

# Extracting and Encoding: Leveraging Large Language Models and Medical Knowledge to Enhance Radiological Text Representation

Pablo Messina<sup>1,4,5</sup>, René Vidal<sup>2</sup>, Denis Parra<sup>1,4,5</sup>, Álvaro Soto<sup>1,5</sup>, Vladimir Araujo<sup>3</sup>

<sup>1</sup>Pontificia Universidad Católica de Chile, <sup>2</sup>University of Pennsylvania, <sup>3</sup>KU Leuven

<sup>4</sup>Millennium Institute for Intelligent Healthcare Engineering (iHEALTH), Chile

<sup>5</sup>National Center for Artificial Intelligence (CENIA), Chile

{pamessina, dparras, vgaraujov}@uc.cl, vidalr@seas.upenn.edu,  
asoto@ing.puc.cl

## Abstract

Advancing representation learning in specialized fields like medicine remains challenging due to the scarcity of expert annotations for text and images. To tackle this issue, we present a novel two-stage framework designed to extract high-quality factual statements from free-text radiology reports in order to improve the representations of text encoders and, consequently, their performance on various downstream tasks. In the first stage, we propose a *Fact Extractor* that leverages large language models (LLMs) to identify factual statements from well-curated domain-specific datasets. In the second stage, we introduce a *Fact Encoder* (CXRFE) based on a BERT model fine-tuned with objective functions designed to improve its representations using the extracted factual data. Our framework also includes a new embedding-based metric (CXRFEScore) for evaluating chest X-ray text generation systems, leveraging both stages of our approach. Extensive evaluations show that our fact extractor and encoder outperform current state-of-the-art methods in tasks such as sentence ranking, natural language inference, and label extraction from radiology reports. Additionally, our metric proves to be more robust and effective than existing metrics commonly used in the radiology report generation literature. The code of this project is available at <https://github.com/PabloMessina/CXR-Fact-Encoder>.

## 1 Introduction

In the context of medical image analysis, radiology reports serve as a rich source of information. Radiologists routinely generate these free-text reports, which typically include sections such as *comparison*, *indication*, *findings*, and *impression*, as illustrated in Figure 1.

Radiology reports have been employed for various purposes, including label extraction for structured supervision in medical image tasks (Irvin et al., 2019; Wu et al., 2021; Jain et al., 2021a),



**Comparison:** Chest radiographs XXXX.  
**Indication:** XXXX-year-old male, chest pain.  
**Findings:** The cardiomeastinal silhouette is within normal limits for size and contour. The lungs are normally inflated without evidence of focal airspace disease, pleural effusion, or pneumothorax. Stable calcified granuloma within the right upper lung. No acute bone abnormality.  
**Impression:** No acute cardiopulmonary process.

Figure 1: Example image and report from the IU X-ray dataset (Demner-Fushman et al., 2015)

training data for models in report generation (Messina et al., 2022; Miura et al., 2021; Tanida et al., 2023), summarization tasks (Chen et al., 2023b; Ma et al., 2023), and the development of multimodal models capable of jointly understanding medical images and text (Wang et al., 2022; Boecking et al., 2022; Bannur et al., 2023).

A crucial aspect in addressing such tasks is the accurate comprehension of factual information within the report. Specifically, the *findings* and *impression* sections can be considered as repositories of factual information regarding the imaging examination. These factual statements encompass various elements, including observations (such as abnormalities, diseases, or devices), interpretations derived from observations, references to anatomical locations, discussions on severity or confidence levels, comparisons to previous studies, and more. For instance, in Figure 1, one factual statement indicates **no acute bone abnormality** (a normal observation), while another describes a **stable calcified granuloma within the right upper lung** (an abnormality found in a specific anatomical site).

Despite the various aforementioned applications that recent research has explored in the use of radiology reports, the persistent absence of effective methods for precise fact extraction and encoding from medical reports remains a critical challenge. As demonstrated in our experimental evaluation (Section 4), existing encoders and label extraction techniques frequently struggle to capture the nuanced details within free-text radiology reports, re-

sulting in incomplete or inaccurate depictions of clinical information. This is evident in various aspects. For instance, existing text encoders developed for the medical domain may struggle to generate consistent representations of paraphrased statements (Table 1) or to differentiate between similar sentences conveying contradictory meanings (Table 2), a crucial requirement to prevent encoding erroneous diagnoses. Similarly, current label extraction methods often rely on rigid labeling schemes based on manually crafted rules, leading to incomplete capture of all factual statements within a report (Table 5). Similar limitations are also observed in commonly used evaluation metrics for radiology text generation (Table 4).

In this work, we propose a novel approach that leverages the capabilities exhibited by Large Language Models (LLMs) such as ChatGPT, which have showcased outstanding performance in medical contexts (Liu et al., 2023c; Katz et al., 2023; Liu et al., 2023b; Adams et al., 2023), to improve factual statement representation. Our methodology also takes advantage of the existing knowledge in expert-annotated datasets. These datasets offer indispensable training data and also serve as a benchmark to enhance our model’s clinical terminology and context comprehension. Concretely, our contributions are three-fold:

- A fact extractor: a novel and simple approach to extracting facts that leverages ChatGPT and a fine-tuned version of T5 (Raffel et al., 2020) to capture relevant information from Chest X-ray radiology reports without requiring annotations from radiologists.
- A fact encoder: **Chest X-ray Fact Encoder (CXRFE)**, a CXR BERT-based model (Boecking et al., 2022) fine-tuned with a multi-task approach that leverages domain expertise from radiologists as well as ChatGPT and T5 generated annotations. CXRFE exhibits significant improvement in fact comprehension, demonstrated on sentence ranking and natural language inference tasks.
- A new evaluation metric for radiology text generation: **CXRFEScore**, which measures the factual accuracy of a generated text relative to a real text, by extracting and comparing the similarity of fact embeddings.

We release the weights of all our models, as well as the data and code necessary to replicate the

results. We also release CXRFEScore as a Python library for ease of use by the research community.

## 2 Related Work

In this section, we discuss prior work on BERT-based approaches to radiology text representation and label extraction from radiology reports, and leave discussion of prior work on evaluation of factual correctness in radiology text generation, applications of LLMs to medical text, and knowledge distillation from LLMs to Appendix A.1.

**BERT-based Approaches for Radiology Text Representation.** The advent of BERT (Devlin et al., 2019) has sparked notable progress in numerous NLP domains. This has inspired researchers to customize BERT for specific applications, including the medical field. Pioneering works such as BioClinicalBERT (Alsentzer et al., 2019), PubMedBERT (Gu et al., 2020), and BioLinkBERT (Yasunaga et al., 2022) have applied the masked language modeling (MLM) objective introduced by BERT to domain-specific corpora, such as PubMed paper abstracts and MIMIC-III (Johnson et al., 2016), an electronic health records dataset.

More recently, specialized variants like CXRBERT (Boecking et al., 2022) and BioViL-T (Banur et al., 2023) have been developed, targeting the unique challenges posed by CXR reports. CXRBERT provides both a general version, pretrained with MLM on PubMed abstracts and MIMIC-III documents, and a specialized version, fine-tuned with MLM coupled with a radiology section matching loss specifically tailored for reports from the MIMIC-CXR dataset (Johnson et al., 2019b). BioViL-T adopts the same pretraining strategy as CXRBERT but is subsequently fine-tuned using global and local multi-modal contrastive learning and image-informed MLM objectives. By combining reports with temporally sequenced image pairs, this approach enhances the understanding of radiological sentences with temporal descriptions.

Drawing inspiration from these works, we adopt BERT as our base model for text encoding. However, unlike prior approaches that aim to improve BERT’s representations with a single pre-text task (Reimers and Gurevych, 2019; Araujo et al., 2023), we employ a novel domain-specific multi-task learning protocol. This protocol leverages LLMs to generate large-scale supervision alongside expert-curated annotations from domain-specific datasets.

**Label Extraction from Radiology Reports.** Our work is also related to the problem of extracting information, usually in the form of labels, from free-text radiology reports. A well-known example in the literature is the CheXpert labeler (Irvin et al., 2019), which uses a rule-based system to infer the presence or absence of 13 observations (plus the label "No findings"). CheXbert (Smit et al., 2020) and VisualCheXbert (Jain et al., 2021b) are subsequent versions that follow the same labeling standard of CheXpert but are based on BERT.

The Chest ImaGenome (Wu et al., 2021) dataset is another example that used a rule-based NLP system to label reports to construct scene graphs for the corresponding frontal images in the MIMIC-CXR dataset (Johnson et al., 2019a). RadGraph (Jain et al., 2021a) proposed a labeling scheme of entities and relations for radiology reports and trained a variant of BERT, DyGIE++ (Wadden et al., 2019), for entity and relation extraction on examples annotated by radiologists. PadChest (Bustos et al., 2019) followed a similar approach, by labeling reports with an LSTM that was previously trained on examples annotated by physicians.

Our work contributes to this field by introducing a more flexible, open-vocabulary approach to information extraction, focused on extracting the essential factual information contained in the report, without imposing constraints that are too rigid. Specifically, we propose extracting factual statements, referred to as "facts," from reports, by leveraging the proven effectiveness of recent LLMs.

### 3 Method

We introduce a two-stage method for encoding the information within a CXR report. In the first stage, called fact extraction (Section 3.1), we utilize LLMs to extract facts from the original sentences of the report. In the second stage, called fact encoding (Section 3.2), we employ a BERT-based text encoder to generate sentence embeddings for each extracted fact. When used in tandem, these two stages form a cohesive system capable of producing vectorial representations of the factual statements found within a CXR report.

#### 3.1 Fact Extraction

Figure 2 outlines our method for extracting facts from radiology reports, with examples taken from the MIMIC-CXR dataset (Johnson et al., 2019b). Initially, we use regular expressions and simple

rules to pinpoint relevant radiological sections in MIMIC-CXR reports, mainly *Findings* and *Impression*, but we also handle alternate headings. These sections are then divided into sentences. For simplicity, we use the *sent\_tokenize* function from the NLTK library<sup>1</sup>. Next, we proceed to extract concise factual statements from each sentence. The main reason for doing this is that radiologists often write sentences that are noisy or complicated. Figure 2 shows two examples of such sentences. The first example contains multiple factual statements in one sentence, which can be simplified into shorter phrases. The second example is overly verbose, but the essential observation can be summarized in a brief phrase. We provide more examples in Table 6. Given the recent success of LLMs, an effective strategy to achieve this sort of extraction is by directing ChatGPT using a custom prompt.

**T5 as an alternative to ChatGPT.** In theory, this entire process could be executed using off-the-shelf LLMs. However, the expenses associated with accessing the API to annotate the entire dataset can be prohibitive. Therefore, the alternative approach we adopted was to annotate only a strategically selected subset of sentences and then transfer the acquired knowledge from these annotations to a more cost-effective sequence-to-sequence model, such as T5, through fine-tuning. This approach mirrors the strategy employed by Yang et al. (2023), where a T5 is fine-tuned to condense verbose descriptions from GPT-3 in LLM-assisted image classification.

We provide detailed implementation steps for this fact extraction procedure in Appendix A.2

#### 3.2 Fact Encoding

After we extract facts, we encode them to obtain vectors in a latent space via a text encoder model, called CXRFE. In this work, we rely on CXR-BERT (Boecking et al., 2022) to implement our fact encoder. Specifically, we use the CXR-BERT-specialized variant available on the Huggingface hub<sup>2</sup>, leveraging its built-in [CLS] token projection, which yields a 128-D vector serving as the final representation of the text.

Building on top of CXR-BERT-specialized, we explore 6 different training tasks to enhance the latent representation of radiological sentences: triplet loss for sentence ranking (T), natural language in-

<sup>1</sup><https://www.nltk.org/>

<sup>2</sup><https://huggingface.co/microsoft/BiomedVLP-CXR-BERT-specialized>

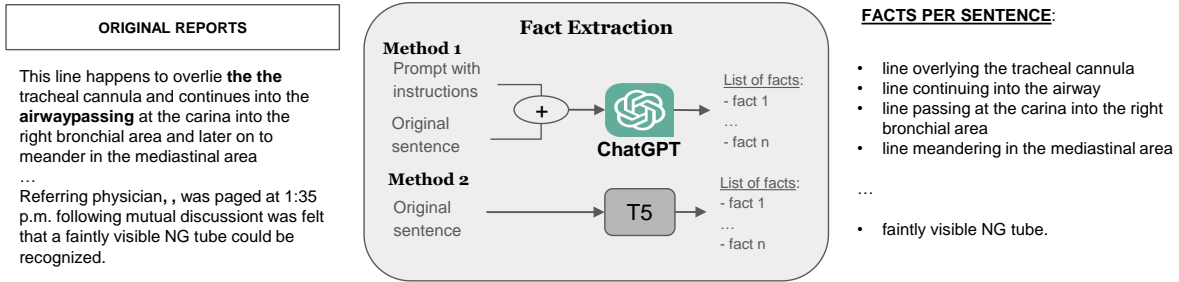


Figure 2: Fact extraction from radiology reports, designed to handle noisy input, including repeated words, incorrect sentence tokenization, typos, and verbose sentences, among other issues. When employing ChatGPT, context with instructions is used, whereas T5 is directly applied to the original sentence.

ference (NLI), quadruplet loss to enforce a separation between entailment and contradiction pairs (EC), entity and relation extraction (ER), sentence classification tasks (C), and sentence decoding (SD). We provide details on the implementation of each task in Appendix A.3.

Thus, by combining the two stages, the whole framework can accurately extract and encode facts from CXR reports, thus providing a rich and convenient representation of the factual information for downstream applications.

### 3.3 CXRFEScore

A potential application of our framework is in evaluating report generation from chest X-rays. We introduce CXRFEScore, an embedding-based metric that leverages both T5 and our CXRFE model for evaluation. Figure 3 illustrates how the metric works.

Given a reference report and a candidate report, we extract facts from each and represent them as embedding vectors, denoting the sets for the reference and candidate reports as  $R$  and  $C$  respectively. The cosine similarity matrix  $M$  of size  $|R| \times |C|$  is formed, where  $M_{i,j}$  represents the cosine similarity between the  $i^{th}$  vector of  $R$  and the  $j^{th}$  vector of  $C$ . This allows us to calculate  $S_{row}$ ,  $S_{col}$ , and the final CXRFEScore as follows:

$$S_{row} = \frac{\sum_i \max_j M_{i,j}}{|R|}$$

$$S_{col} = \frac{\sum_j \max_i M_{i,j}}{|C|}$$

$$CXRFEScore = \frac{S_{row} + S_{col}}{2}$$

The equations of this metric resemble those of BERTScore (Zhang et al., 2020a). The key difference lies in CXRFEScore’s comparison of fact embeddings rather than token embeddings. This metric illustrates the fusion of the two proposed

stages: fact extraction (implemented with T5) and fact encoding (implemented with CXRFE). We provide evidence of the robustness of this metric compared to many existing metrics in Section 4.

### 3.4 Dataset construction

In our experiments, we primarily utilize the MIMIC-CXR dataset (Johnson et al., 2019b), which comprises 227,827 radiology reports associated with 377,110 chest X-ray images. However, we focus solely on utilizing the reports for our experiments, deferring the exploration of images and multi-modality for future work. Additionally, we incorporate annotations from the Chest ImaGenome dataset (Wu et al., 2021), which provides scene graphs linking observations to anatomical image locations for each frontal view image in MIMIC-CXR. These annotations serve two main purposes: facilitating our creation of a binary multi-label classification task and introducing a radiologist-informed annotation standard covering various observation types and anatomical locations. Similarly, we utilize RadGraph (Jain et al., 2021a), which offers an entity and relation annotation scheme for radiology reports in MIMIC-CXR, and datasets such as MedNLI (Romanov and Shivade, 2018), RadNLI (Miura et al., 2021) and MS-CXR-T (Bannur et al., 2023) for experiments on Natural Language Inference (NLI). To assess the performance and generalization ability of our two-stage framework, we also evaluate CXRFEScore using the 3955 reports and associated tag annotations of the IU X-ray dataset (Demner-Fushman et al., 2015). It is important to highlight that radiologists or doctors have partially or fully annotated these datasets, which adds significant value to training and evaluation. We direct the reader to Appendix A.3 for further details on these datasets



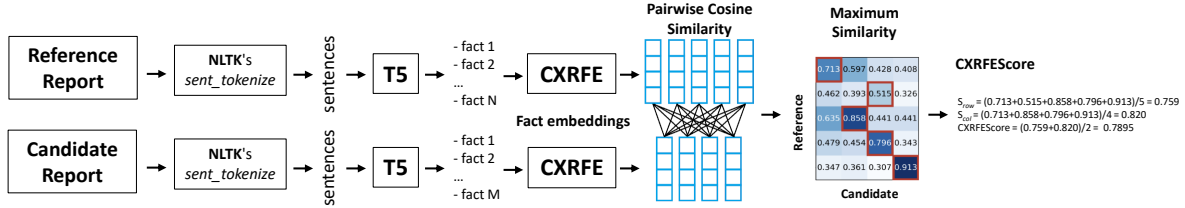


Figure 3: Illustration of the computation of CXRFEScore. Given a reference report and a candidate report, we employ NLTK’s *sent\_tokenize* to extract sentences from each, extract facts from these sentences using T5, generate embeddings from the facts using CXRFE, and finally compute the pairwise cosine similarity matrix. Greedy matching is highlighted in red, with the calculation of the final score explicitly shown on the right.

and their utilization in our experiments.

