

# Making the Most of Text Semantics to Improve Biomedical Vision–Language Processing

Benedikt Boecking<sup>\*†</sup>, Naoto Usuyama<sup>\*</sup>, Shruthi Bannur, Daniel C. Castro, Anton Schwaighofer, Stephanie Hyland, Maria Wetscherek, Tristan Naumann, Aditya Nori, Javier Alvarez-Valle, Hoifung Poon, and Ozan Oktay<sup>‡</sup>

*Microsoft Health Futures, Redmond, WA, USA*

*Microsoft Research, Cambridge, UK*

## Abstract

Multi-modal data abounds in biomedicine, such as radiology images and reports. Interpreting this data at scale is essential for improving clinical care and accelerating clinical research. Biomedical text with its complex semantics poses additional challenges in vision-language modelling compared to the general domain, and previous work has used insufficiently adapted models that lack domain-specific language understanding. In this paper, we show that principled textual semantic modelling can substantially improve contrastive learning in self-supervised vision–language processing. We release a language model that achieves state-of-the-art results in radiology natural language inference through its improved vocabulary and novel language pretraining objective leveraging semantics and discourse characteristics in radiology reports. Further, we propose a self-supervised joint vision–language approach with a focus on better text modelling. It establishes new state of the art results on a wide range of publicly available benchmarks, in part by leveraging our new domain-specific language model. We release a new dataset with locally-aligned phrase grounding annotations by radiologists to facilitate the study of complex semantic modelling in biomedical vision–language processing. A broad evaluation, including on this new dataset, shows that our contrastive learning approach, aided by textual-semantic modelling, outperforms prior methods in segmentation tasks, despite only using a global-alignment objective.

## 1 Introduction

Advances in deep learning have enabled automated diagnosis systems that operate near or above expert-level performance, paving the way for the use of machine learning systems to improve healthcare workflows, for example by supporting fast triaging and assisting medical professionals to reduce errors and omissions [9, 19, 53, 71]. A major hurdle to the widespread development of these systems is a requirement for large amounts of detailed ground-truth clinical annotations for supervised training, which are expensive and time-consuming to obtain. Motivated by this challenge, there has been a rising interest in multi-modal self-supervised learning [44, 30] and cross-modal weak supervision [71, 75, 32, 18, 20] (using partial and imperfect image labels derived from the auxiliary modality), in particular for paired image–text data. Such data is collected daily in routine clinical practice, and common examples are X-ray images [18, 32, 75] or computed tomography scans [9, 18, 20, 71] paired with reports written by qualified medical experts. Importantly, while many remain private, some paired clinical datasets have been released to the research community including MIMIC-CXR [33], Open-I [14], and PadChest [3].

This article focuses on self-supervised vision–language processing (VLP) for paired image and text data in the biomedical domain. The goal is to jointly learn good image and text representations that can be leveraged by downstream applications such as zero-/few-shot image classification, report generation and error detection, and disease localisation. Self-supervised VLP has several advantages over supervised learning, not

---

<sup>\*</sup>These authors contributed equally.

<sup>†</sup>The work was conducted during Benedikt Boecking’s summer internship at Microsoft Research.

<sup>‡</sup>Corresponding author: [ozan.oktay@microsoft.com](mailto:ozan.oktay@microsoft.com)

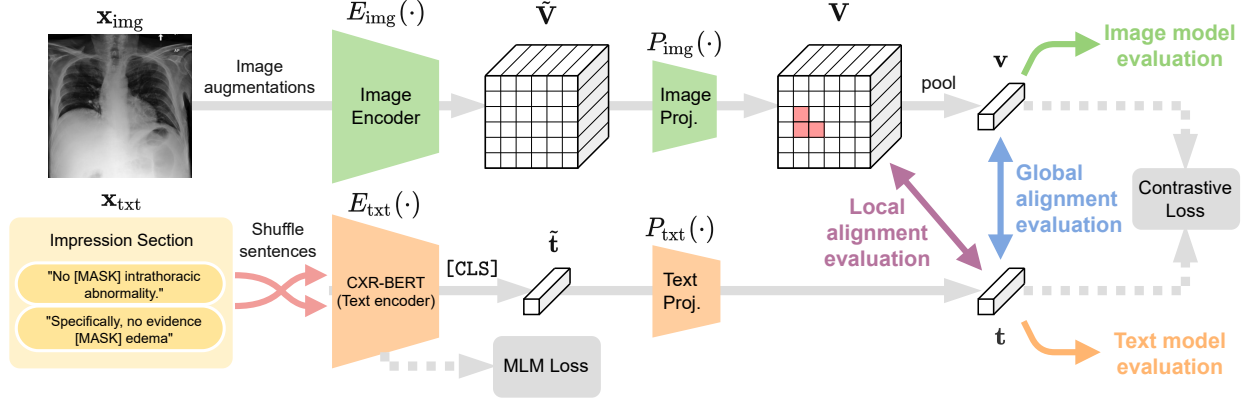


Figure 1: BioViL leverages our radiology-specific text encoder (CXR-BERT), text augmentation, regularisation, and maintains language model quality via a masked language modelling (MLM) loss. We conduct a broad evaluation of models and representations that includes zero-shot classification, phrase grounding, and natural language inference.

just because it does not require laborious manual annotations, but also because it does not operate on a fixed number of predetermined conditions or object categories, since the joint latent space is learned from raw text.

