	-	_	0.68	0.67	_
	Finding-pos Prec.	-	-	_	0.40	0.41	_
	Finding-pos Rec.	-	-	_	0.29	0.26	_
	Finding-pos F1	-	-	_	0.39	0.39	_
	Region Prec.	-	-	_	0.67	0.67	_
	Region Rec.	-	-	_	0.66	0.66	_
	Region F1	-	-	_	0.66	0.67	_
	Region-pos Prec.	-	-	_	0.29	0.34	_
(SS	Region-pos Rec.	-	-	_	0.21	0.21	_
RadStrucVQA (Tags)	Region-pos F1	-	-	_	0.29	0.32	_
Ą	Main-category Prec.	-	-	_	0.73	0.73	_
Š	Main-category Rec.	-	-	_	0.70	0.70	_
ĮĮ	Main-category F1	-	-	_	0.72	0.72	_
adS	Main-category-pos Prec.	-	-	_	0.49	0.52	_
23	Main-category-pos Rec.	-	-	_	0.36	0.34	_
	Main-category-pos F1	-	-	_	0.47	0.49	_
	Sub-category Prec.	-	-	_	0.71	0.71	-
	Sub-category Rec.	-	-	_	0.67	0.67	_
	Sub-category F1	-	-	_	0.69	0.69	_
	Sub-category-pos Prec.	-	-	-	0.47	0.50	_
	Sub-category-pos Rec.	-	-	_	0.34	0.32	-
	Sub-category-pos F1	-	-	_	0.45	0.46	-
	Bbox-pos-entity Prec.	-	-	-	0.31	0.32	-
	Bbox-pos-entity Rec.	-	-	-	0.22	0.20	_
	Bbox-pos-entity F1	-	-	_	0.26	0.26	_

\*Our RadStrucVQA implementation.

## 1539 F.2 RadStrucVQA Metric

Following the RadFact [18] metric, we split the predictions into individual elements. In our case, we treat each answer part as its own element, ignoring the hierarchy level and order. For each QA-sample, this results in a set of prediction elements  $\hat{\mathcal{Y}}$ , where  $|\hat{\mathcal{Y}}|$  is the number of answer parts in the predicted answer, and a set of target elements  $\mathcal{Y}$ , where  $|\mathcal{Y}|$  is the number of answer parts in the target answer. For each RadStrucVQ sub-metric (sub  $\in$  {logical, grounding, finding, . . . }), we compute a sample-level precision  $p_{\text{sub}}$  and recall  $r_{\text{sub}}$  score individually:

$$p_{\text{sub}}(\hat{\mathcal{Y}}, \mathcal{Y}) = s_{\text{sub}}(\hat{\mathcal{Y}}, \mathcal{Y}),$$
 (1)

$$r_{\text{sub}}\left(\hat{\mathcal{Y}}, \mathcal{Y}\right) = s_{\text{sub}}\left(\mathcal{Y}, \hat{\mathcal{Y}}\right),$$
 (2)

where  $s_{\mathrm{sub}}\left(\mathcal{H},\mathcal{C}\right)\in\left[0,1\right]$  is a sub-metric specific scoring function considering the hypothesis set  $\mathcal{H}$  given the context set  $\mathcal{C}$ . For precision  $\mathcal{H}=\hat{\mathcal{Y}}$  is the prediction set and  $\mathcal{C}=\mathcal{Y}$  is the target set, while for recall  $\mathcal{H}=\mathcal{Y}$  and  $\mathcal{C}=\hat{\mathcal{Y}}$ .

The score  $s_{\mathrm{sub}}$  is computed as the fraction of relevant hypothesis elements  $h \in \mathcal{H}$  that are entailed, using a sub-metric specific entailment definition, given the context  $\mathcal{C}$ . More precisely:

$$s_{\text{sub}}(\mathcal{H}, \mathcal{C}) = \frac{\left| \left\{ h \in \mathcal{H} \mid \text{entailed}_{\text{sub}}(h, \mathcal{C}[h]) \land \text{relevant}_{\text{sub}}(h) \right\} \right|}{\left| \left\{ h \in \mathcal{H} \mid \text{relevant}_{\text{sub}}(h) \right\} \right|}, \tag{3}$$

where  $\mathcal{C}[h]$  is the evidence from  $\mathcal{C}$  for h defined as

$$C[h] = \left\{ c \in C \middle| h \text{ is logically entailed with } C \land c \text{ provides evidence for } h \right\}. \tag{4}$$