### 3.4.1 LLM-assisted data augmentations for training supervision

As mentioned in Section 3.2, we use six distinct tasks to train CXRFE. For certain tasks, we leverage ChatGPT to generate additional training data. We elaborate in depth on these aspects in Appendix A.3, so here we only offer a concise overview.

For the triplet loss task (detailed in Section 3.4.2), some triplets incorporate paraphrased facts, which we generate using ChatGPT (an example is depicted in Figure 15). Similarly, we employ ChatGPT to produce challenging triplets, as indicated by the prompt in Figure 21.

In classification-related tasks, each fact is annotated with a JSON metadata object using ChatGPT. This object encompasses fields such as *anatomical location*, *detailed observation*, *short observation*, *category*, *health status*, and *comparison status*. The respective prompt for this process is shown in Figure 17. To refine the "comparison status" field, we use another prompt showcased in Figure 18. Also, we employ ChatGPT to label according to the scheme of the Chest ImaGenome dataset, producing additional observation and anatomical location labels, as shown in Figures 19 and 20.

For the task of natural language inference, we extensively utilize GPT-4 to generate training examples with distinct prompts (refer to Figures 22, 23, 24, 25, and 26).

### 3.4.2 Triplet Sampling

CXRFE is trained to generate sentence embeddings that cluster semantically similar sentences in the embedding space through a triplet ranking task with binary cross-entropy loss. This approach uses a dataset of triplets, each one with an anchor, a positive sample (akin to the anchor), and a negative one. The difference in similarities is computed as

$\Delta\text{sim}(a, p, n) = \text{sim}(a, p) - \text{sim}(a, n)$  from their embeddings’ dot product. By minimizing the binary cross-entropy loss, the encoder ensures closely related sentences are nearer and unrelated ones are more distant in the embedding space.

We define six triplet sampling rules to guide the selection process. **Rule 1** prioritizes paraphrases generated with ChatGPT. **Rule 2** involves sampling triplets based on the consensus of BioVIL-T and Levenshtein distance, with the anchor and positive sample sharing the same health status. **Rule 3** ensures proximity between short observations, detailed observations, and original facts, along with their paraphrases. **Rule 4** samples triplets based on Chest ImaGenome labels, ensuring that the anchor and positive sample share at least one label and that BioVIL-T and Levenshtein distance agree. **Rule 5** ranks triplets according to the overlap of entities and relations from RadGraph. **Rule 6** includes hard triplets generated by ChatGPT. For each rule, we sampled around 3 to 4 million training triplets and 1,000 each for validation and testing.

These rules encapsulate specific intuitions and heuristics regarding the ranking of sentence embeddings within the semantic space. The design of these sampling rules and the construction of the triplets dataset involve several technical details, which are elaborated upon in Appendix A.4.

### 3.4.3 Natural Language Inference

Natural Language Inference (NLI) aims to classify the relationship between a premise and a hypothesis into one of three categories: "entailment", "neutral", or "contradiction". For instance, consider a premise stating “*There are no evident signs of pleural effusion*”, while a hypothesis asserts “*There are evident signs of pleural effusion*”. Despite their structural and lexical similarities, these sentences express contradictory meanings. A robust

Table 1: Triplet and sentence ranking results. **Triplet ranking:** 1000 samples per rule. For each rule, we report the fraction of correctly ranked triplets. (\*) denotes perfect scores achieved by BioViL-T due to an unfair advantage from how triplets were heuristically sampled in Rules 2 and 4. (o) stands for observations and (a) for anatomical locations. **Sentence ranking:** 2412 sentences annotated by radiologists, sourced from the Chest ImaGenome Gold dataset. Notation: a@k represents the mean average accuracy of the top k ranked sentences (larger is better), while c@k represents the mean number of contradictory sentences among the top k ranked sentences (smaller is better).

ID	Text Model	Triplet Ranking							Sentence Ranking				
		R1(o)	R1(a)	R2	R3	R4	R5	R6	AUC	a@50	a@100	c@50	c@100
1	BioLinkBERT	0.753	0.725	0.786	0.756	0.644	0.774	0.520	0.717	0.951	0.945	5.523	12.158
2	PubMedBERT	0.901	0.853	0.905	0.873	0.767	0.834	0.603	0.775	0.954	0.947	4.552	10.134
3	BioClinicalBERT	<b>0.922</b>	<b>0.864</b>	0.933	0.912	0.834	<b>0.948</b>	0.601	0.830	0.957	0.950	3.823	8.615
4	CheXbert	0.855	0.771	0.908	0.884	0.760	0.937	0.635	<b>0.864</b>	<b>0.962</b>	<b>0.955</b>	<b>1.914</b>	<b>4.299</b>
5	CXR-BERT-specialized	0.880	0.804	0.992	0.914	0.904	0.932	0.717	0.837	0.953	0.947	2.905	6.230
6	BioViL-T	0.910	0.851	<b>1.000*</b>	<b>0.938</b>	<b>1.000*</b>	0.944	<b>0.765</b>	0.845	0.956	0.949	3.158	6.903
7	CXRFE(T)	0.968	0.955	0.925	0.964	0.798	0.952	0.946	0.896	0.963	0.957	1.668	3.940
8	CXRFE(T+C)	0.967	0.945	0.967	0.982	<b>0.926</b>	<b>0.988</b>	0.937	<b>0.919</b>	<b>0.975</b>	<b>0.967</b>	2.955	9.157
9	CXRFE(T+ER)	0.962	0.946	0.917	0.961	0.798	0.954	0.927	0.888	0.964	0.957	1.403	3.543
10	CXRFE(T+SD)	<b>0.981</b>	<b>0.966</b>	0.954	0.977	0.875	0.981	0.898	0.897	0.961	0.955	2.427	5.465
11	CXRFE(T+EC)	0.957	0.953	0.950	0.971	0.809	0.965	0.943	0.840	0.957	0.951	1.237	2.522
12	CXRFE(T+NLI)	0.910	0.898	0.958	0.970	0.861	0.968	0.903	0.820	0.951	0.945	1.378	3.240
13	CXRFE(T+EC+NLI)	0.928	0.950	0.923	0.961	0.777	0.934	0.925	0.812	0.950	0.943	<b>0.969</b>	<b>2.478</b>
14	CXRFE(T+C+EC+NLI)	0.971	0.932	0.974	0.980	0.892	0.982	0.945	0.890	0.969	0.960	1.437	3.789
15	CXRFE(T+C+EC+NLI+ER)	0.972	0.944	<b>0.978</b>	<b>0.983</b>	0.911	0.984	0.936	0.917	0.972	0.964	1.719	4.694
16	CXRFE(T+C+EC+NLI+SD)	0.977	0.953	0.972	0.980	0.883	0.980	0.951	0.906	0.971	0.962	1.445	3.603
17	CXRFE(T+C+EC+NLI+ER+SD)	0.976	0.958	0.976	0.982	0.880	0.980	<b>0.952</b>	0.907	0.970	0.962	1.502	3.804

sentence embedding should be able to distinguish between sentences conveying contradictory diagnoses. Thus, our objective in leveraging NLI during training is to refine sentence embeddings to discern such subtle distinctions accurately.

All MedNLI splits (Romanov and Shivade, 2018) are used for training, amounting to 14,049 annotated sentence pairs. Radiology-specific datasets include RadNLI (Miura et al., 2021) with 960 pairs and MS-CXR-T (Bannur et al., 2023), an evaluation set with 361 pairs emphasizing condition evolution over time. Given the dearth of training data, we build a custom NLI dataset combining MedNLI (14,049), RadNLI development set (480) and GPT-4 generated examples (154,498), resulting in a total of 169,025 pairs (entailment: 25,175, neutral: 44,729, contradiction: 99,121), while RadNLI test set (480) and MS-CXR-T (361) are set apart for evaluation.

## 4 Experimental Results

Our experiments evaluate different versions of CXRFE, each defined by a subset of the six tasks (T, C, NLI, EC, ER, SD) outlined in Section 3.2. This results in a total of 64 potential combinations. However, for the sake of simplicity, we heuristically assess only 11 combinations. Further details on our rationale are provided in Appendix A.3.

**Triplet Ranking.** We evaluate CXRFE and mul-

tiples from the literature on triplet ranking accuracy, using a separate test set of 1000 triplets per rule sampled according to the sampling rules detailed in Section 3.4.2. The left side of Table 1 presents these results. Notably, all different versions of CXRFE outperform all the baselines in triplet rules where ChatGPT is heavily involved, namely, paraphrases (R1(o), R1(a), R3) and hard triplets (R6). The hard triplets are especially challenging for the baselines, with BioViL-T only achieving 0.765 accuracy (row 6), while the best performing version of CXRFE achieves 0.952 (row 17). In addition to triplet loss (T), which is key to learning an embedding for these rules, we notice that sentence decoding (SD) and classification (C) appear to be helpful auxiliary tasks since most of the best scores are achieved by variants that include them (rows 8, 10, 15, 17).

**Sentence Ranking.** To complement the triplet ranking evaluation, which is based on heuristic sampling rules, we conduct a sentence ranking evaluation using 2,412 carefully annotated sentences provided by radiologists from the gold standard of the Chest ImaGenome dataset (Wu et al., 2021). These sentences are annotated with a vocabulary of 70 observations (*yes* (1), *no* (0), or *omitted* (-1)), and 38 anatomical locations (*mentioned* (1) or *unmentioned* (0)), resulting in a discrete vector of size 108 for each sentence.

Table 2: NLI results as sentence similarity.  $a_{E+C}^*$  denotes the upper bound accuracy with optimal similarity threshold tuning on the same evaluation data. BT stands for best threshold (tuned in NLI custom dataset).

ID	Text Model	NLI custom dataset Ent. pairs: 25,175 Cont. pairs: 99,121				RadNLI test set Ent. pairs: 102 Cont. pairs: 98				MS-CXR-T Ent. pairs: 141 Cont. pairs: 220			
		BT	$a_E$	$a_C$	$a_{E+C}$	$a_E$	$a_C$	$a_{E+C}$	$a_{E+C}^*$	$a_E$	$a_C$	$a_{E+C}$	$a_{E+C}^*$
1	PubMedBERT	1.000	1.9	99.8	50.8	0.0	100.0	50.0	63.8	0.0	100.0	50.0	56.5
2	BioLinkBERT-large	0.999	1.9	100.0	50.9	0.0	100.0	50.0	65.9	0.0	100.0	50.0	58.1
3	BioClinicalBERT	1.000	1.9	100.0	50.9	0.0	100.0	50.0	69.2	0.0	100.0	50.0	69.5
4	CheXbert	0.556	88.3	33.8	61.0	96.1	83.7	<b>89.9</b>	<b>90.4</b>	100.0	1.8	50.9	63.3
5	CXR-BERT-specialized	0.619	76.9	44.6	60.8	73.5	85.7	79.6	82.3	100.0	8.6	<b>54.3</b>	77.5
6	BioViL-T	0.709	70.9	51.4	<b>61.2</b>	58.8	89.8	74.3	78.0	100.0	7.7	53.9	<b>87.8</b>
7	CXRFE(T)	0.640	79.8	57.2	68.5	69.6	89.8	79.7	87.3	100.0	18.6	59.2	78.0
8	CXRFE(T+C)	0.934	43.3	69.9	56.6	44.1	95.9	70.0	75.4	97.9	12.3	55.1	62.6
9	CXRFE(T+ER)	0.606	81.6	54.3	67.9	70.6	89.8	80.2	85.9	100.0	16.8	58.4	78.5
10	CXRFE(T+SD)	0.620	68.1	50.5	59.3	81.3	79.6	80.5	81.9	99.3	10.9	55.1	70.3
11	CXRFE(T+EC)	0.308	97.0	93.7	<b>95.3</b>	98.0	93.8	96.0	96.4	96.5	69.5	83.0	93.3
12	CXRFE(T+NLI)	0.233	89.3	78.6	84.0	99.0	93.9	<b>96.4</b>	<b>97.0</b>	100.0	75.9	87.9	95.7
13	CXRFE(T+EC+NLI)	0.267	96.7	93.3	95.0	98.0	93.9	96.0	<b>97.0</b>	99.3	84.1	<b>91.7</b>	97.6
14	CXRFE(T+C+EC+NLI)	0.288	94.9	82.4	88.6	99.0	91.8	95.4	96.0	100.0	75.9	88.9	97.2
15	CXRFE(T+C+EC+NLI+ER)	0.431	95.2	83.0	89.1	99.0	89.8	94.4	95.5	100.0	59.0	79.5	<b>98.5</b>
16	CXRFE(T+C+EC+NLI+SD)	0.480	96.6	87.3	92.0	99.0	89.8	94.4	95.5	100.0	70.5	85.2	93.6
17	CXRFE(T+C+EC+NLI+ER+SD)	0.455	97.1	87.0	92.0	98.0	91.8	94.9	96.0	100.0	72.7	86.4	96.0

We report the results of this evaluation on the right side of Table 1. In this evaluation, each sentence serves as a query, against which all other sentences are ranked based on the cosine similarity of their embeddings. AUC,  $acc@k$ , and  $cont@k$  are computed by comparing the labels of the query sentence with those of each ranked sentence. AUC requires defining which sentences are relevant (1) or irrelevant (0) for the query. For two sentences  $a$  and  $b$ ,  $a$  is deemed relevant for  $b$  if  $a$ 's labels logically entail  $b$ 's or vice versa; otherwise, they are considered irrelevant.  $acc@k$  represents the mean average accuracy up to the  $k$ th sentence in the ranking, while  $cont@k$  represents the mean number of sentences contradicting the query up to the  $k$ th sentence, by having contradictory values in at least one observation (1 vs. 0).

CheXbert (row 4) consistently performs the strongest among the baselines, likely due to its training in a similar multi-label classification task. CXRFE achieves the best overall performance. Notably, the T+C variant (row 8), combining triplet loss (T) and classification (C), attains the highest AUC and accuracy, while T+EC+NLI (row 13), containing tasks designed for pulling contradictory sentences apart, yields the smallest  $cont@k$  scores (where smaller is better). Rows 14-17 represent intermediate points between these two extremes.

**NLI.** Table 2 presents the NLI results using cosine similarity between sentence vectors, following a methodology akin to Bannur et al. (2023). This methodology specifically focuses on *entailment*

Table 3: RadNLI test set accuracy. Results for CXR-BERT, IFCC, PTUnifier and DoT5 are from the original papers. CoT stands for Chain-of-Thought prompting.

ID	Text Model	Accuracy
1	CXR-BERT (Boecking et al., 2022)	65.2
2	IFCC (Miura et al., 2021)	77.8
3	PTUnifier (Chen et al., 2023a)	80.0
4	DoT5 (Liu et al., 2023a)	<b>82.1</b>
5	GPT-4 w/ simple prompt	82.3
6	GPT-4 w/ CoT + examples	<b>89.0</b>
7	Meta-Llama-3-8B w/ simple prompt	58.1
8	Meta-Llama-3-8B w/ CoT + examples	61.5
9	CXRFE(T+NLI)	84.2
10	CXRFE(T+EC+NLI)	81.3
11	CXRFE(T+C+EC+NLI)	86.7
12	CXRFE(T+C+EC+NLI+SD)	85.2
13	CXRFE(T+C+EC+NLI+ER)	88.1
14	CXRFE(T+C+EC+NLI+ER+SD)	88.5
15	CXRFE(NLI fine-tuning)	<b>89.8</b>

and *contradiction* pairs, aiming to assess the efficacy of a text embedding in distinguishing between the two given a similarity threshold. We present results across three datasets: our NLI custom dataset mentioned in Section 3.4.3, the RadNLI test set, and MS-CXR-T. Our reported results are based on thresholds fine-tuned in the NLI custom dataset, alongside upper bounds obtained by tuning thresholds within the same data used for evaluation.