However, in contrast to the general domain setting, self-supervised VLP with biomedical data poses additional challenges. Take radiology as an example, publicly available datasets [33, 14, 3] are usually smaller, on the order of a few hundred thousand pairs rather than millions in general-domain vision–language processing (e.g. [60] collected 400M text–image pairs on the Internet for self-supervision). Furthermore, linguistic challenges are different in biomedical settings, including common usage of negations, expressions of uncertainty, long-range dependencies, more frequent spatial relations, the use of domain-specific modifiers, as well as scientific terminology rarely found in the general domain. Taking negation as an example, “there is no dog in this picture” would be a highly unusual caption on social media, but “there is no evidence of pneumonia in the left lung” or “there are no new areas of consolidation to suggest the presence of pneumonia” are descriptions commonly found in radiology reports. Moreover, pretrained models including object detectors often used in general domain visual grounding are typically unavailable or under-perform in domain-specific applications (see also Supp. in [30]). Additionally, imbalance in underlying latent entities of interest (e.g., pulmonary findings) can cause larger numbers of false negatives in contrastive learning objectives that sample at random, which can lead models to degrade and memorise irrelevant text and image aspects. For example, radiology images and text reports with normal findings occur much more frequently compared to exams that reveal abnormal conditions such as pneumonia or pneumothorax (also see [10]). Supp. B.1 provides further discussion of these challenges.

Related self-supervised VLP work [29, 84, 44, 30, 55] has achieved impressive downstream classification and zero-shot classification performance. However, our study reveals that suboptimal text modelling due to insufficient vocabulary adjustment, fine-tuning, and language grounding appears to have gone unnoticed, all of which are shown to degrade the quality of joint latent representations. In particular, a more thorough benchmarking of the text, image, and shared embeddings, across a multitude of downstream benchmarks, reveals that large improvements in performance are possible by taking care to build highly specialised text models and by maintaining their performance during joint training. Free-text image descriptions provide a semantically dense learning signal compared to image-only contrastive methods and supervised classification [15]. Further, extracting shared semantics of images and text pairs is easier for text, as the modality is already discretised. Thus, making the most of text modelling before and during joint training can lead to large improvements in not just the text model, but also of the image model and joint representations. We present the following contributions in this work:

1. We introduce and release a new chest X-ray (CXR) domain-specific language model, CXR-BERT<sup>1</sup>(Fig. 2). Through an improved vocabulary, a novel pretraining procedure, regularisation, and text augmentation, the model considerably improves radiology natural language inference [53], radiology masked

<sup>1</sup>The pretrained model weights of CXR-BERT will soon be released via HuggingFace.

token prediction [16, 47], and downstream VLP task performance.

2. We propose and release a simple but effective self-supervised VLP approach for paired biomedical data which we name BioViL (Fig. 1) and evaluate in the radiology setting. Through improvements in text modelling, text model grounding, augmentation, and regularisation, the approach yields new state-of-the-art performance on a wide range of public downstream benchmarks. Our large-scale evaluation (see Table 2) includes phrase grounding, natural language inference [53], as well as zero-/few-shot classification and zero-shot segmentation via the RSNA Pneumonia dataset [65, 75]. Notably, our approach achieves improved segmentation performance despite only using a global alignment objective during training.
3. We also release a novel biomedical phrase grounding dataset, namely **MS-CXR**<sup>2</sup>, to encourage reproducible evaluation of shared latent semantics learned by biomedical image-text models. This large, well-balanced phrase grounding benchmark dataset contains carefully curated image regions annotated with descriptions of eight radiology findings, as verified by board-certified radiologists. Unlike existing chest X-ray benchmarks, this challenging phrase grounding task evaluates joint, local image-text reasoning while requiring real-world language understanding, e.g. to parse domain-specific location references, complex negations, and bias in reporting style.

## 2 Making the Most of Free-Text Supervision

We assume that we are given a set  $\mathcal{D}$  of pairs of radiology images and reports  $(\mathbf{x}_{\text{img}}, \mathbf{x}_{\text{txt}})$ . Let  $\mathbf{w} = (w_1, \dots, w_T)$  denote a vector of  $T$  (sub-)word tokens of a text document  $\mathbf{x}_{\text{txt}}$  (after tokenisation). Recall that a BERT [72] encoder  $E_{\text{txt}}$  outputs a feature vector for each input token  $w_t$  as well as a special global [CLS] token used for downstream classification. Let  $\tilde{\mathbf{t}} = [E_{\text{txt}}(\mathbf{w})]_{[\text{CLS}]}$  denote the [CLS] token prediction by  $E_{\text{txt}}$  based on input  $\mathbf{w}$ , and  $\mathbf{t} = P_{\text{txt}}(\tilde{\mathbf{t}})$  its lower-dimensional projection by a model  $P_{\text{txt}}$ .