We compute  $\mathcal{C}[h]$  by prompting an LLM to (i) identify entailment of h given all context elements in  $\mathcal{C}$ , where h can be ENTAILED or NOT\_ENTAILED (neutral or contradicting); and (ii) provide the relevant evidence for entailment, i.e. the context units  $c \in \mathcal{C}$  that support h. The LLM is given only the textual descriptions of each element (answer part), i.e. the entailment classification is purely logical and does not consider localization or any tags. Note that  $\mathcal{C}[h] = \{\}$  if h is not entailed.

Given the hypothesis h and its evidence C[h], the sub-metric entailment is computed individually by

entailed<sub>sub</sub>
$$(h, C[h]) \in \{\text{true}, \text{false}\},$$
 (5)

while the relevant subset of hypothesis elements is identified using the sub-metric specific

$$relevant_{sub}(h) \in \{true, false\}.$$
 (6)

The definitions for each sub-metric can be found in Tab. 14.

Table 14: RadStrucVQA sub metric definitions. The logical, grounding, and spatial sub-metrics follow the same principles as the corresponding sub-metrics in RadFact [18].

Sub-metric	$\mathrm{entailed_{sub}}\big(h,\mathcal{C}[h]\big)$	$\mathrm{relevant}_{\mathrm{sub}}(h)$	
logical	$\mathcal{C}[h]$ is not empty i.e. there is positive evidence for $h$ in $\mathcal C$ and $h$ does not contradict $\mathcal C$	always true	
grounding	$ \begin{array}{c} \operatorname{entailed_{logical}}\left(h,\mathcal{C}[h]\right) \wedge \operatorname{IoH}(h,\mathcal{C}[h]) \geq 0.5 \\ \text{where IoH is the Intersection between boxes in $h$ and boxes in $\mathcal{C}[h]$} \\ \operatorname{over the total box area in $h$,} \\ \text{with intersection/area computed based on box-masks (unions of boxes)} \\ \end{array} $	$h$ has bounding boxes $\land$ entailed $_{ ext{logical}}ig(h,\mathcal{C}[h]ig)$	
spatial	$\operatorname{entailed}_{\operatorname{grounding}} ig( h, \mathcal{C}[h] ig)$	h has bounding boxes	
finding	entailed <sub>logical</sub> $(h, \mathcal{C}[h]) \land$ each of the finding tags in $h$ is present in any of $\mathcal{C}[h]$ , only considering the subset of $\mathcal{C}[h]$ with the same positivity	h has finding tags	
finding-pos	$\mathrm{entailed}_{\mathrm{finding}}\big(h,\mathcal{C}[h]\big)$	$ ext{relevant}_{ ext{finding}}(h) \land h  ext{ is positive}$	
region	entailed <sub>logical</sub> $(h, \mathcal{C}[h]) \land$ each of the region tags in $h$ is present in any of $\mathcal{C}[h]$ , only considering the subset of $\mathcal{C}[h]$ with the same positivity	h has region tags	
region-pos	$\mathrm{entailed}_{\mathrm{region}}\big(h,\mathcal{C}[h]\big)$	$ ext{relevant}_{ ext{region}}(h) \land h  ext{ is positive}$	
main-category	entailed <sub>logical</sub> $(h, \mathcal{C}[h]) \land$ each of the finding main category tags in $h$ is present in any of $\mathcal{C}[h]$ , only considering the subset of $\mathcal{C}[h]$ with the same positivity	h has finding main category tags	
main-category-pos	$\mathrm{entailed_{main-category}}\big(h,\mathcal{C}[h]\big)$	$relevant_{main-category}(h) \land h$ is positive	
sub-category	entailed <sub>logical</sub> $(h, \mathcal{C}[h]) \land$ each of the finding sub category tags in $h$ is present in any of $\mathcal{C}[h]$ , only considering the subset of $\mathcal{C}[h]$ with the same positivity	h has finding sub category tags	
sub-category-pos	$\mathrm{entailed_{sub-category}}ig(h,\mathcal{C}[h]ig)$	$ ext{relevant}_{ ext{sub-category}}(h) \land h  ext{ is positive}$	
bbox-pos-entity	$\operatorname{entailed}_{\operatorname{finding}} ig( h, \mathcal{C}[h] ig) \wedge \operatorname{entailed}_{\operatorname{grounding}} ig( h, \mathcal{C}[h] ig)$	$\operatorname{relevant_{finding-pos}}(h) \land \\ \operatorname{relevant_{spatial}}(h)$	