Notably, employing the quadruplet entailment/contradiction loss (EC) and natural language inference (NLI) (rows 11 to 17) is essential to achieve high performance, significantly outperforming all baselines. In contrast, variants lacking EC and NLI (rows 7 to 10) exhibit a weaker result.

Table 4: Comparative evaluation of text generation metrics. Notation: a@k denotes the mean average accuracy up to the  $k$ th sentence, c@k represents the mean number of contradictory sentences up to the  $k$ th sentence, and j@k represents the mean average Jaccard index up to the  $k$ th sentence.

ID	Metric	Chest ImaGenome Gold Sentences (2412)					Chest ImaGenome Gold Reports (500)		IU X-Ray Reports (3955)		RadNLI + MS-CXR-T
		AUC	a@50	a@100	c@50	c@100	j@20	j@50	j@20	j@50	AUC
1	BLEU (Papineni et al., 2002)	0.767	0.951	0.944	4.432	9.148	0.509	0.469	0.392	0.336	0.537
2	ROUGE-L (Lin, 2004)	0.773	0.953	0.946	4.228	9.029	0.508	0.466	0.391	0.335	0.547
3	METEOR (Banerjee and Lavie, 2005)	0.829	0.954	0.947	4.282	10.100	0.514	0.471	0.390	0.328	0.596
4	CIDEr-D (Vedantam et al., 2015)	0.778	0.952	0.946	4.883	10.234	0.482	0.441	0.415	0.366	0.556
5	BERTScore (Zhang et al., 2020a)	<b>0.840</b>	0.960	0.952	5.012	11.328	0.531	0.489	<b>0.429</b>	<b>0.383</b>	0.559
6	CheXpert Accuracy (Irvin et al., 2019)	0.764	0.941	0.939	3.263	6.263	0.467	0.446	0.375	0.322	0.592
7	CheXpert F1 (Irvin et al., 2019)	0.742	0.939	0.938	3.658	6.877	0.460	0.441	0.312	0.327	0.582
8	CheXbert Accuracy (Smit et al., 2020)	0.778	0.941	0.941	<b>2.667</b>	<b>5.103</b>	0.485	0.452	0.356	0.357	0.592
9	CheXbert F1 (Smit et al., 2020)	0.753	0.939	0.939	2.709	5.301	0.472	0.445	0.365	0.369	0.583
10	RadGraph F1 (Full) (Jain et al., 2021a)	0.831	<b>0.961</b>	<b>0.953</b>	2.881	7.189	0.546	0.500	<b>0.429</b>	0.371	<b>0.610</b>
11	RadGraph F1 (Partial) (Jain et al., 2021a)	0.789	0.960	0.951	3.453	8.833	<b>0.549</b>	<b>0.501</b>	0.415	0.361	0.574
12	CXRFEScore(T+C+EC+NLI)	0.912	0.967	0.958	1.563	3.953	0.562	0.510	0.504	0.460	0.927
13	CXRFEScore(T+C+EC+NLI+ER)	<b>0.921</b>	<b>0.968</b>	<b>0.959</b>	1.575	4.122	<b>0.563</b>	<b>0.511</b>	0.513	0.471	0.931
14	CXRFEScore(T+C+EC+NLI+SD)	0.911	0.967	0.958	<b>1.518</b>	<b>3.763</b>	0.553	0.505	0.517	0.471	0.936
15	CXRFEScore(T+C+EC+NLI+ER+SD)	0.911	0.967	0.958	1.606	3.874	0.555	0.507	<b>0.518</b>	<b>0.473</b>	<b>0.938</b>

Table 5: Template-based report generation metrics on MIMIC-CXR test set (3269 reports) for different label extraction methods. CXRFEScore was calculated using the T+C+EC+NLI+ER+SD variant. Notation: FE = Fact Extraction; BS-F1 = BERTScore F1; B-4 = BLEU-4; C-D = CIDEr-D; R-L = ROUGE-L; MET = METEOR.

ID	Label Extraction Method	CXRFEScore	RadGraph F1		CheXpert F1		CheXbert F1		BS-F1	B-4	C-D	R-L	MET
			Full	Partial	Micro	Macro	Micro	Macro					
1	CheXpert labeler	0.644	0.119	0.161	<b>0.998</b>	<b>0.990</b>	0.939	0.854	<b>0.470</b>	<b>0.007</b>	<b>0.023</b>	<b>0.123</b>	<b>0.179</b>
2	CheXbert	0.647	<b>0.120</b>	0.162	0.948	0.921	<b>0.983</b>	<b>0.907</b>	<b>0.470</b>	<b>0.007</b>	<b>0.023</b>	<b>0.123</b>	<b>0.179</b>
3	Chest ImaGenome	<b>0.677</b>	0.104	<b>0.237</b>	0.767	0.693	0.776	0.751	0.257	0.003	0.002	0.086	0.170
4	FE (T5-small)	<b>0.983</b>	<b>0.784</b>	0.745	<b>0.973</b>	<b>0.964</b>	<b>0.959</b>	<b>0.947</b>	<b>0.789</b>	0.275	0.672	<b>0.533</b>	<b>0.650</b>
5	FE (GPT-4)	0.937	0.773	<b>0.751</b>	0.969	0.959	0.957	0.946	0.777	0.275	0.853	0.503	0.635
6	FE (Meta-Llama-3-8B)	0.922	0.737	0.738	0.944	0.940	0.932	0.925	0.775	<b>0.286</b>	<b>0.902</b>	0.522	0.636

Among the baselines, CheXbert (row 4) demonstrates superior performance on RadNLI, while BioVil-T (row 6) is the clear victor on MS-CXR-T. However, all baselines struggle considerably in our NLI custom dataset and are outperformed across all three datasets by variants 11-17 of CXRFE.

Additionally, Table 3 presents the accuracy achieved on the RadNLI test set in the context of the typical 3-class classification task encompassing entailment, contradiction, and neutral classes. In this evaluation, we exclusively assess variants of CXRFE equipped with an NLI classification head (rows 9-15). For insights into implementing NLI classification, please consult Figure 7.

Within the existing literature, the strongest baseline identified is DoT5 (Liu et al., 2023a) (82.1), employing a sophisticated sequence-to-sequence approach based on T5. Furthermore, we conducted evaluations on GPT-4 (rows 5-6) and Meta-Llama-3-8B (AI@Meta, 2024) (rows 7-8) utilizing two distinct prompts: a simple prompt (Figure 24) and a prompt with Chain-of-Thought (CoT) + exam-

ples (Figure 25). Notably, the second prompt led to a significant performance boost for GPT-4 (from 82.3 to 89.0), whereas Meta-Llama-3-8B, an open-source LLM from Meta, only experienced a moderate improvement (from 58.1 to 61.5), with very low accuracies overall. Consequently, GPT-4 with the second prompt was selected as our "oracle" for generating additional training data (more details on this in Appendix A.3).

Most versions of CXRFE showed superior performance compared to the baselines. Surprisingly, a version fine-tuned explicitly for NLI (row 15) even outperforms GPT-4 with CoT (row 6) by a narrow margin (89.8).

**CXRFEScore vs. existing metrics.** To assess the quality of our proposed metric, we conduct an evaluation of CXRFEScore alongside multiple metrics from the literature, as shown in Table 4. This assessment encompasses four settings: (1) a sentence ranking evaluation using 2412 sentences, (2) a report ranking evaluation with 500 reports, both sourced from the gold standard of Chest Im-



aGenome, (3) a report ranking evaluation with 3955 reports from the IU X-ray dataset (Demner-Fushman et al., 2015) leveraging the manual and automatic tags associated with each report, and (4) a natural language inference evaluation utilizing entailment (336) and contradiction (424) pairs from RadNLI and MS-CXR-T. Note that all these datasets are annotated by radiologists, thus serving as gold standards for metric comparison.

Among the baseline metrics, RadGraph F1 (rows 10-11) emerges as one of the most promising based on its performance on Chest ImaGenome Gold, which aligns with the findings of Yu et al. (2022). However, BERTScore’s results (row 5) on Chest ImaGenome Gold are quite similar, achieving the highest AUC among the baselines (0.840). Additionally, BERTScore achieves the highest Jaccard index scores on IU X-ray among the baselines. Notably, CheXbert (rows 8-9), closely followed by CheXpert (rows 6-7), shows the fewest contradictions on Chest ImaGenome Gold.

All baseline metrics, however, are surpassed by CXRFEScore (rows 12-15) in all the evaluation metrics. A particularly striking observation is that the baseline metrics struggle significantly to differentiate between entailed and contradictory sentences, as indicated by the AUC results in the last column of Table 4. RadGraph F1 Full achieves an AUC of only 0.610, whereas the best version of CXRFEScore (row 15) achieves an AUC of 0.938. This suggests that current metrics assign elevated scores to pairs of sentences with contradictory semantics, highlighting the necessity for improved metrics to discern these subtleties—precisely what CXRFEScore is designed to accomplish.

Similarly, CXRFEScore outperforms all the baselines on IU X-ray, a dataset not used to develop the metric. This provides valuable evidence of the metric’s ability to generalize to radiology reports from a different institution.

We provide additional details and results about these metrics in Appendix A.5.

**Fact Extraction Quality.** To evaluate the quality of the fact extraction stage, we interpret these facts as *open-vocabulary* labels and compare them against three existing radiology report label extraction methods: CheXpert labeler (Irvin et al., 2019), CheXbert (Smit et al., 2020), and Chest ImaGenome (Wu et al., 2021). For Chest ImaGenome, we use the labels from the dataset’s scene graphs, as the original NLP algorithm is not publicly available. For fact extraction, we compare T5-small,

fine-tuned specifically for this task, against GPT-4 and Meta-Llama-3-8B. The latter two models use the prompt shown in Figure 14. Our evaluation protocol involves the following: for each MIMIC-CXR test set report and label extraction method, labels are extracted, converted into a report, and evaluated against the original report by several metrics. We adopt the templates suggested by Pino et al. (2021) for CheXpert labeler and CheXbert, while for Chest ImaGenome, we utilize basic templates such as “(no) {observation} in {anatomical location}”. For fact extraction, we simply concatenate the facts. These template-based reports are illustrated in Tables 7 and 8.

Table 5 presents these results. Notably, fact-based reports generated by T5 (row 4) achieve most of the top scores, even slightly outperforming GPT-4 and Meta-Llama-3-8B, demonstrating the efficacy of the fact extraction process. Interestingly, CXRFEScore suggests that Chest ImaGenome outperforms both the CheXpert labeler and CheXbert, which is reasonable given Chest ImaGenome’s broader range of labels, although the improvement is rather marginal. However, it is evident that all three baseline methods fall short of fully capturing the factual information within the reports, likely due to their rigid annotation rules. This shortcoming is highlighted by their results in CXRFEScore, RadGraph F1, BERTScore, and most of the other metrics, compared to the fact extraction methods.

## 5 Conclusions & Future Work

In this work, we present a novel two-stage framework for extracting and encoding factual information from radiology reports. The first stage, fact extraction, uses ChatGPT and T5 to extract factual statements. The second stage, fact encoding, introduces CXRFE, a specialized variant of CXRBERT, fine-tuned through multitask learning by incorporating tasks that support representation improvement. Our system’s effectiveness is validated through comprehensive evaluations. Additionally, we introduce CXRFEScore, a novel metric for evaluating radiology text generation, leveraging our two-stage system. We anticipate that our work will stimulate further research in enhanced fact extraction and representation, LLM utilization, advanced training methodologies, and improved evaluation metrics. For future work, we aim to expand our framework to integrate visual modality, focusing on image-based fact detection and visual grounding.

## Limitations

Our study acknowledges several limitations and highlights areas for improvement. First, more expert evaluations, particularly from radiologists, are needed to refine the use of large language models (LLMs) in radiology. Although we extensively utilized publicly available gold standards, such as those from the Chest ImaGenome dataset, RadNLI, and MS-CXR-T, there remains room for improvement. For instance, involving radiologists in the prompt engineering process and developing more rigorous evaluation protocols are two strategies we believe will enhance the evaluation and utilization of LLMs for radiological text.

We also see potential in designing better triplet sampling heuristics, especially with input from radiologists. Optimizing LLM prompts for triplet sampling and incorporating more advanced auxiliary embeddings could further enhance our approach.

Furthermore, while our study focuses on text-only analysis, we recognize the importance of integrating visual data, such as chest X-ray images, into a multimodal framework. Devising a training protocol that effectively combines supervision from both images and text is an area of potential improvement for future work.

In this work, we limited our experiments to the sections "findings," "impression," and similar headings providing factual statements about the imaging exam. However, other sections, such as "comparison," "indication," and "history," were left out of the analysis, yet they can provide deeper insights into patient information and context. Investigating how this broader information can be extracted and encoded to enhance downstream applications is another avenue for future exploration and potential improvement.

Lastly, we acknowledge that our fact extraction algorithm may be limited due to its reliance on the `sent_tokenize` function of the NLTK library, which we use to obtain a preliminary division of the report into coarse sentences (before fact extraction). This method could falter when a fact spans multiple sentences connected through co-reference. While such occurrences are relatively uncommon in our observations, a deeper exploration of this linguistic aspect could guide the development of a more refined fact extraction mechanism that overcomes this challenge.

## Acknowledgements

This work was funded by the Chilean National Agency for Research and Development (ANID), including Instituto Milenio en Ingeniería e Inteligencia Artificial para la Salud (iHEALTH) ICN2021\_004; Centro Nacional de Inteligencia Artificial (CENIA) FB210017; Fondecyt regular 1231724; Fondecyt 1221425; and the ANID Scholarship Program / Doctorado Becas Chile / 2019 - 21191569. Additionally, Pablo was supported by the National Institutes of Health (NIH) grant 1R01AG067396. We are grateful for the support from all funding sources mentioned above.

## References

- Lisa C Adams, Daniel Truhn, Felix Busch, Avan Kader, Stefan M Niehues, Marcus R Makowski, and Keno K Bressen. 2023. Leveraging gpt-4 for post hoc transformation of free-text radiology reports into structured reporting: a multilingual feasibility study. *Radiology*, 307(4):e230725.
- AI@Meta. 2024. [Llama 3 model card](#).
- Emily Alsentzer, John R Murphy, Willie Boag, Weihung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. 2019. Publicly available clinical bert embeddings. *arXiv preprint arXiv:1904.03323*.
- Vladimir Araujo, Marie-Francine Moens, and Alvaro Soto. 2023. [Learning sentence-level representations with predictive coding](#). *Machine Learning and Knowledge Extraction*, 5(1):59–77.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proc of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. ACL.
- Shruthi Bannur, Stephanie Hyland, Qianchu Liu, Fernando Perez-Garcia, Maximilian Ilse, Daniel C Castro, Benedikt Boecking, Harshita Sharma, Kenza Bouzid, Anja Thieme, et al. 2023. Learning to exploit temporal structure for biomedical vision-language processing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15016–15027.
- Benedikt Boecking, Naoto Usuyama, Shruthi Bannur, Daniel C Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, et al. 2022. Making the most of text semantics to improve biomedical vision-language processing. In *European conference on computer vision*, pages 1–21. Springer.

- Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria de la Iglesia-Vayá. 2019. Padchest: A large chest x-ray image dataset with multi-label annotated reports. *arXiv:1901.07441*.
- Zhihong Chen, Shizhe Diao, Benyou Wang, Guanbin Li, and Xiang Wan. 2023a. Towards unifying medical vision-and-language pre-training via soft prompts. *arXiv preprint arXiv:2302.08958*.
- Zhihong Chen, Maya Varma, Xiang Wan, Curtis Langlotz, and Jean-Benoit Delbrouck. 2023b. [Toward expanding the scope of radiology report summarization to multiple anatomies and modalities](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 469–484, Toronto, Canada. Association for Computational Linguistics.
- Jean-Benoit Delbrouck, Pierre Chambon, Christian Bluethgen, Emily Tsai, Omar Almusa, and Curtis Langlotz. 2022. [Improving the factual correctness of radiology report generation with semantic rewards](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4348–4360, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jean-Benoit Delbrouck, Maya Varma, Pierre Chambon, and Curtis Langlotz. 2023. Overview of the radsum23 shared task on multi-modal and multi-anatomical radiology report summarization. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 478–482.
- Dina Demner-Fushman, Marc D. Kohli, Marc B. Rosenman, Sonya E. Shooshan, Laritza Rodriguez, Sameer Antani, George R. Thoma, and Clement J. McDonald. 2015. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Markus Eberts and Adrian Ulges. 2020. [Span-based joint entity and relation extraction with transformer pre-training](#). In *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, pages 2006–2013. IOS Press.
- Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2020. [Domain-specific language model pretraining for biomedical natural language processing](#).
- Yu Gu, Sheng Zhang, Naoto Usuyama, Yonas Woldesenbet, Cliff Wong, Praneeth Sanapathi, Mu Wei, Naveen Valluri, Erika Strandberg, Tristan Naumann, et al. 2023. Distilling large language models for biomedical knowledge extraction: A case study on adverse drug events. *arXiv preprint arXiv:2307.06439*.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597.
- Saahil Jain, Ashwin Agrawal, Adriel Saporta, Steven Truong, Du Nguyen Duong Nguyen Duong, Tan Bui, Pierre Chambon, Yuhao Zhang, Matthew Lungren, Andrew Ng, Curtis Langlotz, Pranav Rajpurkar, and Pranav Rajpurkar. 2021a. [Radgraph: Extracting clinical entities and relations from radiology reports](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran.
- Saahil Jain, Akshay Smit, Steven QH Truong, Chanh DT Nguyen, Minh-Thanh Huynh, Mudit Jain, Victoria A. Young, Andrew Y. Ng, Matthew P. Lungren, and Pranav Rajpurkar. 2021b. [Visualchexpert: Addressing the discrepancy between radiology report labels and image labels](#). In *Proceedings of the Conference on Health, Inference, and Learning, CHIL '21*, page 105–115, New York, NY, USA. Association for Computing Machinery.
- Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. 2019a. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(1):317.
- Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. 2019b. [MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs](#). *arXiv:1901.07042*.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. 2016. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9.
- Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. 2023. GPT-4 passes the bar exam. *Available at SSRN 4389233*.



- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. ACL.
- Fangyu Liu, Qianchu Liu, Shruthi Bannur, Fernando Pérez-García, Naoto Usuyama, Sheng Zhang, Tristan Naumann, Aditya Nori, Hoifung Poon, Javier Alvarez-Valle, Ozan Oktay, and Stephanie L. Hyland. 2023a. [Compositional Zero-Shot Domain Transfer with Text-to-Text Models](#). *Transactions of the Association for Computational Linguistics*, 11:1097–1113.
- Jiaxiang Liu, Tianxiang Hu, Yan Zhang, Xiaotang Gai, YANG FENG, and Zuozhu Liu. 2023b. [A chatGPT aided explainable framework for zero-shot medical image diagnosis](#). In *ICML 3rd Workshop on Interpretable Machine Learning in Healthcare (IMLH)*.
- Qianchu Liu, Stephanie Hyland, Shruthi Bannur, Kenza Bouzid, Daniel Castro, Maria Wetscherek, Robert Tinn, Harshita Sharma, Fernando Pérez-García, Anton Schwaighofer, Pranav Rajpurkar, Sameer Khanna, Hoifung Poon, Naoto Usuyama, Anja Thieme, Aditya Nori, Matthew Lungren, Ozan Oktay, and Javier Alvarez-Valle. 2023c. [Exploring the boundaries of GPT-4 in radiology](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14414–14445, Singapore. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Chong Ma, Zihao Wu, Jiaqi Wang, Shaochen Xu, Yaonai Wei, Zhengliang Liu, Lei Guo, Xiaoyan Cai, Shu Zhang, Tuo Zhang, et al. 2023. [Impressiongpt: an iterative optimizing framework for radiology report summarization with chatgpt](#). *arXiv preprint arXiv:2304.08448*.
- Pablo Messina, Pablo Pino, Denis Parra, Alvaro Soto, Cecilia Besa, Sergio Uribe, Marcelo Andia, Cristian Tejos, Claudia Prieto, and Daniel Capurro. 2022. A survey on deep learning and explainability for automatic report generation from medical images. *ACM Computing Surveys (CSUR)*, 54(10s):1–40.
- Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. 2021. [Improving factual completeness and consistency of image-to-text radiology report generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5288–5304, Online. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proc of the 40th Annual Meeting of the ACL*, pages 311–318, Philadelphia, Pennsylvania, USA. ACL, ACL.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in pytorch.
- Pablo Pino, Denis Parra, Cecilia Besa, and Claudio Lagos. 2021. Clinically correct report generation from chest x-rays using templates. In *Machine Learning in Medical Imaging: 12th International Workshop, MLMI 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Proceedings 12*, pages 654–663. Springer.
- Pablo Pino, Denis Parra, Pablo Messina, Cecilia Besa, and Sergio Uribe. 2020. Inspecting state of the art performance and nlp metrics in image-based medical report generation. *arXiv preprint arXiv:2011.09257*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Alexey Romanov and Chaitanya Shivade. 2018. [Lessons from natural language inference in the clinical domain](#).
- Yucheng Shi, Hehuan Ma, Wenliang Zhong, Gengchen Mai, Xiang Li, Tianming Liu, and Junzhou Huang. 2023. [Chatgraph: Interpretable text classification by converting chatgpt knowledge to graphs](#). *arXiv preprint arXiv:2305.03513*.
- Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y. Ng, and Matthew P. Lungren. 2020. [Chexbert: Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT](#). *CoRR*, abs/2004.09167.
- Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel Rueckert. 2023. Interactive and explainable region-guided radiology report generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7433–7442.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proc of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575.
- David Wadden, Ulme Wennberg, Yi Luan, and Hananeh Hajishirzi. 2019. [Entity, relation, and event extraction with contextualized span representations](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5784–5789, Hong Kong, China. Association for Computational Linguistics.



Fuying Wang, Yuyin Zhou, Shujun Wang, Varut Vardhanabhuti, and Lequan Yu. 2022. Multi-granularity cross-modal alignment for generalized medical visual representation learning. *Advances in Neural Information Processing Systems*, 35:33536–33549.

Joy T Wu, Nkechinyere Agu, Ismini Lourentzou, Ismini Lourentzou, Arjun Sharma, Joseph Alexander Paguio, Jasper Seth Yao, Edward C Dee, William Mitchell, Satyananda Kashyap, Andrea Giovannini, Leo Anthony Celi, and Mehdi Moradi. 2021. [Chest imagenome dataset for clinical reasoning](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran.

Justin Xu, Zhihong Chen, Andrew Johnston, Louis Blankemeier, Maya Varma, Curtis Langlotz, and Jean-Benoit Delbrouck. 2024. Overview of the first shared task on clinical text generation: Rrg24 and “discharge me!”. In *The 23rd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, Bangkok, Thailand. Association for Computational Linguistics.

Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin, Chris Callison-Burch, and Mark Yatskar. 2023. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19187–19197.

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. Linkbert: Pretraining language models with document links. In *Association for Computational Linguistics (ACL)*.

Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser Ururahy Nunes Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein Abad, Andrew Y Ng, et al. 2022. Evaluating progress in automatic chest x-ray radiology report generation. *medRxiv*, pages 2022–08.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020a. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Yuhao Zhang, Derek Merck, Emily Tsai, Christopher D. Manning, and Curtis Langlotz. 2020b. [Optimizing the factual correctness of a summary: A study of summarizing radiology reports](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5108–5120, Online. Association for Computational Linguistics.

## A Appendix

### A.1 Additional Related Work

**Evaluation of Factual Correctness in Radiology Text Generation.** One important area of application motivating this work is the evaluation of

factual correctness in systems that generate radiological text, usually from input medical imaging. Recent research emphasizes enhancing the accuracy of generated facts in applications like report generation (Miura et al., 2021; Delbrouck et al., 2022; Pino et al., 2020, 2021) and summarization (Zhang et al., 2020b; Delbrouck et al., 2023). Pino et al. (2020) conducted an evaluation of several trivial report generation baselines using established metrics such as BLEU (2002), ROUGE-L (2004), and CIDEr-D (2015). They achieved results comparable to state-of-the-art papers at the time. However, when assessed using the CheXpert labeler (Irvin et al., 2019), a domain-specific NLP tool designed to detect 13 findings, the performance was notably poor, underscoring the urgent necessity for standardizing improved evaluation metrics among researchers. More recently, Delbrouck et al. (2022) repurposed RadGraph’s entity and relation extraction model (Jain et al., 2021a) to create a factual correctness reward. This reward measures the overlap of entities and relations between real and generated reports, serving as the guiding signal to optimize a report generation model through reinforcement learning. Interestingly, their proposed reward aligns functionally with the RadGraph F1 metric introduced by Yu et al. (2022). Yu et al. conducted a study on metrics for radiology report generation, determining that RadGraph F1 and BLEU show the highest correlation with radiologists’ judgement. Recently, the RadSum23 challenge (Delbrouck et al., 2023) evaluated multimodal radiology report summarization quality using BLEU-4 (2002), ROUGE-L (2004), BERTScore (2020a), CheXbert F1 (2020), and RadGraph F1 (2021a).

Our work is highly relevant in this domain because of our development of the CXRFEScore metric. This metric leverages the strengths of both stages within our framework: fact extraction and encoding. As a result, CXRFEScore is specifically designed to assess the factual accuracy of generated radiological text against a reference text. Section 4 offers a comprehensive evaluation, showcasing the effectiveness of our two-stage system. This includes a comparison of CXRFEScore with commonly used metrics in the literature, with very favorable results.

**LLMs in Medicine.** Our work falls under the category of applications of LLMs to the medical domain. Specifically, in this work we make use of ChatGPT versions GPT-3.5 and GPT-4 through

OpenAI's API<sup>3</sup>. Recent works have shown the effectiveness of ChatGPT applied to medical tasks. Most notably, recent work by Liu et al. (2023c) explored the boundaries of GPT-4 in radiology, evaluated in tasks such as classification, summarization, and natural language inference, with remarkable performance. Liu et al. (2023b) employed ChatGPT to generate short sentences with plausible symptoms of medical conditions for interpretable zero-shot medical image diagnosis. Adams et al. (2023) used GPT-4 to transform free-text radiology reports into structured templates, with remarkable results. GPT-4 is also known for having passed the bar exam (Katz et al., 2023). Inspired by these results, we make extensive use of ChatGPT to produce abundant annotations through diverse prompts.

**Knowledge Distillation from LLMs.** Our approach can also be viewed as a form of LLM knowledge distillation, where a large language model (LLM) "teacher" generates annotations for training a more compact "student" model. Shi et al. (2023) illustrated this by using ChatGPT to extract knowledge graphs from text to train a smaller model for text classification. Similarly, Gu et al. (2023) applied this concept in the biomedical field, distilling knowledge from GPT-3.5 for adverse drug event extraction with student models like PubMedBERT and BioGPT. In line with these works, our research presents a form of knowledge distillation as we fine-tune T5 (student) using annotations generated by ChatGPT (teacher). This process trains an imitator that is more cost-effective and easier to run on our machines. We apply this method to the tasks of fact extraction from reports and the generation of metadata and labels for each fact. Additionally, we extensively use GPT-4 to produce high-quality NLI labels. Although we don't use T5 in this case, these labels serve as part of the training supervision for CXRFE, which can also be considered a form of knowledge distillation.

## A.2 Fact Extraction Implementation Details

In our experiments, we processed the 227,827 radiology reports provided by the MIMIC-CXR dataset (Johnson et al., 2019b). To pinpoint relevant radiological sections in the MIMIC-CXR reports, such as *Findings*, *Impression*, and various other headings, we employed a combination of regular expressions and simple rules. These sections

were then segmented into sentences using NLTK's `sent_tokenize` function, resulting in 677,694 unique sentences after processing the entire dataset. Subsequently, we extracted facts from each sentence.

Extracting factual statements from a sentence of a free-text radiology report using traditional approaches, such as regular expressions, hand-designed rules, and similar heuristics, poses significant challenges due to the complexity and diversity of vocabulary used by radiologists. A more promising alternative is to leverage the capabilities of powerful LLMs like ChatGPT to tackle this task. Table 6 presents several examples of facts extracted by GPT-4 from challenging sentences. This is achieved by providing the model with a specific set of instructions: (Refer to Figure 14 for a screenshot of OpenAI's web interface displaying the same prompt.)

Relevant facts:

1. observations of abnormalities 2. observations of diseases 3. observations of strange visual patterns 4. observations of devices 5. observations of foreign bodies 6. observations of specific anatomical regions that look normal or healthy 7. absences of abnormalities (usually expressed with a negation) 8. comparisons with respect to a previous study (something changed or remained the same)

Task:

Given a sentence taken from a chest x-ray report, generate a JSON list of relevant facts. Each fact should be about one observation. If a sentence mentions multiple observations, each observation should be extracted as a separate fact. Each fact should include the anatomical location where it was observed. If multiple facts occur in the same location, repeat the location in each fact.

If no relevant facts are mentioned, return [] (an empty array).

Examples:

Opacity and density in the right lobe

[ "opacity in the right lobe", "density in the right lobe" ]

Lungs are well inflated without evidence of focal airspace consolidation to suggest pneumonia.

[ "well inflated lungs", "lungs without evidence of focal airspace consolidation", "lungs without evidence of pneumonia" ]

Taken together, compared with less than 1 hr earlier, the findings are suggestive of worsening of CHF, with new or significantly increased left greater right pleural effusions and underlying bibasilar collapse and/or consolidation, particularly on the left.

[ "worsening of CHF", "new or significantly increased left pleural effusions", "new or significantly increased right pleural effusions", "underlying bibasilar collapse on the left", "underlying consolidation on the left", ]

<sup>3</sup><https://platform.openai.com/>

No acute cardiopulmonary abnormality  
 [ "no acute cardiopulmonary abnormality" ]

Given the relatively high cost of both GPT-4 and GPT-3.5 models, we opted to annotate a subset of the MIMIC-CXR sentences with these LLMs. To select this subset, we employed a two-pronged approach. Firstly, we ranked sentences based on the sum of the inverse frequency of their tokens, thus prioritizing longer and more complex sentences, which often include infrequent abnormalities, typos, and symbols. Secondly, to ensure diversity of topics, we clustered the sentences into 200 groups using K-Means, utilizing embedding representations obtained with BioViL-T, and sampled equally from each cluster in order of difficulty. This combined strategy allowed us to curate a subset that is both diverse and challenging. GPT-4-0613 extracted facts from 24,998 sentences, while GPT-3.5-turbo-0613 processed 69,936. Subsequently, we trained T5-small for fact extraction using a total of 94,934 training examples, reserving 200 examples for validation. Once trained, T5-small annotated the remaining 582,760 sentences, yielding a total of 1,341,830 facts, including duplicates, of which 591,920 were unique after duplicate removal.

### A.3 CXRFE’s Tasks Details

CXRFE is a fine-tuned version of CXR-BERT-specialized, accessible for download from <https://huggingface.co/microsoft/BiomedVLP-CXR-BERT-specialized>. This fine-tuning process entails multi-task learning, incorporating six distinct tasks: triplet loss for sentence ranking (T), sentence classification (C), sentence decoding (SD), natural language inference (NLI), quadruplet loss for enforcing separation between entailment and contradiction pairs (EC), and entity and relation extraction (ER). While there exist a total of 64 possible task combinations, we focus on 11 combinations in our experiments, namely: T, T+C, T+ER, T+SD, T+EC, T+NLI, T+EC+NLI, T+C+EC+NLI, T+C+EC+NLI+ER, T+C+EC+NLI+SD, T+C+EC+NLI+ER+SD. This selection is guided by the following heuristics:

- Triplet loss (T) is always included, given our belief, as discussed in Appendix A.4, that the developed triplets dataset captures many desirable properties for learned sentence embeddings.

- We explore each combination of T with the other tasks individually: T+C, T+ER, T+SD, T+EC, T+NLI.
- Recognizing the complementary nature of EC and NLI in exploiting natural language inference data, we explore T+EC+NLI.
- According to our experimental results, both C and EC+NLI serve as effective auxiliary tasks; thus, we keep them fixed while varying other combinations of SD and ER: T+C+EC+NLI, T+C+EC+NLI+ER, T+C+EC+NLI+SD, T+C+EC+NLI+ER+SD.

Next, we delve into the implementation details of each task in our experiments.

**Triplet loss for sentence ranking (T).** One of the tasks we explore for model fine-tuning is sentence ranking via triplet loss. Figure 4 illustrates this task. Concretely, we forward 3 sentences (anchor, positive, negative) through CXR-BERT-specialized with weight sharing, obtaining three vectors  $a$ ,  $b$ , and  $c$  each of dimension 128, and compute  $\Delta\text{sim}(a, p, n) = a \cdot p - a \cdot n$ . This serves as the input logit for a binary cross-entropy loss.

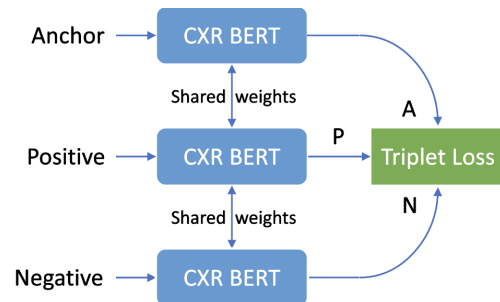


Figure 4: Triplet loss (T)

We provide details on the construction of the triplets dataset used in this task in Section A.4.