### 2.1 CXR-BERT: Domain-Specific Language Model Pretraining

We introduce and publicly release CXR-BERT (Fig. 2), a specialised CXR language model with an adjusted vocabulary, pretrained in three phases to capture dense semantics in radiology reports [4]. To achieve this specialisation to the CXR report domain despite limited data availability, our approach includes pretraining on larger data from closely related domains. The phases proceed as follows:

(I) First, we construct a custom WordPiece [79] vocabulary of 30k tokens from PubMed abstracts<sup>3</sup> (15 GB), MIMIC-III [34] clinical notes (3.5 GB), and MIMIC-CXR radiology reports (0.1 GB). With this custom vocabulary, our model produces fewer sub-word breakdowns (Table 1).

(II) Second, we pretrain a randomly initialised BERT model via Masked Language Modelling (MLM) on the PubMed + MIMIC-III + MIMIC-CXR corpora. We largely follow RoBERTa [47] pretraining configurations, i.e. dynamic whole-word masking for MLM and packing of multiple sentences into one input sequence. This phase aims to build an initial domain-specific BERT model in the biomedical and clinical domains. (III) Third, we continue pretraining on MIMIC-CXR only to further specialise our CXR-BERT to the CXR domain. Here, we also add a novel sequence prediction task to the objective to obtain better sequence representations, as explained below.

Note that a raw radiology report  $\mathbf{x}_{\text{txt}}$  typically consists of several sections, including a ‘FINDINGS’ section that details clinical observations, and an ‘IMPRESSION’ section summarising the clinical assessment [73, 76]. Our sequence prediction objective of phase (III) aims to take advantage of this structure. Specifically, we continually run MLM pretraining on MIMIC-CXR radiology reports and propose to add a radiology section matching (RSM) pretraining task, formulated to match IMPRESSION to FINDINGS sections of the same study.

Table 1: Vocabulary comparison of common radiology terms with ClinicalBERT (Wiki/Book, cased), PubMedBERT (PubMed, uncased), and CXR-BERT (PubMed+MIMIC-III/CXR, uncased). ✓ marks that a word appears in the vocabulary, otherwise its sub-tokens are shown.

| Full word    | ClinicalBERT     | PubMedBERT      | CXR-BERT |
|--------------|------------------|-----------------|----------|
| pneumonia    | ✓                | ✓               | ✓        |
| opacity      | op-acity         | ✓               | ✓        |
| effusion     | e-ff-usion       | ✓               | ✓        |
| pneumothorax | p-ne-um-oth-orax | ✓               | ✓        |
| atelectasis  | ate-lect-asis    | ate-le-ct-asis  | ✓        |
| cardiomegaly | card-io-me-gal-y | cardio-me-gal-y | ✓        |
| bibasilar    | bi-bas-ila-r     | bib-asi-la-r    | ✓        |

<sup>2</sup>MS-CXR dataset will soon be released in PhysioNet.

<sup>3</sup>Available at <https://pubmed.ncbi.nlm.nih.gov/download/>

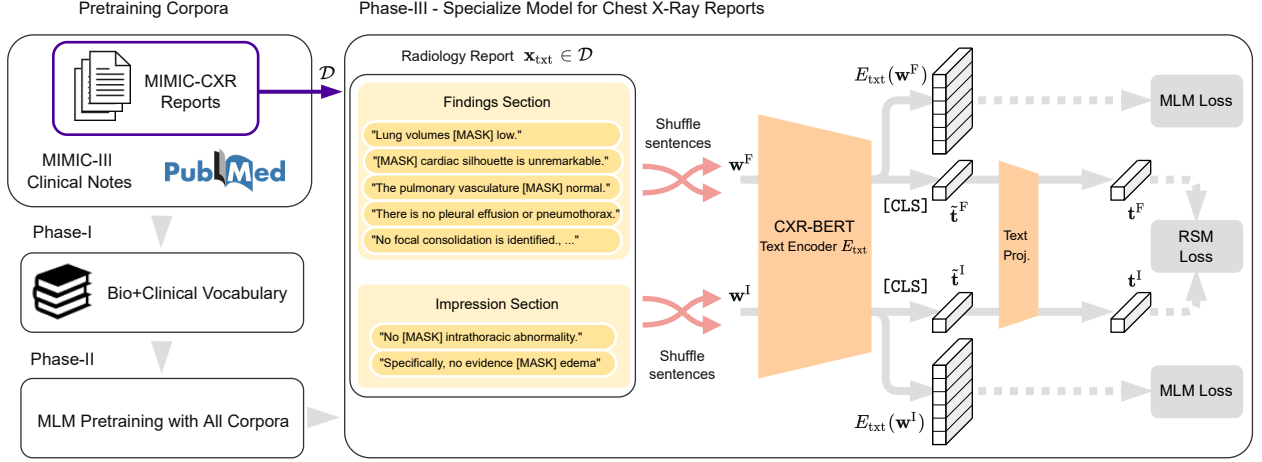


Figure 2: The proposed CXR-BERT text encoder has three phases of pretraining and uses a domain-specific vocabulary, masked language modelling (MLM) and radiology section matching (RSM) losses, regularisation, and text augmentations.