**Implementation Details** The final precision/recall scores are computed by averaging the sample-level scores. F1 scores can also be computed by first taking the per-sample harmonic mean of precision and recall before averaging the sample-level F1 scores. Invalid answers, samples with LLM parse errors during entailment computation, as well as samples without relevant hypotheses are ignored during averaging. We use the same entailment prompts and few-shot examples as in RadFact [18] but use the Llama 3.1 8B [26] model, allowing us to compute the metric locally.

#### F.3 Experimental Setup

Vision-Language Model Training Our vision-language model follows the Llava architecture [37], using Rad-DINO [38] (microsoft/rad-dino) for image encoding and the 3B Llama 3.2 language model (https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct) connected via an MLP projection layer. We freeze the image encoder and all existing language model parameters but add new special tokens (with trainable embeddings) and apply LoRA [42] to the language model. Therefore, we only train the projection layer, the LoRA parameters, and the newly added token embeddings (keeping the existing token embeddings frozen). We train for one epoch on 1M samples from our CXR-QBA fine-tuning grade dataset (MIMIC-CXR's train split), where we use autoregressive training but only apply the loss to answer tokens. For image encoding and projection, we adopt the hyperparameters of MAIRA-2 [18]: We square-crop the images and resize them to  $518 \times 518$ , leading to  $37 \times 37 = 1369$  image patches (i.e. image tokens), then we use the features of the last image encoder layer, and project the image tokens using 4 projection layers with GeLU activations. For LoRA, we use r = 64,  $\alpha = 16$ , and dropout 0.05. The maximum number of tokens for the language model is restricted to 2048. We use the AdamW optimizer with cosine annealing scheduling with 500 warmup steps, maximum learning rate 1e - 3, no weight decay, a batch size

of 4 with 16 accumulation steps, gradient norm clipping at 1.0, and bf16 precision. The model is evaluated on the test split (following MIMIC-CXR) of our CXR-QBA fine-tuning grade set.

Prompt and Special Tokens Our question prompt follows the template shown in Listing 12, where <boi> (begin of image), <eoi> (end of image), and <imgref1> (first image reference) are newly added special tokens, <img> tokens are replaced by image token features, and {QUESTION} is replaced by the specific question.

# Listing 12: Question prompt.

```
Consider the following chest X-ray image: <br/>
Consider the following chest X-ray image: <br/>
Sequently consider the following chest X-ray image: <br/>
Consider the follow
```

The answers are formatted into sequences using XML-style structures and special tokens to represent tags and bounding boxes. An example is given in Listing 13.

Listing 13: Answer prompt.

```
1594
1595
    <answer>
       <regions><bilateral><lungs></regions>
1596
       obability > < certain > < neg > < probability >
1597
       <categories>
1598
         <ANATOMICAL_FINDING><DISEASE>
1599
         <subcat>LUNG FIELD</subcat><subcat>PULMONARY DISEASES</subcat>
1600
1601
       </categories>
       <entities><entity>pneumothorax</entity></entities>
1602
       <modifiers></modifiers>
1603
       <box><imgref1><x51><y18><x90><y87><box>
1604
       <box><imgref1><x09><y19><x52><y93></box>
1605
      No, there is no indication of pneumothorax.
1606
    </answer>
1682
```