**Sentence classification tasks (C).** A second task group is classification tasks (Figure 5). These include category (5 classes: *anatomical finding, disease, technical assessment, tubes and lines and device*), health status (4 classes: *normal, abnormal, ambiguous, unknown*), comparison status (15 classes, see Figure 18), Chest ImaGenome observations (74 classes, see Figure 19) and anatomical locations (38 classes, see Figure 20). Category, Health Status and Comparison Status are single-label multi-class classification tasks, whereas Chest ImaGenome observations and anatomical locations are multi-label binary classification tasks. Implementing these tasks require attaching fully con-

Table 6: Examples of facts extracted by GPT-4, with the prompt shown in Figure 14

Original sentence	Extracted facts	Comment
2.Atelectasis of the left lower lobe 3.Stable left lower lobe laceration.	<ul style="list-style-type: none"> <li>Atelectasis of the left lower lobe</li> <li>Stable left lower lobe laceration</li> </ul>	In this example, NLTK's sent_tokenize failed to split the text into two sentences, but GPT-4 has no problems identifying the facts.
A CT of the abdomen showing the lower chest performed on adds the detection of acinar-sized nodules, some of which are confluent at the lung periphery and confirms the presence of small bilateral pleural effusions.	<ul style="list-style-type: none"> <li>detection of acinar-sized nodules at the lung periphery</li> <li>some nodules are confluent at the lung periphery</li> <li>presence of small bilateral pleural effusions</li> </ul>	A long complicated sentence broken down into 3 shorter factual statements.
NO chf, focal infiltrate, or gross effusionis identified	<ul style="list-style-type: none"> <li>no CHF identified</li> <li>no focal infiltrate identified</li> <li>no gross effusion identified</li> </ul>	A negated sentence becomes 3 short facts, each one preserving the negation. GPT-4 also correctly fixes a typo.
Severe cardiomegaly is again noted as well as aortic valve replacement including stent material at the aortic root compatible with appearance of Carevalve Aortic bioprosthesis.	<ul style="list-style-type: none"> <li>Severe cardiomegaly</li> <li>aortic valve replacement</li> <li>stent material at the aortic root</li> <li>appearance of Carevalve Aortic bioprosthesis at the aortic root</li> </ul>	Complex sentence becomes 4 shorter facts.
A right-sided internal jugular central venous catheter is seen, which appears slightly high in position, terminating in the region of the right innominate vein/right internal jugular/innominate vein junction.	<ul style="list-style-type: none"> <li>right-sided internal jugular central venous catheter is seen</li> <li>catheter appears slightly high in position'</li> <li>catheter terminating in the region of the right innominate vein</li> <li>catheter terminating in the region of the right internal jugular</li> <li>catheter terminating in the region of the innominate vein junction</li> </ul>	Complex sentence becomes 5 shorter facts. Note GPT-4's ability to handle forward slashes (/).
THEY REPORT TEXT FOLLOWS: The heart and great vessels are normal.	<ul style="list-style-type: none"> <li>normal heart</li> <li>normal great vessels</li> </ul>	GPT-4 successfully ignores a distracting header (in uppercase) and correctly distributes the word "normal" to each extracted fact.

nected heads on top of CXR-BERT-specialized's built-in projection layer in order to perform the classification.

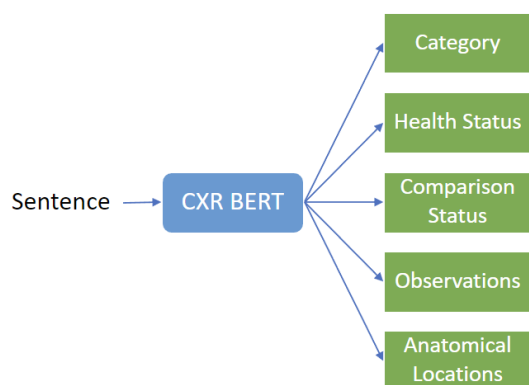


Figure 5: Sentence classification (C)

The data for the classification tasks comes from two primary sources. The first source is the Chest ImaGenome dataset, which provides observation and anatomical location labels in scene graphs, conveniently linked to specific sentences from MIMIC-CXR reports. The second source is ChatGPT,

which we leverage to annotate the extracted facts with details such as category, health status, comparison status, observations, and anatomical locations. This process involves a teacher-student approach similar to our method for fact extraction: we use ChatGPT to annotate a subset of facts, incurring some monetary cost, with the prompts shown in Figures 17, 18, 19, and 20, and then fine-tune T5 to annotate the remaining facts at no additional cost.

Concretely, from the Chest ImaGenome scene graphs, we retrieved 556,111 sentences, each annotated with observations (74 classes) and anatomical locations (38 classes). In addition, facts were annotated as follows:

- 5,000 facts were annotated with observations by GPT-4-0613.
- 84,708 facts were annotated with observations by GPT-3.5-turbo-0613.
- 2,816,982 facts were annotated with observations by T5-small fine-tuned with the teacher-student approach. These include facts extracted from reports plus additional facts obtained via paraphrases. (For more details on



paraphrases, see Section A.4 on the construction of the triplets dataset.)

Similarly, anatomical locations were annotated as follows:

- 72,400 facts were annotated with anatomical locations by GPT-3.5-turbo-0613.
- 2,598,778 facts were annotated with anatomical locations by T5-small fine-tuned.

For other classification tasks (Category, Health Status, Comparison Status), we generated a JSON object with metadata from each fact using the prompt shown in Figure 17:

- 59,921 facts were annotated by GPT-3.5-turbo-0613.
- 535,959 facts were annotated by T5-small fine-tuned with the teacher-student approach.

**Sentence decoding (SD).** Another task is sentence decoding (Figure 6). We attach a lightweight, shallow Transformer Decoder to CXR-BERT-specialized’s projection layer in order to generate back the original sentence. This can be viewed a sort of text autoencoder, forcing the projection layer to capture as much information as possible of the input sentence to facilitate the reconstruction of the sentence by the Transformer Decoder. We use a Transformer Decoder with embedding, hidden and feedforward dimension 256, only one self-attention head and only one layer.



Figure 6: Sentence decoding (SD)

Because this task relies on self-supervision and does not require specialized annotations, any sentence or fact can serve as a training instance.

**Natural language inference (NLI).** For NLI, we adopt an approach similar to that of SBERT (Reimers and Gurevych, 2019), by concatenating the embeddings of the premise, hypothesis and their element-wise multiplication, followed by a fully connected layer and a softmax layer for NLI classification (see Figure 7).

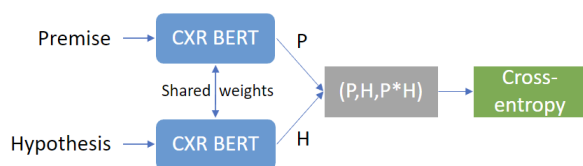


Figure 7: Natural language inference (NLI)

As discussed in Section 3.4.3, our NLI dataset was compiled by merging data from MedNLI (14,049 pairs), the RadNLI development set (480 pairs), and examples generated by GPT-4 (154,498 pairs), resulting in a comprehensive collection of 169,025 pairs. These consist of 25,175 entailment pairs, 44,729 neutral pairs, and 99,121 contradiction pairs. Additionally, the RadNLI test set (480 pairs) and MS-CXR-T (361 pairs) were reserved for evaluation purposes. To generate NLI examples using GPT-4, we employed four distinct prompts:

- A prompt that generates entailment, neutral, and contradiction sentences from a reference sentence (Figure 22).
- A prompt that aims to produce examples analogous to a given NLI reference example (Figure 23).
- A prompt that predicts the correct label for a given premise and hypothesis, incorporating Chain-of-Thought (CoT) reasoning and examples (Figure 25). A simpler version of this prompt, which only requests the correct label, was also considered (Figure 24). However, as indicated in Table 3, CoT is crucial for achieving significantly more accurate predictions.
- A prompt designed to generate contradictory sentences relative to a reference sentence (Figure 26).

**Quadruplet loss: enforcing separation between entailment and contradiction pairs (EC).**

The next task is what we refer to as entailment-/contradiction quadruplet loss (Figure 8). The goal of this task is to promote a generalized separation of entailment and contradiction sentence pairs in the latent space, by randomly sampling entailment and contradiction pairs and requiring that the entailment pair have greater similarity than the contradiction pair. This loss was a key contributor to achieving the state-of-the-art results presented in Table 2.

Since this task complements the standard NLI task, it utilizes the same dataset but excludes neutral pairs, focusing solely on entailment and contradiction pairs.

**Entity and relation extraction (ER).** Lastly, for entity and relation extraction we augment CXR-BERT-specialized with the layers proposed by SpERT (Eberts and Ulges, 2020), as illustrated in Figure 9. This adaptation was relatively straightforward, since the authors of SpERT released an implementation (<https://github.com/lavis-nlp/spert/>) that is compatible with Hug-

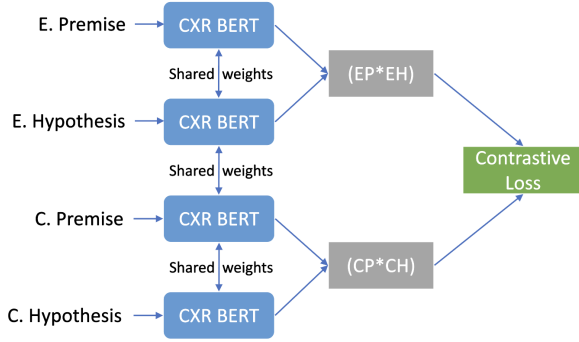


Figure 8: Entailment/contradiction quadruplet loss (EC)

giface models like CXRT-BERT-specialized.

To implement this task, we utilize the gold data provided by the RadGraph dataset (Jain et al., 2021a). This dataset comprises 500 MIMIC-CXR radiology reports, annotated with an entity-and-relation schema by board-certified radiologists. Additionally, a test set containing 50 MIMIC-CXR and 50 CheXpert reports, annotated in the same manner, is included. We use all of this data for training.

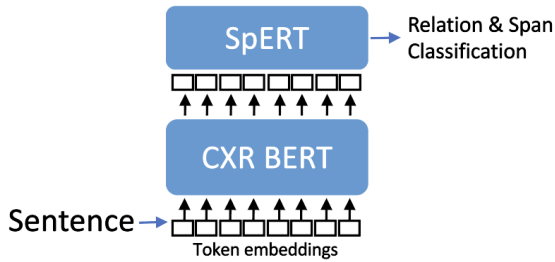


Figure 9: Entity and relation extraction (ER) with SpERT

#### A.4 Triplets Dataset Construction Details

CXRFE is trained to generate sentence embeddings that organize semantically similar sentences into clusters within the embedding space via a triplet ranking task. This task involves a dataset comprising triplets  $(a, p, n)$ , where  $a$  represents an anchor sentence,  $p$  is a positive sentence, and  $n$  denotes a negative sentence. The objective is to train a text encoder in a manner that ensures  $\text{sim}(a, p) > \text{sim}(a, n)$  holds true for each triplet, with  $\text{sim}(x, y)$  computed as the cosine similarity between sentence embeddings (or dot product if they are already normalized).

The rationale behind selecting triplet loss lies in its versatility, offering flexibility in terms of desired properties for the learned embedding function.

For instance, if  $a$  and  $p$  denote identical medical observations using different vocabulary, while  $n$  presents an unrelated observation, it’s logical to enforce  $\text{sim}(a, p) > \text{sim}(a, n)$ . Similarly, if there’s substantial topic overlap between  $a$  and  $p$ , whereas  $n$  is largely unrelated or contradictory to  $a$ , it’s reasonable to aim for  $\text{sim}(a, p) > \text{sim}(a, n)$  to hold true.

To guide a text encoder towards learning sentence embeddings consistent with triplets, we can define the difference in similarities as  $\Delta\text{sim}(a, p, n) = \text{sim}(a, p) - \text{sim}(a, n)$ . Then, a straightforward approach is to utilize  $\text{sim}(a, p, n)$  as the input for a binary cross-entropy loss, where the ground-truth label is consistently set to 1.

**Notation.** To define our triplet sampling heuristics, we use the notation  $E(x)$  for the embedding of sentence  $x$ ,  $\cos(E(x), E(y))$  for the cosine similarity between embeddings of  $x$  and  $y$ ,  $\text{lev}(x, y)$  for the Levenshtein string distance between them, and  $\text{levsim}(x, y) = 1 - \text{lev}(x, y) / \max(\text{len}(x), \text{len}(y))$ .  $c(x)$  indicates the cluster id for sentence  $x$  after running a clustering algorithm like K-Means on the sentence embeddings. This requires having an auxiliary text encoder capable of producing these auxiliary embeddings and clusters. Specifically, we use BioViL-T (Bannur et al., 2023), a state-of-the-art BERT-based model for radiological text, available on Huggingface<sup>4</sup>.

With this, we sample triplets based on the following heuristics:

**Rule 1: Rank paraphrases highly.**  $\Delta\text{sim}(a, p, n) > 0$  if  $p$  is a paraphrase of  $a$  generated by ChatGPT and  $n$  is any other sentence (unless  $\cos(E(a), E(p)) < \cos(E(a), E(n))$  and  $\text{lev}(a, p) > \text{lev}(a, n)$ ). To generate paraphrases, we employ the prompts shown in Figures 15 and 16, one for paraphrasing facts and another for paraphrasing anatomical locations. We decided to paraphrase anatomical locations too in order to strengthen the model’s understanding of their vocabulary. As a reminder, the anatomical locations are obtained from facts as part of the metadata generated with the prompt shown in Figure 17.

**Rule 2: Sample triplets according to the consensus of E and lev, while anchor and positive share the same health status.**  $\Delta\text{sim}(a, p, n) > 0$  if  $\text{HS}(a) = \text{HS}(p)$ ,  $c(p) = c(a)$ ,  $c(p) \neq c(n)$ ,

<sup>4</sup><https://huggingface.co/microsoft/BiomedVLP-BioViL-T>

Table 7: Examples of template-based generated reports for different label extraction algorithms.