Let  $\theta$  denote the weights of our language model and  $m \subset \{1, \dots, T\}$  denote mask indices for  $M$  masked tokens, randomly sampled for each token vector  $\mathbf{w}$  at every iteration. Given a batch  $\mathcal{B}$  of token vectors  $\mathbf{w} = (w_1, \dots, w_T)$ , we write the MLM loss as the cross-entropy for predicting the dynamically masked tokens:

$$\mathcal{L}_{\text{MLM}} = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{w} \in \mathcal{B}} \log p_{\theta}(\mathbf{w}_m | \mathbf{w}_{\setminus m}). \quad (1)$$

Further, let  $(\tilde{\mathbf{t}}_i^F, \tilde{\mathbf{t}}_i^I)$  denote a pair of [CLS] tokens corresponding to the FINDINGS and IMPRESSION sections of the same  $i^{\text{th}}$  report, and let  $(\mathbf{t}_i^F, \mathbf{t}_i^I)$  denote the pair projected to a lower dimension via a two-layer perceptron  $P_{\text{txt}}$ . We introduce a contrastive loss on the text modality that favours IMPRESSION and FINDINGS text pair from the same report over unmatched ones. Specifically, for a batch of  $N$  such pairs, the RSM loss is defined as

$$\mathcal{L}_{\text{RSM}} = -\frac{1}{N} \sum_{i=1}^N \left( \log \frac{\exp(\mathbf{t}_i^F \cdot \mathbf{t}_i^I / \tau_1)}{\sum_{j=1}^N \exp(\mathbf{t}_i^F \cdot \mathbf{t}_j^I / \tau_1)} + \log \frac{\exp(\mathbf{t}_i^I \cdot \mathbf{t}_i^F / \tau_1)}{\sum_{j=1}^N \exp(\mathbf{t}_i^I \cdot \mathbf{t}_j^F / \tau_1)} \right), \quad (2)$$

where  $\tau_1$  is a scaling parameter to control the margin. The resulting total loss of the specialisation phase (III) is  $\mathcal{L}_{\text{III}} = \mathcal{L}_{\text{RSM}} + \lambda_{\text{MLM}} \mathcal{L}_{\text{MLM}}$ . An additional important component for regularising the RSM loss is the use of increased dropout (25%), including on attention. We set  $\tau_1 = 0.5$  and  $\lambda_{\text{MLM}} = 0.1$ , determined by a limited grid-search measuring  $\mathcal{L}_{\text{GA}}$  (Eq. (3)) of the joint model on a validation set. We also note that similar losses to the RSM loss, over the same or separate text segments, have been explored successfully for sentence representation learning [22, 49] in other settings. As such, we empirically observed that an objective as in [22] using masked FINDINGS to FINDINGS matching can achieve similar performance and may be an appropriate replacement in other biomedical settings with differing text structure.

**Text Augmentation.** As domain-specific datasets are often quite small, effective text augmentation can induce large benefits. In the radiology domain, the sentences of the FINDINGS and IMPRESSION sections, which contain the detailed description and summary of the radiological findings, are usually permutation-invariant on the sentence level (cf. [59]). We thus find that randomly shuffling sentences within each section is an effective text-augmentation strategy for both pretraining of CXR-BERT as well as during joint model training.

## 2.2 BioViL: Vision-Language Representation Learning

We now introduce BioViL, a simple but effective self-supervised VLP setup for the biomedical domain (Fig. 1), which we study in a chest X-ray (CXR) application setting. BioViL uses a convolutional neural



network (CNN) [37] image encoder  $E_{\text{img}}$ , our CXR-BERT text encoder  $E_{\text{txt}}$ , and projection models  $P_{\text{img}}$  and  $P_{\text{txt}}$  to learn representations in a joint space. The CNN model allows us to obtain a grid of local image embeddings  $\tilde{\mathbf{V}} = E_{\text{img}}(\mathbf{x}_{\text{img}})$ , which is fine-grained enough to be useful for segmentation (e.g.  $16 \times 16$ ). Each encoder is followed by a modality-specific two-layer perceptron projection model  $P$ , which projects the encoded modality to a joint space of 128 dimensions—e.g.  $\mathbf{V} = P_{\text{img}}(\tilde{\mathbf{V}})$ —where the representation is  $\ell_2$ -normalised. Note that projection should be applied to local embeddings before mean-pooling  $\mathbf{v} = \text{pool}(P_{\text{img}}(\tilde{\mathbf{V}}))$ , which gives us the global image embedding  $\mathbf{v}$ . The text branch uses the IMPRESSION section’s projected [CLS] token  $\mathbf{t}^I$  as the text representation in the joint space, as it contains a succinct summary of radiological findings. To align the representations and learn a joint embedding, we propose to use two loss terms. For a batch of size  $N$ , a symmetric contrastive loss [57] for *global alignment* of the image and text projections helps us learn the shared latent semantics:

$$\mathcal{L}_{\text{GA}} = -\frac{1}{N} \sum_{i=1}^N \left( \log \frac{\exp(\mathbf{v}_i \cdot \mathbf{t}_i^I / \tau_2)}{\sum_{j=1}^N \exp(\mathbf{v}_i \cdot \mathbf{t}_j^I / \tau_2)} + \log \frac{\exp(\mathbf{t}_i^I \cdot \mathbf{v}_i / \tau_2)}{\sum_{j=1}^N \exp(\mathbf{t}_i^I \cdot \mathbf{v}_j / \tau_2)} \right). \quad (3)$$

where  $\tau_2$  is a scaling parameter. Further, we maintain the  $\mathcal{L}_{\text{MLM}}$  loss (Eq. (1)) during joint training, resulting in the final joint loss  $\mathcal{L}_{\text{joint}} = \lambda_{\text{GA}} \mathcal{L}_{\text{GA}} + \mathcal{L}_{\text{MLM}}$ . We set  $\tau_2 = 0.5$  and  $\lambda_{\text{GA}} = 0.5$ , determined by a limited grid search measuring  $\mathcal{L}_{\text{GA}}$  on a validation set.

**Augmentations, Regularisation, and Image Encoder Pretraining.** Due to the small dataset sizes expected in biomedical applications, we use image and text augmentations to help learn known invariances. We use a ResNet-50 [28] architecture as our image encoder and pretrain the model on MIMIC-CXR images using SimCLR [6] with domain-specific augmentations as detailed in Section 4.1. For text, we use the same sentence-shuffling augmentation as in pretraining of CXR-BERT (see Section 4.1 for details). Furthermore, as in phase (III) of CXR-BERT training, we apply higher text encoder dropout (25%) than in standard BERT settings [16, 72]. We find that the combination of all these components, including continuous MLM optimisation, is important to improve downstream performance across the board (see ablation in Table 4).

**Zero-shot Classification.** After joint training, we use text prompts to cast the zero-shot classification problem into an image-text similarity task as in [30, 60, 61]. For  $C$  classes, subject-matter experts design  $C$  text prompts representing the target labels  $c \in \{1, \dots, C\}$ , e.g. for presence or absence of pneumonia (see Section 4.5). Each class prompt is represented as a vector of tokens  $\mathbf{w}^c$  and passed to the text encoder and projector of BioViL to obtain  $\ell_2$ -normalised text features  $\mathbf{t}^c = P_{\text{txt}}(E_{\text{txt}}(\mathbf{w}^c)) \in \mathbb{R}^{128}$ . For each input image  $\mathbf{x}_{\text{img}} \in \mathbb{R}^{H \times W}$ , we use the image encoder and projection module to obtain patch embeddings  $\mathbf{V} = P_{\text{img}}(E_{\text{img}}(\mathbf{x}_{\text{img}})) \in \mathbb{R}^{\frac{H}{16} \times \frac{W}{16} \times 128}$  for segmentation tasks or the pooled embedding  $\mathbf{v} = \text{pool}(\mathbf{V}) \in \mathbb{R}^{128}$  for instance-classification. We use dilated convolutions [81] to obtain higher-resolution feature maps. Probabilities for classes/regions can then be computed via a softmax over the cosine similarities between the image (or region) and prompt representations.

**Few-shot Tasks with BioViL.** To further assess the representation quality, linear probing is applied to local ( $\mathbf{V}$ ) and global ( $\mathbf{v}$ ) image representations, by learning  $\beta \in \mathbb{R}^{128 \times C}$  weights and a bias term. Unlike [30, 84], we leverage the pretrained projectors and class text embedding  $\mathbf{t}^c$  from the zero-shot setting by using them for initialisation, which leads to improved performance and further reduces the need for manual label collection. Specifically, in few-shot classification settings, the weights and bias are initialised with  $\beta = [\mathbf{t}^1, \dots, \mathbf{t}^C]$  and zeros, respectively.

### 3 Evaluating Self-Supervised Biomedical VLP

Accurate local alignment between modalities is an important characteristic of successful joint image-text training in healthcare, in particular since image and report samples often contain multiple clinical findings, each of which correspond to distinct image regions. Standard global-alignment approaches may attain high classification accuracy by overfitting to spurious image features for a given finding (e.g., chest tubes in images correlating with mentions of pneumothorax in reports). Image classification, the most frequently evaluated

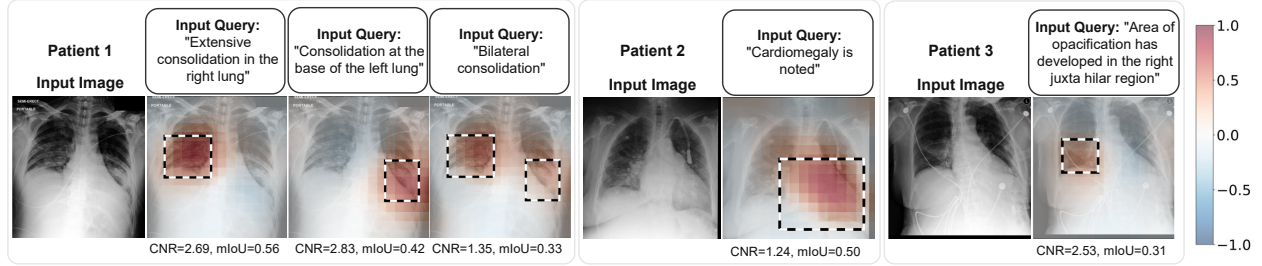


Figure 3: Examples from the newly released **MS-CXR** phrase grounding dataset with BioViL latent vector similarity for different input text queries superimposed as heatmaps. Dashed boxes are ground-truth annotations by radiologists. X-ray images are mirrored horizontally.