1609

1611

1612

1613

1614

1615

1616

1617

1618

1619

1620

1621

1622

1623

1626

1627 1628

1629

1630

1631

1632

1633

We use special start and end tokens for answer parts (<answer> / </answer>), bounding //</categories>,<entities>,// / </modifiers>). For some tags we use individual special tokens, namely for laterality (e.g. <bilateral>), regions (e.g. <lungs>), certainty (e.g. <certain>), positivity (e.g. <neg>), and main categories (e.g. <ANATOMICAL\_FINDING>). For others we use start/end tokens and normal text, namely for sub-categories (<subcat> / </subcat>) and finding entities (<entity> / </entity>). Bounding boxes are listed after all other tags, where we use <box> / </box> tokens and refer back to the image using <imgref1>. Inside the box-tokens we use special relative coordinate tokens (following MAIRA-2 [18]) that represent the normalized  $(x_1, y_1, x_2, y_2)$  coordinates of the bounding box, each quantized to 100 different tokens per dimension. We use different tokens for the x- and y-dimensions but share them for both corners (e.g.  $x_1$  and  $x_2$  share the same token set). The textual description is the last part of each answer part and consists of plain text without special tokens. If an answer consists of multiple answer parts, then each answer part uses an individual block as in Listing 13. All new token embeddings are initialized close to the existing token embeddings, where we try to initialize them based on keywords defined for each token. More precisely, given a set of keywords for a new token, we tokenize the keywords using the old vocabulary and compute the average embedding of all these tokens. This is then used as the initialization for the new token.

MAIRA-2 Baseline We use the MIARA-2 [18] checkpoint available at https://huggingface.co/microsoft/maira-2. We freeze the full model but modify the prompt. More precisely, we use their original prompt for grounded report generation but slightly modify it, asking the model to answer to the question (included in the modified prompt) instead of reporting all findings in the image. The rest of the prompt is kept unchanged. This model is then evaluated on the same test set as our vision-language model. It is capable of generating individual answer parts, each with bounding boxes, but does not generate bounding boxes for negative answers and cannot generate any tags.

1634 **Computation Costs** We train on a single Nvidia A100 GPU (with 48GB of memory) for about 8 GPU days.

# NeurIPS Paper Checklist

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We claim to contribute a VQA dataset construction pipeline and a resulting VQA dataset with a specific size and answer structure and that we showcase the utility on an example task. These claims match the descriptions in Secs. 3 to 5.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss limitations in Sec. 6.2.

# Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
  only tested on a few datasets or with a few runs. In general, empirical results often
  depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

## 1688 Answer: [NA]

Justification: This work does not contribute theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

#### Answer: [Yes]

Justification: We describe the construction of the dataset in Sec. 3 with further details being provided in Appendices D and E, while we describe the example VQA task in Sec. 5 with further details in Appendix F. Additionally, we release (and provide it to the reviewers) the created dataset as well as the code to construct the dataset and to train on it.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

 Justification: We will release the dataset through the PhysioNet platform (https://physionet.org/). There the dataset will be openly accessible to all researches after being credentialed by PhysioNet, which requires a short online training. This type of hosting is required because we derived our dataset from the MIMIC-CXR dataset. Code to construct the dataset and to train on it is made openly available. See also Appendix E.5.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so No is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: An overview of the training procedure is provided in Sec. 5 with further details in Appendix F. Additional details are provided in the published code.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

#### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: For evaluation/results tables, we include the 95% confidence intervals over the samples in the (test) datasets, computed using bootstrapping with n=1000 (see Secs. 4 and 5).

#### Guidelines:

The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: For dataset construction, we provide details about compute resources in Appendix E.4.3. For training the VQA model, we describe the compute resources in Appendix F.3.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

# 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]
Justification:

# Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We briefly discuss societal impacts in Appendix E.5 and also mention potential use cases and limitations in Sec. 6.

# Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal
  impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: Discussed for dataset release in Appendix E.5.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We provide the links and URLs of all datasets and models/code used in this work in Appendix E.4.

# Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.

- The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The dataset contains a README-file describing the data structure. Additionally, we provide code to load and use the dataset with additional documentation. We also provide data construction and training code with READMEs on how to use them.

## Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This work does not directly involve human subjects and does not collect new human subject data. All human subject data used (such as chest X-rays or radiology reports) is derived from existing, public datasets, which where collected independently of and prior to this work.

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

# 1950 Answer: [NA]

Justification: This work does not directly involve human subjects and does not collect new human subject data. All human subject data used (such as chest X-rays or radiology reports) is derived from existing, public datasets, which where collected independently of and prior to this work.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

#### Answer: [Yes]

Justification: LLM usage during dataset construction and assessment is described in Secs. 3.1 and 3.3 with further details being provided in Appendix E. The usage as a component of the proof-of-concept model is described in Sec. 5 with further details in Appendix F.

#### Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.