Ground-truth Report	FE (T5-small)	FE (GPT-4)	FE (Meta-Llama-3-8B)
New PICC line on the right is projecting with its tip somewhere in the mediastinum. Appears to cross the midline, there is concern for potential arterial location. The initial line concerns were communicated over the telephone at the time of the wet read. Repeat PA and lateral radiograph, taken approximately an hour after the radiograph demonstrated the PICC line in the mid SVC. Potential small right pleural effusion. Stable moderate cardiomegaly.	new PICC line on the right. tip of PICC line in the mediastinum. potential arterial location crossing the midline. PICC line in the mid SVC. potential small right pleural effusion. stable moderate cardiomegaly  <b>CXRFEScore:</b> 1.000 <b>RadGraph F1 Full:</b> 0.750 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0	new PICC line on the right. PICC line tip projecting in the mediastinum. potential arterial location crossing the midline. repeat PA and lateral radiograph taken approximately an hour after the previous radiograph. PICC line observed in the mid SVC. potential small right pleural effusion. stable moderate cardiomegaly  <b>CXRFEScore:</b> 0.891 <b>RadGraph F1 Full:</b> 0.899 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0	new PICC line on the right. PICC line tip in the mediastinum. appears to cross the midline. concern for potential arterial location. PA radiograph. lateral radiograph. PICC line in the mid SVC. potential small right pleural effusion. stable cardiomegaly  <b>CXRFEScore:</b> 0.966 <b>RadGraph F1 Full:</b> 0.813 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0
	<b>CheXbert</b>	<b>CheXpert labeler</b>	<b>Chest ImaGenome</b>
	the heart is enlarged. the cardiomeastinal silhouette is enlarged. no focal consolidation. the lungs are free of focal airspace disease. no atelectasis. a device is seen. pleural effusion is seen. no fibrosis. no pneumonia. no pneumothorax is seen. no pulmonary edema. no pulmonary nodules or mass lesions identified. no fracture is seen  <b>CXRFEScore:</b> 0.481 <b>RadGraph F1 Full:</b> 0.017 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0	the heart is enlarged. the cardiomeastinal silhouette is enlarged. no focal consolidation. the lungs are free of focal airspace disease. no atelectasis. a device is seen. pleural effusion is seen. no fibrosis. no pneumonia. no pneumothorax is seen. no pulmonary edema. no pulmonary nodules or mass lesions identified. no fracture is seen  <b>CXRFEScore:</b> 0.481 <b>RadGraph F1 Full:</b> 0.017 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0	enlarged cardiac silhouette in cardiac silhouette. abnormal cardiac silhouette. picc in left shoulder. picc in mediastinum. lung opacity in right costophrenic angle. pleural effusion in right costophrenic angle. abnormal right costophrenic angle. lung opacity in right lung. pleural effusion in right lung. abnormal right lung. picc in right shoulder. picc in svc. enlarged cardiac silhouette. lung opacity. pleural effusion. picc  <b>CXRFEScore:</b> 0.660 <b>RadGraph F1 Full:</b> 0.121 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 1.0
Ground-truth Report	FE (T5-small)	FE (GPT-4)	FE (Meta-Llama-3-8B)
In comparison with the earlier study of this date, the Dobhoff tube can be seen only to the level of the distal esophagus, were crosses the lower margin of the image, to evaluate the tip, a view of the upper abdomen must be obtained using abdominal technique. The area of increased opacification at the right base seen previously is not appreciated at this time.	Dobhoff tube seen only to the level of the distal esophagus. Dobhoff tube crosses the lower margin of the image. evaluation of the tip of Dobhoff tube. suggestion for a view of the upper abdomen using abdominal technique. area of increased opacification at the right base not appreciated  <b>CXRFEScore:</b> 0.898 <b>RadGraph F1 Full:</b> 0.633 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 0.929	Dobhoff tube seen only to the level of the distal esophagus. distal esophagus crosses the lower margin of the image. need to obtain a view of the upper abdomen to evaluate the tip of the Dobhoff tube. previous increased opacification at the right base. current absence of increased opacification at the right base  <b>CXRFEScore:</b> 0.674 <b>RadGraph F1 Full:</b> 0.575 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 0.929	observations of devices in the esophagus. device crosses the lower margin of the image in the esophagus. device tip cannot be evaluated in the esophagus. device tip must be evaluated in the upper abdomen. absence of increased opacification at the right base. absence of increased opacification at the right base (compared to previous study)  <b>CXRFEScore:</b> 0.489 <b>RadGraph F1 Full:</b> 0.193 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 0.929
	<b>CheXbert</b>	<b>CheXpert labeler</b>	<b>Chest ImaGenome</b>
	heart size is normal. the mediastinal contour is normal. no focal consolidation. one or more airspace opacities are seen. no atelectasis. no pleural effusion. no fibrosis. no pneumonia. no pneumothorax is seen. no pulmonary edema. no pulmonary nodules or mass lesions identified. no fracture is seen  <b>CXRFEScore:</b> 0.344 <b>RadGraph F1 Full:</b> 0.0 <b>CheXpert Acc:</b> 0.929 <b>CheXbert Acc:</b> 1.0	heart size is normal. the mediastinal contour is normal. no focal consolidation. one or more airspace opacities are seen. no atelectasis. a device is seen. no pleural effusion. no fibrosis. no pneumonia. no pneumothorax is seen. no pulmonary edema. no pulmonary nodules or mass lesions identified. no fracture is seen  <b>CXRFEScore:</b> 0.394 <b>RadGraph F1 full:</b> 0.0 <b>CheXpert Acc:</b> 1.0 <b>CheXbert Acc:</b> 0.929	enteric tube in abdomen. enteric tube in mediastinum. enteric tube in neck. no lung opacity in right lower lung zone. no abnormal right lower lung zone. no lung opacity in right lung. no abnormal right lung. no lung opacity. enteric tube  <b>CXRFEScore:</b> 0.561 <b>RadGraph F1 Full:</b> 0.066 <b>CheXpert Acc:</b> 0.857 <b>CheXbert Acc:</b> 0.786

Table 8: More examples of template-based generated reports for different label extraction algorithms.

Ground-truth Report	FE (T5-small)	FE (GPT-4)	FE (Meta-Llama-3-8B)
<p>Frontal and lateral chest radiographs were obtained. There are persistent, stable bilateral upper lung reticular nodular opacities consistent with history of sarcoidosis. No focal consolidation, pleural effusion, pneumothorax, or pulmonary edema is seen. The heart size is normal. Mediastinal and hilar contours are stable. 1. No focal consolidation to suggest pneumonia. 2. Stable bilateral upper lungs zone fibrosis consistent with history of sarcoidosis.</p>	<p>frontal chest radiograph obtained. lateral chest radiograph obtained. persistent bilateral upper lung reticular nodular opacities. stable bilateral upper lung reticular nodular opacities. opacities consistent with history of sarcoidosis. no focal consolidation. no pleural effusion. no pneumothorax. no pulmonary edema. normal heart size. stable mediastinal contours. stable hilar contours. no focal consolidation to suggest pneumonia. stable bilateral upper lungs zone fibrosis. fibrosis consistent with history of sarcoidosis</p> <p><b>CXRFEScore:</b> 0.982  <b>RadGraph F1 Full:</b> 0.862  <b>CheXpert Acc:</b> 1.0  <b>CheXbert Acc:</b> 0.857</p>	<p>persistent bilateral upper lung reticular nodular opacities. stable bilateral upper lung reticular nodular opacities. bilateral upper lung reticular nodular opacities consistent with history of sarcoidosis. no focal consolidation. no pleural effusion. no pneumothorax. no pulmonary edema. normal heart size. stable mediastinal contours. stable hilar contours. normal size heart. unremarkable mediastinum. clear lungs. no focal consolidation to suggest pneumonia. observation of diseases. stable bilateral upper lungs zone fibrosis. history of sarcoidosis</p> <p><b>CXRFEScore:</b> 0.911  <b>RadGraph F1 Full:</b> 0.723  <b>CheXpert Acc:</b> 1.0  <b>CheXbert Acc:</b> 1.0</p>	<p>persistent bilateral upper lung reticular nodular opacities. consistent with history of sarcoidosis. no focal consolidation. no pleural effusion. no pneumothorax. no pulmonary edema. normal heart size. mediastinal contours are stable. hilar contours are stable. no focal consolidation to suggest pneumonia. diseases present. diseases such as pneumonia present. diseases such as pulmonary embolism present. diseases such as CHF present. stable bilateral upper lung zone fibrosis. fibrosis in the upper lung zone. consistent with history of sarcoidosis</p> <p><b>CXRFEScore:</b> 0.858  <b>RadGraph F1 Full:</b> 0.743  <b>CheXpert Acc:</b> 0.857  <b>CheXbert Acc:</b> 0.714</p>
	<b>CheXbert</b>	<b>CheXpert labeler</b>	<b>Chest ImaGenome</b>
	<p>heart size is normal. the cardiomeastinal silhouette is enlarged. no focal consolidation. the lungs are free of focal airspace disease. no atelectasis. no pleural effusion. pleural thickening is present. no pneumonia. no pneumothorax is seen. no pulmonary edema. there are pulmonary nodules or mass identified. no fracture is seen</p> <p><b>CXRFEScore:</b> 0.691  <b>RadGraph F1 Full:</b> 0.363  <b>CheXpert Acc:</b> 0.857  <b>CheXbert Acc:</b> 0.929</p>	<p>heart size is normal. the cardiomeastinal silhouette is enlarged. no focal consolidation. one or more airspace opacities are seen. no atelectasis. no pleural effusion. pleural thickening is present. there is evidence of pneumonia. no pneumothorax is seen. no pulmonary edema. there are pulmonary nodules or mass identified. no fracture is seen</p> <p><b>CXRFEScore:</b> 0.641  <b>RadGraph F1 Full:</b> 0.359  <b>CheXpert Acc:</b> 1.0  <b>CheXbert Acc:</b> 0.786</p>	<p>no enlarged cardiac silhouette in cardiac silhouette. no abnormal cardiac silhouette. no consolidation in left costophrenic angle. no pleural effusion in left costophrenic angle. no pneumothorax in left costophrenic angle. no pulmonary edema or hazy opacity in left hilar structures. no consolidation in left lung. lung lesion in left lung. lung opacity in left lung. multiple masses or nodules in left lung. no pleural effusion in left lung. pleural or parenchymal scarring in left lung. no pneumothorax in left lung. no pulmonary edema or hazy opacity in left lung. interstitial lung disease in left lung. no pneumonia in left lung. abnormal left lung. lung lesion in left upper lung zone. lung opacity in left upper lung zone. multiple masses or nodules in left upper lung zone. pleural or parenchymal scarring in left upper lung zone. interstitial lung disease in left upper lung zone. abnormal left upper lung zone. no consolidation in right costophrenic angle. no pleural effusion in right costophrenic angle. no pneumothorax in right costophrenic angle. no pulmonary edema or hazy opacity in right hilar structures. no consolidation in right lung. lung lesion in right lung. lung opacity in right lung. multiple masses or nodules in right lung. no pleural effusion in right lung. pleural or parenchymal scarring in right lung. no pneumothorax in right lung. no pulmonary edema or hazy opacity in right lung. interstitial lung disease in right lung. no pneumonia in right lung. abnormal right lung. lung lesion in right upper lung zone. lung opacity in right upper lung zone. multiple masses or nodules in right upper lung zone. pleural or parenchymal scarring in right upper lung zone. interstitial lung disease in right upper lung zone. abnormal right upper lung zone. no consolidation. no enlarged cardiac silhouette. lung lesion. lung opacity. multiple masses or nodules. no pleural effusion. pleural or parenchymal scarring. no pneumothorax. no pulmonary edema or hazy opacity. interstitial lung disease. no pneumonia</p> <p><b>CXRFEScore:</b> 0.691  <b>RadGraph F1 Full:</b> 0.057  <b>CheXpert Acc:</b> 0.714  <b>CheXbert Acc:</b> 0.857</p>



$\cos(E(a), E(p)) > \cos(E(a), E(n)) + \text{margin}_{\cos}$  and  $\text{levsim}(a, p) > \text{levsim}(a, n) + \text{margin}_{\text{lev}}$ . This rule is very heuristic, and can be interpreted as a sort of mixture of experts. Essentially, we are looking for multiple signals to agree, increasing the likelihood that the resulting triplet passing those tests is in fact a high quality triplet. Levenshtein distance is useful to capture the notion of string similarity, while BioViL-T captures more the semantics (to some extent). The health status condition is adding an additional constraint. As a reminder, the health status is also obtained with the prompt of Figure 17.

**Rule 3: Short observation, detailed observation and the original fact (and their paraphrases) should be close to each other.** Given a fact  $f$ ,  $\Delta\text{sim}(a, p, n) > 0$  if  $a$  and  $p \in S(f)$ ,  $n \notin S(f)$  and  $c(a) \neq c(n)$  (unless  $\cos(E(a), E(p)) < \cos(E(a), E(n))$  and  $\text{lev}(a, p) > \text{lev}(a, n)$ ). Here,  $S(f)$  stands for the union of  $f$ , its detailed observation, its short observation and all the paraphrases (if any) generated for all of them with ChatGPT. The intuition here is that they are all closely related semantically, as they are all derived from the same fact (see Figure 17).

**Rule 4: Sample triplets according to Chest ImaGenome labels.**  $\Delta\text{sim}(a, p, n) > 0$  if  $\text{CIGL}(a) \cap \text{CIGL}(p) \neq \emptyset$ ,  $\text{CIGL}(a) \cap \text{CIGL}(n) = \emptyset$ ,  $\text{CIGL}(p) \cap \text{CIGL}(n) = \emptyset$ , and if  $(\cos(E(a), E(p)) > \cos(E(a), E(n)) \text{ AND } \text{levsim}(a, p) > \text{levsim}(a, n))$ . Here,  $\text{CIGL}(x)$  stands for the set of Chest ImaGenome labels of the sentence  $x$ . With this rule, our aim is to heuristically utilize the Chest ImaGenome labels to identify semantically similar sentences that warrant clustering together. However, we enhance this approach by incorporating signals from Levenshtein distance and BioViL-T to bolster our confidence in the triplet quality.

**Rule 5: Rank triplets according to the overlap of entities and relations from RadGraph.**  $\Delta\text{sim}(a, p, n) > 0$  if  $c(a) = c(p)$ ,  $c(a) \neq c(n)$ , and  $J(\text{RG}(a), \text{RG}(p)) > J(\text{RG}(a), \text{RG}(n)) + \text{margin}_{\text{RG}}$ . Here,  $\text{RG}(x)$  stands for the set of RadGraph entities and relations for the sentence  $x$ , and  $J$  for Jaccard similarity. Following the same idea of the previous rule, we seek to utilize the entities and relations provided by the RadGraph dataset as valuable cues for identifying semantically similar sentences that ought to be clustered together in the embedding space.

**Rule 6: Hard triplets generated by ChatGPT.**  $\Delta\text{sim}(a, p, n) > 0$  if  $(a, p, n)$  is a hard triplet gen-

erated by ChatGPT. The intuition behind this rule is very simple: we want to leverage ChatGPT’s remarkable skills to produce challenging triplets, requiring a good understanding of the text to be ranked correctly. Figure 21 shows the prompt used to generate these triplets along with an example.

For each rule, we sample approximately 3 to 4 million training triplets, along with 1,000 triplets for validation and 1,000 triplets for testing.

## A.5 Additional Metric Evaluation Details and Results

We conduct a thorough comparison involving CXRFEScore against several established metrics, including BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), METEOR (Banerjee and Lavie, 2005), CIDEr-D (Vedantam et al., 2015), BERTScore (Zhang et al., 2020a), CheXpert labeler (Irvin et al., 2019), CheXbert (Smit et al., 2020), and RadGraph F1 (Delbrouck et al., 2022). The last three are considered domain-specific metrics, tailored to the radiology domain, while the others serve as general-purpose evaluation metrics.

Both the CheXpert labeler and CheXbert output a 14-dimensional discrete vector representing 13 observations along with a label indicating "no findings". The values in this vector denote presence (1), absence (0), uncertainty (-1), and unknown (-2). We binarize this vector, treating both presence and uncertainty as 1, and absence and unknown as 0. This enables the computation of CheXpert F1, CheXbert F1, CheXpert Accuracy, and CheXbert Accuracy.

Regarding the RadGraph F1 metric, drawing inspiration from the methodology of Delbrouck et al. (2022), we employ the pretrained entity and relation extraction model provided within the dataset. This model has been conveniently made accessible as an installable package (<https://pypi.org/project/radgraph/>). Specifically, we assess the RadGraph F1 Partial variant included in the package, as it was the officially designated variant for the First Shared Task on Clinical Text Generation: RRG24 (Xu et al., 2024). Additionally, we explore another variant, which we term RadGraph F1 Full. This variant is based on the underlying model of the package. We achieve this by creating a "bag" comprising entities, relations without type, and relations with type, and then computing the F1 score between the "bag" of a referenced report and that of a generated report.

In addition to the results that were already pre-

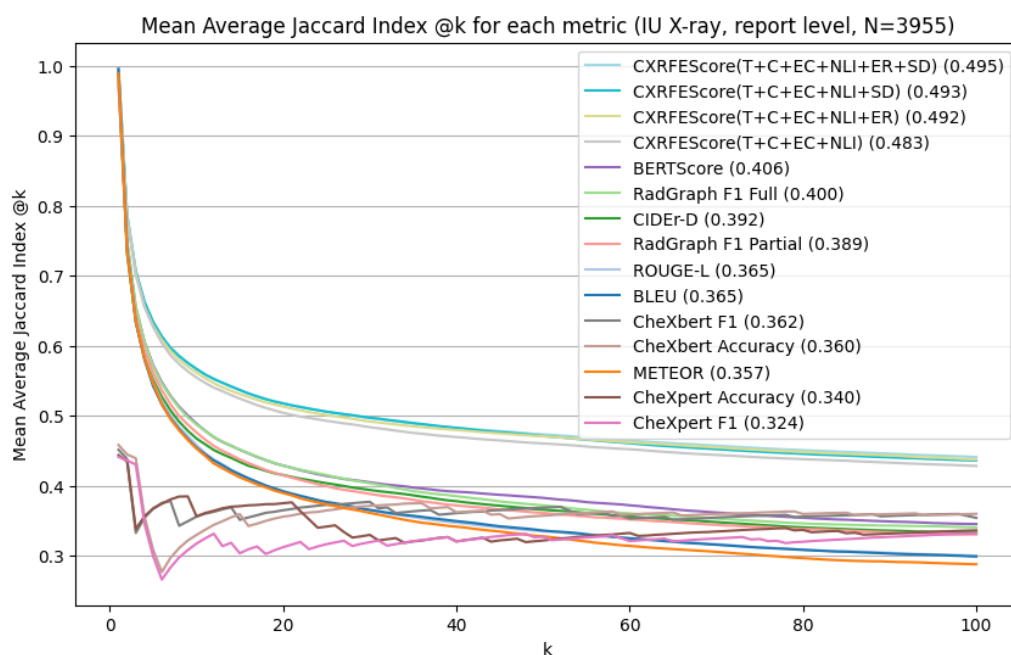


Figure 10: Mean average Jaccard Index at k, for the 3955 reports of the IU X-ray dataset. Larger Jaccard Index is better. The Jaccard Index is calculated by comparing bags of words obtained from the manual and automatic tags associated with each report in the IU X-ray dataset.