downstream task in related work [84, 30, 44, 55], requires only scene-level labels, hence a less sophisticated understanding of natural-language image descriptions. Image classification tasks can largely be solved by simply detecting a small set of words and maintaining some understanding of negation, as exemplified by the development of automated, rule-based text-labellers such as CheXpert [32]. Instance-level image-text retrieval tasks address some evaluation limitations, but do not require the level of language reasoning needed to solve local correspondence between phrases and image regions.

With this motivation in mind, we design a healthcare equivalent of general domain visual-grounding benchmarks, whilst being mindful of domain-specific radiology language (e.g. paraphrasing and negations). To name just a few challenges, a phrase grounding task requires the ability to parse domain specific location modifiers, the ability to deal with reporting style biases, and understanding of complex negations, all while relating the correct findings to specific image regions. To the best of our knowledge, existing public CXR benchmark datasets to evaluate local aspects of VLP have one or more of the following limitations (see Section 5 and Supp. C,D for more details): bounding boxes without corresponding free text descriptions, a limited number of samples, a limited number of abnormalities, and non-curated phrases impacting evaluation quality.

### 3.1 MS-CXR – A Chest X-ray Phrase Grounding Benchmark

We publicly release **MS-CXR**, a new dataset containing Chest X-ray bounding box labels paired with radiology text descriptions, annotated and verified by two board-certified radiologists (see some examples with BioViL outputs in Fig. 3, and more in Fig. C.1). With a large number of samples covering eight findings—balanced to ensure good coverage for all findings, and curated to ensure gold-standard evaluation of phrase grounding—our dataset is a valuable addition to the benchmarking landscape. The phrases in **MS-CXR** are not simple short captions, but genuine descriptions of radiological findings from original radiology reports [33] and dictated transcripts [36]. Thus, compared to existing evaluation datasets, the proposed benchmark is a more challenging real-world image-text reasoning task.

The **MS-CXR** dataset provides 1162 image-sentence pairs of bounding boxes and corresponding phrases, collected across eight different cardiopulmonary radiological findings, with an approximately equal number of pairs for each finding (see Table C.2). The dataset is released with instances chosen from the public MIMIC-CXR v2 [23, 33] image-text dataset. To obtain and verify bounding-box annotations, we first obtain MIMIC-CXR samples from a set of studies with preexisting region proposals, such as ellipses, based on data released in [70, 36]. To link each proposal region with candidate phrases, we sample sentences from the report of each study by extracting the highest matching sentences to the annotated labels using scores of the CheXbert sentence classifier [68], and also use transcriptions of dictations when available [36]. Next, to better balance findings, we sample additional MIMIC-CXR studies at random as well as MIMIC-CXR samples used in the ImaGenome dataset [78], the latter being a dataset of annotations of anatomical regions. These sampled studies do not have preexisting region proposals.

Radiologists then manually review separate sets of candidates. If a bounding box is not available, the radiologists manually annotate the corresponding region(s) in the image with new bounding boxes. Radiologists reject studies where no correct phrase candidates are available and where existing bounding boxes were placed incorrectly (e.g. covering too large an area). To ensure a high quality, consistent benchmark, the phrase-image samples that do not adhere to our guidelines (see Supp. C.1) are filtered out, such as phrases containing multiple abnormalities in distinct lung regions.

Table 2: Comparing evaluations conducted in recent CXR image-text alignment studies.

| Downstream task                            | Used in ref.* | Image encoder | Text encoder | Phrase reasoning | Findings localisation | Latent alignment | Annotation availability |
|--|---------------|---------------|--------------|------------------|-----------------------|------------------|-------------------------|
| Natural language inference                 | [B]           | -             | ✓            | ✓                | -                     | -                | Scarce                  |
| Phrase grounding                           | [B]           | ✓             | ✓            | ✓                | ✓                     | ✓                | Scarce                  |
| Image classification                       | [B,C,G,L,M]   | ✓             | -            | -                | -                     | -                | High                    |
| Zero-shot image classif.                   | [B,G]         | ✓             | ✓            | -                | -                     | ✓                | Moderate                |
| Dense image prediction (e.g. segmentation) | [B,G,L]       | ✓             | -            | -                | ✓                     | -                | High                    |
| Global image-text retrieval                | [C,G]         | ✓             | ✓            | -                | -                     | ✓                | High                    |

\*B, BioViL (Proposed); C, ConVIRT [84]; G, GLoRIA [30]; L, LoVT [55]; M, Local MI [44].

## 4 Experiments

We conduct a comprehensive evaluation of our CXR-BERT language model as well as the proposed BioViL self-supervised VLP approach, and compare both to state-of-the-art counterparts. Table 2 shows how our evaluation coverage compares to recent related studies. We begin by demonstrating CXR-BERT’s superior performance and improved vocabulary, including on a radiology-specific NLI benchmark. Next, we assess joint image-and-text understanding of BioViL on our new **MS-CXR** benchmark, which evaluates grounding of phrases describing radiological findings to the corresponding image regions. We also investigate zero-shot classification and fine-tuning performance of BioViL on image- and pixel-level prediction tasks via the RSNA pneumonia dataset [65, 75].

### 4.1 Setup

**Datasets.** We conduct experiments on the MIMIC-CXR v2 [33, 23] chest radiograph dataset, which provides 227,835 imaging studies for 65,379 patients, all collected in routine clinical practice. Each study contains a radiology report and one or more images (377,110 images in total). We only use frontal view scans (AP and PA) and also discard samples without an IMPRESSION section. From this data, we establish a training set of 146.7k samples and a set of 22.2k validation samples, ensuring that all samples used for the different downstream evaluations are kept in a held-out test set. We emphasise that no labels are used during pretraining; for early stopping only a loss on validation data is tracked. For evaluation, we use RadNLI [53] to assess the proposed CXR-BERT text model in isolation, the new **MS-CXR** assesses joint image-text understanding via phrase grounding, and the RSNA Pneumonia dataset [65, 75] to evaluate zero-shot segmentation, as well as zero-shot and fine-tuned classification performance.

**Image and Text Pre-processing.** We downsize and centre crop images to a resolution of  $512 \times 512$  whilst preserving image aspect ratios. We perform image augmentations during training including: random affine transformations, random colour jitter, and horizontal flips (only for image fine-tuning tasks). For text model pre-training we utilise the ‘FINDINGS’ and ‘IMPRESSION’ sections of reports, while joint training is performed using only the latter. During training, we perform sentence shuffling within sections as text-augmentation. Additionally, we perform limited automatic typo correction as in [5].

**Comparison Approaches.** We compare the proposed CXR-BERT text model to the other specialised PubMedBERT [25] and ClinicalBERT [2] models. Note that ClinicalBERT was used in most related studies [30, 44, 84, 55]. We compare BioViL to the closely related, state-of-the-art ConVIRT [84], LoVT [55] and GLoRIA [30] approaches (see Section 5 for more details).

**Metrics.** We report segmentation results via mean intersection over union (mIoU) and contrast-to-noise ratio (CNR), and report the Dice score to compare to [55]. We first compute the cosine similarity between a projected phrase embedding  $\mathbf{t}$  and each element of the local image representation  $\mathbf{V}$ , resulting in a grid of scores between  $[-1, 1]$ . For a given similarity threshold, we compute  $\text{IoU} = |A \cap B| / |A \cup B|$  with  $A$  being the true bounding box and  $B$  the thresholded region. The mIoU is then defined as an average over the thresholds  $[0.1, 0.2, 0.3, 0.4, 0.5]$ . The CNR measures the discrepancy between scores inside and out of the

Table 3: Evaluation of text encoder intrinsic properties and fine-tuning for radiology natural language inference: (1) RadNLI fine-tuning scores (average of 5 runs); (2) Mask prediction accuracy on MIMIC-CXR val. set; (3) Vocabulary comparison, number of tokens vs. original number of words in FINDINGS, increase shown as percentage.

|   | RadNLI accuracy<br>(MedNLI transfer) | Mask prediction<br>accuracy | Avg. # of tokens<br>after tokenization | Vocabulary<br>size |
|---|--------------------------------------|-----------------------------|--|--------------------|
| RadNLI baseline [53]                        | 53.30                                | -                           | -                                      | -                  |
| ClinicalBERT                                | 47.67                                | 39.84                       | 78.98 (+38.15%)                        | 28,996             |
| PubMedBERT                                  | 57.71                                | 35.24                       | 63.55 (+11.16%)                        | 28,895             |
| CXR-BERT (after Phase-III)                  | 60.46                                | 77.72                       | 58.07 (+1.59%)                         | 30,522             |
| CXR-BERT (after Phase-III + Joint Training) | 65.21                                | 81.58                       | 58.07 (+1.59%)                         | 30,522             |

bounding box region, without requiring hard thresholds. This evaluation of local similarities is important as some clinical downstream applications may benefit from heatmap visualisations as opposed to discrete segmentations. For CNR, let  $A$  and  $\bar{A}$  denote the interior and exterior of the bounding box, respectively. We then compute  $CNR = |\mu_A - \mu_{\bar{A}}| / (\sigma_A^2 + \sigma_{\bar{A}}^2)^{\frac{1}{2}}$ , where  $\mu_X$  and  $\sigma_X^2$  are the mean and variance of the similarity values in region  $X$ . Finally, the Dice score, defined as  $2|A \cap B| / (|A| + |B|)$ , is computed at one fixed threshold.