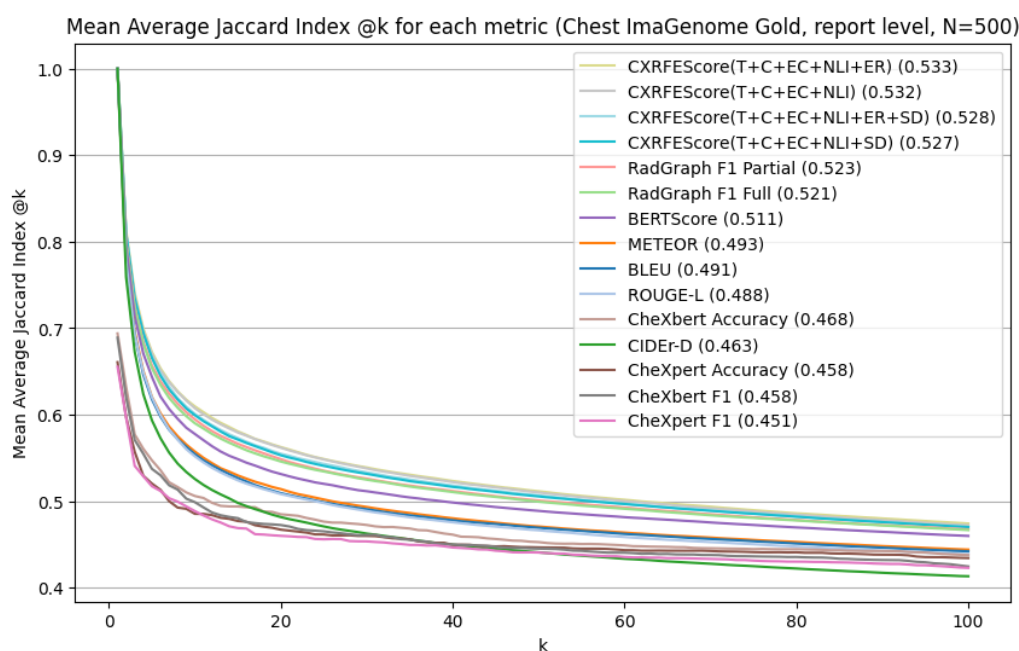


Figure 11: Mean average Jaccard Index at k, for the 500 reports in the gold dataset of Chest ImaGenome. Larger Jaccard Index is better.

sented in Table 4 (Section 4), in this appendix we include the plots shown in Figures 10, 11, 12, 13.

## A.6 Hardware and Other Implementation Details

All of our experiments are implemented using Python 3.10.10 with PyTorch version 1.13.1+cu117 (Paszke et al., 2017). All experiments are conducted on a computing node equipped with a

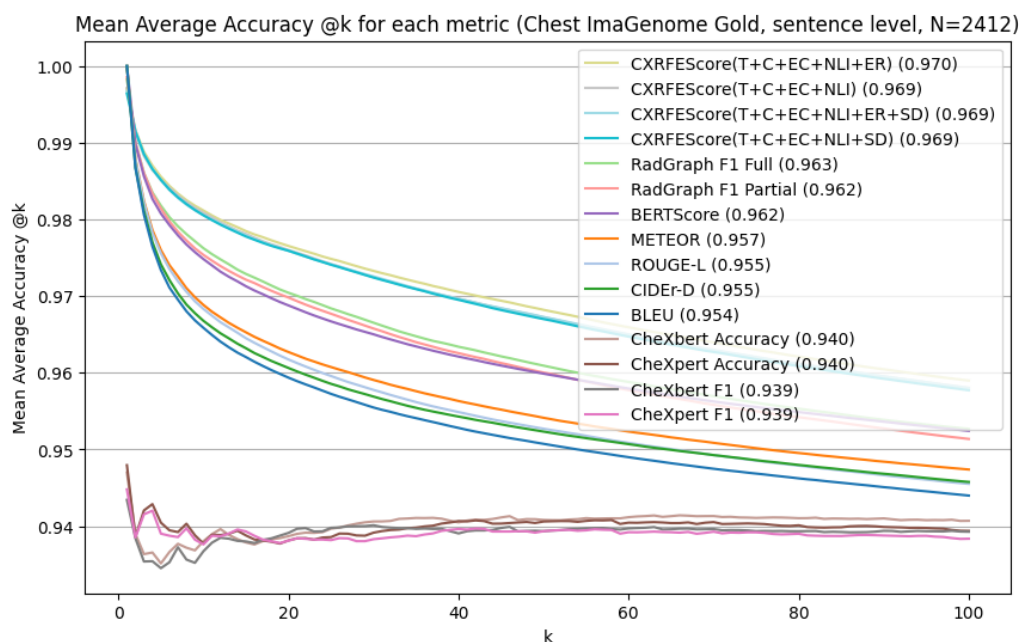


Figure 12: Mean average accuracy at k, for 2412 sentences in the gold dataset of Chest ImaGenome. Larger accuracy is better.

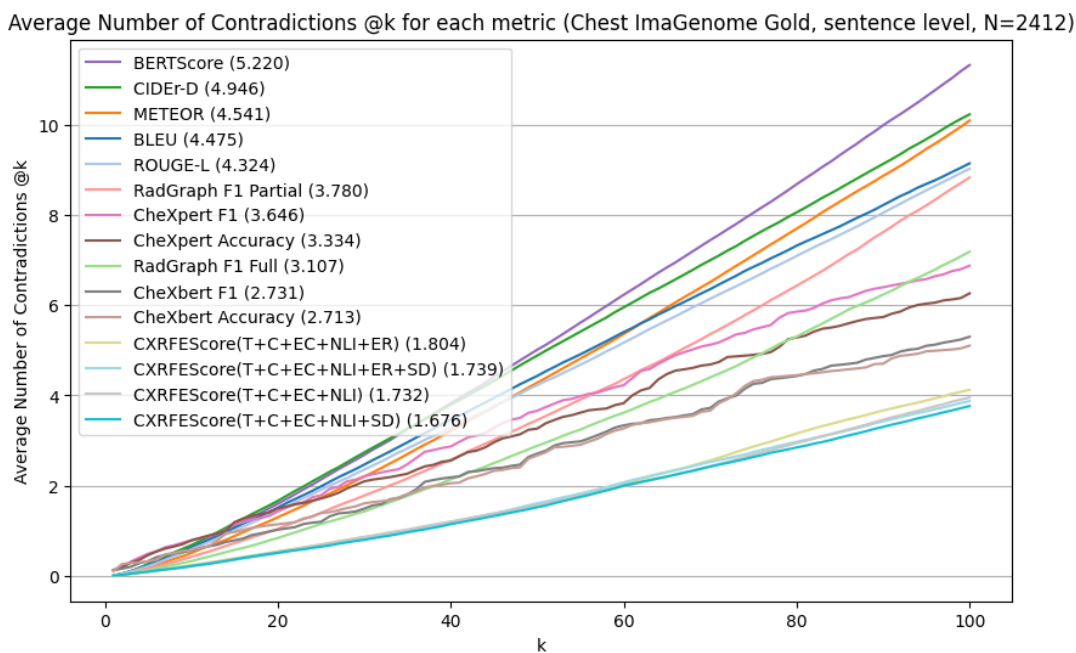


Figure 13: Average number of contradictions at k, for 2412 sentences in the gold dataset of Chest ImaGenome. Smaller number of contradictions is better.

20-core Intel(R) Core(TM) i9-9900X CPU @ 3.50GHz, three NVIDIA GPUs - two GeForce RTX 2080 Ti with 11GB memory and one GeForce RTX 3090 with 24GB memory. The system is complemented by 125GB of RAM.

We implement multitask learning for CXRFE using interleaved dataloaders, multiple model for-

ward passes, and multiple gradient accumulation steps. Specifically, our model features distinct forward functions for each task, with each task assigned its own dataloader. These dataloaders are interleaved according to weights that determine the sampling frequency for each task. To ensure all tasks contribute to the gradients during training, we

employ enough gradient accumulation steps so that each task's batch is sampled at least once before performing backpropagation.

We use the AdamW optimizer (Loshchilov and Hutter, 2019) with a cyclic exponential learning rate that varies from  $8e-5$  to  $1e-6$  over 8 epochs. Here, an epoch consists of roughly 800 batches. Typically, our experiments run for 12-18 hours, after which we observe no significant gains in validation metrics.

## **A.7 ChatGPT prompts**



## Playground

The screenshot displays the OpenAI Playground interface with three panels for system prompts, a chat history, and a settings sidebar.

**System Prompt Panel 1:**

**SYSTEM**

Relevant facts:

1. observations of abnormalities
2. observations of diseases
3. observations of strange visual patterns
4. observations of devices
5. observations of foreign bodies
6. observations of specific anatomical regions that look normal or healthy
7. absences of abnormalities (usually expressed with a negation)
8. comparisons with respect to a previous study (something changed or remained the same)

Task:

Given a sentence taken from a chest x-ray report, generate a JSON list of relevant facts. Each fact should be about one observation. If a sentence mentions multiple observations, each observation should be extracted as a separate fact.

Each fact should include the anatomical

**System Prompt Panel 2:**

**SYSTEM**

Each fact should include the anatomical location where it was observed. If multiple facts occur in the same location, repeat the location in each fact.

If no relevant facts are mentioned, return [] (an empty array).

Examples:

Opacity and density in the right lobe

```
[  
  "opacity in the right lobe",  
  "density in the right lobe"  
]
```

Lungs are well inflated without evidence of focal airspace consolidation to suggest pneumonia.

```
[  
  "well inflated lungs",  
  "lungs without evidence of focal airspace consolidation",  
  "lungs without evidence of pneumonia"  
]
```

**System Prompt Panel 3:**

**SYSTEM**

```
consolidation",  
"lungs without evidence of pneumonia"  
]
```

Taken together, compared with less than 1 hr earlier, the findings are suggestive of worsening of CHF, with new or significantly increased left greater right pleural effusions and underlying bibasilar collapse and/or consolidation, particularly on the left.

```
[  
  "worsening of CHF",  
  "new or significantly increased left pleural effusions",  
  "new or significantly increased right pleural effusions",  
  "underlying bibasilar collapse on the left",  
  "underlying consolidation on the left",  
]
```

No acute cardiopulmonary abnormality

```
[  
  "no acute cardiopulmonary abnormality"  
]
```

**Chat History:**

**USER** There is no pneumothorax, focal consolidation, or pleural effusion

**ASSISTANT** [  
 "no pneumothorax",  
 "no focal consolidation",  
 "no pleural effusion"  
]

**Settings Sidebar:**

- Mode: Chat
- Model: gpt-4-0613
- Temperature: 0
- Maximum length: 512
- Stop sequences: Enter sequence and press Tab
- Top P: 1
- Frequency penalty: 0
- Presence penalty: 0

API and Playground requests will not be used to train our models. [Learn more](#)

Submit Give us feedback

Figure 14: Screenshot of OpenAI's playground web interface showing the prompt used to extract facts from chest X-ray report sentences. The screenshot was manipulated in order to display the full instructions given to ChatGPT.

**SYSTEM**

Given a fact extracted from a Chest X-ray report, output a JSON array of strings. Each string must be a paraphrased fact, with similar semantics, expressing the same finding/diagnosis. These paraphrases will be used to train a fact embedding with triplet loss.

Generate multiple examples (at least 14), covering a wide variety and diversity of radiological terms, synonyms, abbreviations and expressions that radiologists commonly use to communicate the same idea.

Output format:

```
[
  "paraphrase 1",
  "paraphrase 2",
  ...
  "paraphrase 14"
]
```

**USER**      osteoporosis

---

**ASSISTANT**    [

- "Bone thinning is observed",
- "Signs of bone loss are present",
- "Indications of decreased bone density",
- "Evidence of porous bones",
- "Presence of brittle bones",
- "Manifestations of bone weakening",
- "Indications of bone fragility",
- "Signs of demineralized bones",
- "Presence of low bone mass",
- "Manifestations of bone softening",
- "Evidence of bone rarefaction",
- "Indications of skeletal fragility",
- "Signs of reduced bone strength",
- "Presence of bone tissue loss"

]

Add message

Submit ↻

Mode  
Chat

Model  
gpt-4-0613

Temperature      0

Maximum length    512

Stop sequences  
Enter sequence and press Tab

Top P                1

Frequency penalty    0

Presence penalty    0

API and Playground requests will not be used to train our models. [Learn more](#)

[Give us feedback](#)

Figure 15: ChatGPT prompt. Fact to paraphrases

**SYSTEM**

Given a sentence referring to an anatomical location in the context of a chest X-ray, output a JSON array of strings paraphrasing it, covering a wide diverse range of terminology, synonyms and abbreviations that radiologists commonly use to express the same idea in a chest X-ray report.

Examples:

gastroesophageal junction

```
[
  "cardia",
  "cardia of the stomach",
  "cardiac orifice",
  "esophagogastric junction",
  "GE junction",
  "esophagogastric junction",
  "esophagogastric junctional region",
  "gastroesophageal interface",
  "gastroesophageal transition zone",
  "lower esophageal junction",
  "esophagogastric transition",
  "stomach-esophagus interface",
  "junction between stomach and
```

**SYSTEM**

```
"junction between stomach and
esophagus",
"cardioesophageal junction",
"junction of the esophagus and stomach",
"transition between the esophagus and
stomach",
"point where the esophagus meets the
stomach",
"Z-line",
"distal end of the esophagus",
"proximal start of the gastric cardia"
]

costodiaphragmatic recess
[
  "costophrenic angle",
  "costophrenic recess",
  "costophrenic sulcus",
  "costophrenic notch",
  "costophrenic sinus",
  "phrenicocostal sinus",
  "diaphragm-chest wall interface",
  "junction of diaphragm and chest wall",
  "area where the diaphragm meets the
ribs",
```

**SYSTEM**

```
"costal margin of the diaphragm",
"angle formed by diaphragm and rib cage"
]

R>L
[
  "right greater than left"
]

base
[
  "lung base",
  "base of the lung",
  "lower lung area",
  "basal lung region",
  "basal portion of the lung",
  "lung's basal part",
  "bottom part of the lung",
  "lower region of the lung",
  "pulmonary base",
  "base of the pulmonary structure",
  "lower zone of the lung",
  "lower pulmonary region",
  "lower segment of the lung"
]
```

**USER** cardiomediastinal silhouette

---

**ASSISTANT** [

```
"heart and mediastinal contour",
"cardiac and mediastinal silhouette",
"outline of the heart and mediastinum",
"cardiac and mediastinal shadow",
"cardiopulmonary silhouette",
"heart and mediastinal structure",
"cardiac and mediastinal shape",
"cardiac and mediastinal border",
"cardiac and mediastinal profile",
"cardiac and mediastinal image",
"cardiac and mediastinal area",
"cardiac and mediastinal configuration",
"cardiac and mediastinal pattern",
"cardiac and mediastinal appearance",
"cardiac and mediastinal contouring",
"cardiac and mediastinal form",
"cardiac and mediastinal outline",
"cardiac and mediastinal figure",
"cardiac and mediastinal boundary"
]
```

[Submit](#) ↻

Mode

Model

Temperature

Maximum length

Stop sequences

Top P

Frequency penalty

Presence penalty

API and Playground requests will not be used to train our models. [Learn more](#)

Figure 16: ChatGPT prompt. Anatomy location to paraphrases

**SYSTEM**

Given a medical fact, output a JSON object with 7 fields:

- "anatomical location"
- "detailed observation"
- "short observation"
- "category"
- "health status"
- "prev\_study\_comparison?"
- "comparison status"

Definitions:

- "anatomical location" means the anatomical location of the observation. If not given, it should be "".
- "detailed observation" means what was observed, with all the details provided (but excluding location, unless the location itself is the main observation).
- "short observation" should be a brief summary of the previous field, without extra details, but logically consistent.
- "category" can be "anatomical finding", "disease", "technical assessment", "tubes and lines" or "device".

**SYSTEM**

5. "health status" can be "normal", "abnormal", "ambiguous" or "unknown".  
"normal" means healthy.  
"abnormal" means unhealthy.  
"ambiguous" means it is not clear whether it is healthy or unhealthy, leaving room for interpretation.  
"unknown" means the observation provides no information about abnormality/normality.

6. "prev\_study\_comparison?" is a boolean field. Can be either "yes" or "no". If the fact is comparing the current study with a previous study (something improved, worsened, changed or remained the same), return "yes". Otherwise, return "no".

7. "comparison status" can be "resolved", "new", "improved", "worsened", "larger", "smaller", "displaced", or "". If no comparison with a previous study is made, the default is "" (empty).

Examples:

```
small to moderate size left pleural effusion
{
  "anatomical location": "left",
  "detailed observation": "small to moderate size pleural effusion",
  "short observation": "pleural effusion",
```

**SYSTEM**

```
"category": "anatomical finding",
"health status": "abnormal",
"prev_study_comparison?": "no",
"comparison status": ""
}

new finding of mass in the abdomen
{
  "anatomical location": "abdomen",
  "detailed observation": "mass",
  "short observation": "mass",
  "category": "anatomical finding",
  "health status": "abnormal",
  "prev_study_comparison?": "yes",
  "comparison status": "new"
}

clear lung fields
{
  "anatomical location": "lung fields",
  "detailed observation": "clear lung fields",
  "short observation": "clear lung fields",
  "category": "anatomical finding",
  "health status": "normal",
  "prev_study_comparison?": "no",
  "comparison status": ""
}

no ill-definition of right hemidiaphragm
```

**USER**

bilateral effusions equivocally getting worse

⊖

**ASSISTANT**

```
{
  "anatomical location": "bilateral",
  "detailed observation": "effusions equivocally getting worse",
  "short observation": "worsening effusions",
  "category": "anatomical finding",
  "health status": "abnormal",
  "prev_study_comparison?": "yes",
  "comparison status": "worsened"
}
```

Mode

Chat

Model

gpt-4-0613

Temperature

0

Maximum length

512

Stop sequences  
Enter sequence and press Tab

Figure 17: ChatGPT prompt. Fact to metadata



**SYSTEM**

Given a statement extracted from a chest X-ray report, output a comparison category. The category must be one item of the following list:

0. no comparison
1. new finding
2. resolved
3. improved
4. worsened
5. progressed
6. reappeared
7. larger
8. smaller
9. increase
10. decrease
11. position changed
12. stable/unchanged
13. unclear comparison
14. other

Examples:

|  
right pleural effusion similar in size  
12. stable/unchanged

pleural density on the left is somewhat increased  
9. increase

new bibasal consolidations on the left

**SYSTEM**

6. reappeared
7. larger
8. smaller
9. increase
10. decrease
11. position changed
12. stable/unchanged
13. unclear comparison
14. other

Examples:

|  
right pleural effusion similar in size  
12. stable/unchanged

pleural density on the left is somewhat increased  
9. increase

new bibasal consolidations on the left  
1. new finding

small foci of opacity in the left mid lung  
0. no comparison

no ill-definition of right hemidiaphragm  
0. no comparison

heart within upper limits of normal in size  
0. no comparison

---

**USER**      bilateral effusions equivocally getting worse

---

**ASSISTANT**    4. worsened

Mode  
Chat

---

Model  
gpt-4-0613

---

Temperature      0

⊕ Add message

Figure 18: ChatGPT prompt. Fact to comparison status

**SYSTEM**

Instructions:

Apply an annotation standard of observations to raw phrases extracted from chest X-ray reports. The standard considers the following observations:

- airspace opacity
- atelectasis
- bone lesion
- bronchiectasis
- calcified nodule
- clavicle fracture
- consolidation
- costophrenic angle blunting
- cyst/bullae
- diaphragmatic eventration (benign)
- elevated hemidiaphragm
- enlarged cardiac silhouette
- enlarged hilum
- hernia
- hydropneumothorax
- hyperaeration
- increased reticular markings/ijq pattern
- infiltration
- linear/patchy atelectasis
- lobar/segmental collapse
- lung lesion
- lung opacity
- mass/nodule (not otherwise specified)

**SYSTEM**

- mediastinal displacement
- mediastinal widening
- multiple masses/nodules
- pleural effusion
- pleural/parenchymal scarring
- pneumomediastinum
- pneumothorax
- pulmonary edema/hazy opacity
- rib fracture
- scoliosis
- shoulder osteoarthritis
- spinal degenerative changes
- spinal fracture
- sub-diaphragmatic air
- subcutaneous air
- superior mediastinal mass/enlargement
- tortuous aorta
- vascular calcification
- vascular congestion
- vascular redistribution
- aortic graft/repair
- cabg grafts
- cardiac pacer and wires
- prosthetic valve
- alveolar hemorrhage
- aspiration
- coppd/emphysema
- fluid overload/heart failure
- goiter

**SYSTEM**

- granulomatous disease
- interstitial lung disease
- lung cancer
- pericardial effusion
- pneumonia
- artifact
- breast/nipple shadows
- low lung volumes
- rotated
- skin fold
- chest port
- chest tube
- endotracheal tube
- enteric tube
- ij line
- intra-aortic balloon pump
- mediastinal drain
- pigtail catheter
- subclavian line
- swan-ganz catheter
- tracheostomy tube

Output format:

JSON array of strings

Rules:

Only output the observations from the standard that best match or are supported by the phrase. If none match the phrase, output an empty array ([]).

Example:

Increased bibasilar opacities likely representing atelectasis, but cannot exclude aspiration or pneumonia in the correct clinical setting

```
[
  "lung opacity",
  "atelectasis",
  "aspiration",
  "pneumonia"
]
```

**SYSTEM**

- intra-aortic balloon pump
- mediastinal drain
- pigtail catheter
- subclavian line
- swan-ganz catheter
- tracheostomy tube

Output format:

JSON array of strings

Rules:

Only output the observations from the standard that best match or are supported by the phrase. If none match the phrase, output an empty array ([]).

Example:

Increased bibasilar opacities likely representing atelectasis, but cannot exclude aspiration or pneumonia in the correct clinical setting

```
[
  "lung opacity",
  "atelectasis",
  "aspiration",
  "pneumonia"
]
```

**USER**      dual chamber pacemaker

---

**ASSISTANT**    [

"cardiac pacer and wires"

]

Mode: Chat

Model: gpt-4-0613

Temperature: 0

Maximum length: 512

Figure 19: ChatGPT prompt. Fact to Chest ImaGenome observations

**SYSTEM**

Instructions:

Apply an annotation standard of anatomical locations to raw phrases extracted from chest X-ray reports. The standard considers the following anatomical locations:

- abdomen
- aortic arch
- cardiac silhouette
- carina
- cavoatrial junction
- left apical zone
- left arm
- left breast
- left chest wall
- left clavicle
- left costophrenic angle
- left hemidiaphragm
- left hilar structures
- left lower lung zone
- left lung
- left mid lung zone
- left shoulder
- left upper lung zone
- mediastinum
- neck
- right apical zone

**SYSTEM**

- right arm
- right atrium
- right breast
- right chest wall
- right clavicle
- right costophrenic angle
- right hemidiaphragm
- right hilar structures
- right lower lung zone
- right lung
- right mid lung zone
- right shoulder
- right upper lung zone
- spine
- svc
- trachea
- upper mediastinum

Output format:

JSON array of strings

Rules:

Only output the anatomical locations from the standard that best match the given phrase. If none match the phrase, output an empty array ([]).

**SYSTEM**

Only output the anatomical locations from the standard that best match the given phrase. If none match the phrase, output an empty array ([]).

Examples:

cardiomediastinal silhouette

```
[
  "cardiac silhouette",
  "mediastinum"
]
```

A tracheostomy tube is in place

```
[
  "neck",
  "trachea"
]
```

bilateral pulmonary nodules

```
[
  "left lung",
  "right lung"
]
```

left IJV

```
[
  "neck"
]
```

**USER**      dual chamber pacemaker

---

**ASSISTANT**    [

"right atrium",

"right ventricle"

]

Mode: Chat

Model: gpt-4-0613

Temperature: 0

Maximum length: 512

Figure 20: ChatGPT prompt. Fact to Chest ImaGenome anatomical locations

**SYSTEM**

Given an anchor, output a list of positives and a list of negatives.

The anchor will come from a chest X-ray report. The positives must be 6 paraphrases of the anchor, expressing the same idea with different terms and synonyms (cover a wide range of medical terminology).

The negatives must be 10 hard negatives, that is, similar to the anchor but semantically different, for example, by expressing a different diagnosis (if applicable), by referring to a different anatomical location (if applicable), etc.

Output format: a JSON object as follows

```
{
  "positives": [ ... ],
  "negatives": [ ... ]
}
```

**USER**      removal of endotracheal tube

---

**ASSISTANT**    {

```
"positives": [
  "Endotracheal tube has been removed",
  "Extraction of the endotracheal tube",
  "Endotracheal tube has been taken out",
  "Endotracheal tube has been extracted",
  "Endotracheal tube removal has been performed",
  "Discontinuation of the endotracheal tube"
],
"negatives": [
  "Insertion of endotracheal tube",
  "Removal of nasogastric tube",
  "Removal of chest tube",
  "Endotracheal tube is in place",
  "Endotracheal tube repositioned",
  "Removal of central venous catheter",
  "Endotracheal tube is unchanged",
  "Removal of tracheostomy tube",
  "Endotracheal tube is malpositioned",
  "Removal of urinary catheter"
]
}
```

Mode  
Chat

Model  
gpt-4-0613

Temperature 0

Maximum length 512

Stop sequences  
Enter sequence and press Tab

Top P 1

Frequency penalty 0

Presence penalty 0

API and Playground requests will not be used to train our models. [Learn more](#)

Figure 21: ChatGPT prompt. Hard triplets

**SYSTEM**

Given a premise from a Chest X-ray report, output an entailment, a contradiction and a neutral sentence.

The entailment must strictly follow logically from the premise. The neutral can optionally be about a different chest X-ray observation. The output format must be a JSON as follows:

```
{
  "e": ...,
  "c": ...,
  "n": ...
}
```

**USER**      removal of endotracheal tube

---

**ASSISTANT**    {

```
"e": "The patient no longer has an endotracheal tube.",
"c": "The endotracheal tube is still in place.",
"n": "The patient's lungs appear clear."
}
```

⊕ Add message

Mode  
Chat

Model  
gpt-4-0613

Temperature 0

Maximum length 512

Stop sequences  
Enter sequence and press Tab

Figure 22: NLI Prompt 1. Given a premise, output entailment, contradiction, and neutral hypotheses.

**SYSTEM**

You will receive a NLI example from a chest X-ray dataset with premise, hypothesis and label ("entailment", "contradiction" or "neutral"). Generate 10 new NLI examples with the same label and following a similar writing style as the example provided. They have to be about chest X-ray reports. Output the examples as a JSON array of objects. If the premise or the hypothesis (or both) are complex (i.e. they say several things), create challenging examples following a similar structure as the provided example that require a good understanding of logic in order to deduce the label.

Avoid verbosity and words such as "patient" or "X-ray". State observations succinctly.

**USER**      {"P": "No pleural effusions or pneumothorax.", "H": "Previously seen pneumothorax is no longer visualized.", "L": "entailment"}

---

**ASSISTANT**    [

```
{ "P": "No evidence of pulmonary edema or pneumonia.", "H": "The previously noted pneumonia is not present.", "L": "entailment" },
{ "P": "No signs of lung nodules or masses.", "H": "The lung mass seen before is not visible now.", "L": "entailment" },
{ "P": "No visible signs of cardiomegaly or mediastinal widening.", "H": "The previously observed cardiomegaly is not seen.", "L": "entailment" },
{ "P": "No signs of pulmonary embolism or pleural effusion.", "H": "The earlier detected pleural effusion is not visible.", "L": "entailment" },
{ "P": "No evidence of pneumothorax or lung collapse.", "H": "The previously seen lung collapse is not present.", "L": "entailment" },
{ "P": "No signs of rib fractures or lung contusions.", "H": "The rib fracture noted before is not visible now.", "L": "entailment" },
{ "P": "No visible signs of bronchiectasis or lung fibrosis.", "H": "The previously observed bronchiectasis is not seen.", "L": "entailment" },
{ "P": "No signs of lung cancer or tuberculosis.", "H": "The earlier detected tuberculosis is not visible.", "L": "entailment" },
{ "P": "No evidence of aortic dissection or aneurysm.", "H": "The previously seen aortic aneurysm is not present.", "L": "entailment" },
{ "P": "No signs of pericardial effusion or heart failure.", "H": "The
```

Mode  
Chat

Model  
gpt-4-0613

Temperature 0

Maximum length 512

Stop sequences  
Enter sequence and press Tab

Top P 1

Frequency penalty 0

Presence penalty 0

API and Playground requests will not be used to train our models. [Learn more](#)

Figure 23: NLI Prompt 2. Given a ground-truth NLI example, generate multiple similar examples.

**SYSTEM**

Context: natural language inference.

Given a premise and a hypothesis, output "entailment", "contradiction", or "neutral".

Use "entailment" when the facts stated by the premise necessarily entail the truth of the hypothesis.

Use "contradiction" when premise and hypothesis are mutually exclusive/contradictory (both cannot be true at the same time).

Use "neutral", if there is no contradiction (premise and hypothesis are compatible), but the premise does not entail the hypothesis (it's possible for the premise to be true and the hypothesis still be false). In other words, use "neutral" when neither "entailment" nor "contradiction" adequately fit.

**USER**      Premise: Lungs are well-expanded and clear without focal consolidation concerning for pneumonia. | Hypothesis: Lungs are hyperinflated but clear.

---

**ASSISTANT**    Neutral

---

**USER**      Enter a user message here.

⊕ Add message

Mode  
Chat

Model  
gpt-4-0613

Temperature      0

Maximum length    10

Stop sequences  
Enter sequence and press Tab

Top P              1

Frequency penalty   0

Presence penalty    0

Figure 24: NLI Prompt 3. Given a premise and a hypothesis, output the correct label. In this version, GPT-4 has to output the label immediately.

**SYSTEM**

Context: natural language inference.

Given a premise (#P) and a hypothesis (#H), output "Reason: {reason}. Label: {label}" where {reason} is a short sentence and {label} is one of "entailment", "contradiction," or "neutral."

Use "entailment" when the premise necessarily entails the truth of the hypothesis.

Use "contradiction" when premise and hypothesis are mutually exclusive/contradictory. Pay attention to logical inconsistencies, such as expressions suggesting presence vs. absence, etc.

Use "neutral" when there's no contradiction, but the premise doesn't necessarily entail the hypothesis.

Examples:

- #P: Increased pulmonary edema. | #H: worsened pulmonary edema. Label: entailment
- #P: No pulmonary edema, consolidation, or pneumothorax. | #H: No focal consolidation, pleural effusion, or pneumothorax is present. Label: neutral

**USER**      #P: Increased opacity is seen in the left peritracheal region and left lung base | #H: There are signs of abnormality in the left lung base and left peritracheal region

---

**ASSISTANT**    Reason: The premise states that there is increased opacity in the left peritracheal region and left lung base, which is a sign of abnormality. Therefore, the hypothesis is a rephrase of the premise. Label: entailment.

---

⊕ Add message

Model  
gpt-4-0613

Temperature      0

Maximum length    256

Stop sequences  
Enter sequence and press Tab

Top P              1

Frequency penalty   0

Presence penalty    0

ⓘ API and Playground requests will not be used to train our models. [Learn more](#)

Submit
Give us feedback

Figure 25: NLI Prompt 4. Given a premise and a hypothesis, output the correct label. In this version, GPT-4 is given the chance to reason first and then output the label. It is also shown a couple of examples with the expected labels. These simple changes to the prompt significantly boosted GPT-4's performance.

**SYSTEM**

Given a premise from a Chest X-ray report, output 5 statements that explicitly logically contradict the premise, that is, the premise and any generated statement cannot both be true at the same time. Include at least one sentence that is very similar to the premise but with a very slight difference, thus creating a logical contradiction.

For example:

- "There is mid thoracic dextroscoliosis" vs. "There is no mid thoracic dextroscoliosis"
- "Left basal consolidation has slightly improved" vs. "Left basal consolidation has slightly increased"

Output format: a JSON array of strings.

**USER**      Lungs are well-expanded and clear without focal consolidation concerning for pneumonia.

---

**ASSISTANT**    ["Lungs are not well-expanded and clear without focal consolidation concerning for pneumonia.", "Lungs are well-expanded but show focal consolidation concerning for pneumonia.", "Lungs are collapsed and clear without focal consolidation concerning for pneumonia.", "Lungs are well-expanded and clear with focal consolidation not concerning for pneumonia.", "Lungs are well-expanded and clear, but there is focal consolidation concerning for pneumonia."]

---

⊕ Add message

Mode  
Chat

Model  
gpt-4-0613

Temperature      0

Maximum length    500

Stop sequences  
Enter sequence and press Tab

Top P              1

Figure 26: NLI Prompt 5. Given a premise, generate diverse hypotheses that contradict it.