## 4.2 Text Model Evaluation

**Natural Language Understanding.** We use the RadNLI benchmark [53] to evaluate how well the proposed CXR-BERT text model captures domain-specific semantics. The dataset contains labelled hypothesis and premise pairs, sourced from MIMIC-CXR radiology reports, with the following label categories: (1) entailment, i.e. the hypothesis can be inferred from the premise; (2) contradiction, i.e. the hypothesis cannot be inferred from the premise; and (3) neutral, i.e. the inference relation is undetermined. RadNLI provides expert-annotated development and test sets (480 examples each), but no official training set. Thus, following [53], we use MedNLI [66] for training, which has 11k samples sourced from MIMIC-III discharge summaries, with equally distributed NLI labels. We fine-tune the language models up to 20 epochs and use early stopping by monitoring accuracy scores on the RadNLI development set. Table 3 summarises the NLI evaluation, masked token prediction, and subword tokenisation results. Using only MedNLI training samples, our model achieves a good accuracy of 65.21%, and far outperforms fine-tuned ClinicalBERT, PubMedBERT, and the score reported in RadNLI [53]. Another important result is that RadNLI accuracy improves after joint training with images (last row of Table 3).

**Mask Prediction Accuracy.** While mask prediction accuracy does not always translate to downstream application performance, it is an auxiliary metric that captures important aspects of a language model’s grasp of a target domain. We report Top-1 mask prediction accuracy on radiology reports in the MIMIC-CXR validation set (Table 3), and follow the standard masking configuration (15% masking probability). Despite being trained on closely related data, our CXR-BERT displays a much better mask prediction accuracy compared to ClinicalBERT (trained on MIMIC-III, which includes radiology reports) and PubMedBERT (trained on biomedical literature text). This suggests that radiology text significantly differs from other clinical text or biomedical literature text, highlighting the need for specialised text encoder models.

Table 4: CXR-BERT ablation. CNR and mIoU are macro averages of BioViL performance on all categories of **MS-CXR**. *Syn. sim.* denotes the average cosine similarity between RadNLI entailments. *Cont. gap* is the average similarity gap of RadNLI entailment and contradiction pairs. CXR-BERT is the combination of all components below the first row.

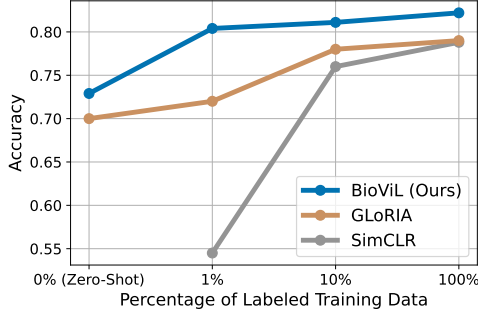
| Model or pretraining stage        | RadNLI    |           | Grounding |       |
|-----------------------------------|-----------|-----------|-----------|-------|
|                                   | Syn. sim. | Cont. gap | mIoU      | CNR   |
| ClinicalBERT                      | .657      | .609      | .182      | 0.791 |
| Pretrain & Vocab (I-II)           | .749      | .646      | .194      | 0.796 |
| + MLM grounding during joint      | .871      | .745      | .209      | 0.860 |
| + Use of attention drop-out (III) | .893      | .802      | .217      | 0.945 |
| + RSM Pretrain (III)              | .877      | .779      | .220      | 1.012 |
| + Sentence shuffling (CXR-BERT)   | .884      | .798      | .220      | 1.031 |

**Ablation.** We also conduct an ablation of the various aspects of CXR-BERT, measuring the impact after joint training. Table 4 shows that all components of CXR-BERT contribute to improved downstream and NLI performance, both in terms of alignment between related sentences (entailments) and of discrimination of contradictions. In particular, note the substantial improvement on these scores due to keeping the MLM objective during joint finetuning.

Table 5: Contrast-to-noise ratio (CNR) obtained on the MS-CXR dataset, averaged over four runs with different seeds. The results are collected using different text encoder and training objectives (e.g., G&L: Global and local loss).

| Method       | Objective | Text encoder | Atelectasis | Cardiomegaly | Consolidation | Lung opacity | Edema | Pneumonia | Pneumothorax | Pl. effusion | Avg.  |
|--------------|-----------|--------------|-------------|--------------|---------------|--------------|-------|-----------|--------------|--------------|-------|
| Baseline     | Global    | ClinicalBERT | 0.70        | 0.53         | 1.15          | 0.75         | 0.83  | 0.85      | 0.29         | 1.05         | 0.769 |
| Baseline     | Global    | PubMedBERT   | 0.72        | 0.64         | 1.22          | 0.69         | 0.80  | 0.91      | 0.21         | 0.99         | 0.773 |
| ConVIRT [84] | Global    | ClinicalBERT | 0.86        | 0.64         | 1.25          | 0.78         | 0.68  | 1.03      | 0.28         | 1.02         | 0.818 |
| GLORIA [30]  | G&L       | ClinicalBERT | 0.98        | 0.53         | 1.38          | 1.05         | 0.66  | 1.18      | 0.47         | 1.20         | 0.930 |
| BioViL       | Global    | CXR-BERT     | 1.02        | 0.63         | 1.42          | 1.05         | 0.93  | 1.27      | 0.48         | 1.40         | 1.027 |
| BioViL-L     | G&L       | CXR-BERT     | 1.17        | 0.95         | 1.45          | 1.19         | 0.96  | 1.19      | 0.74         | 1.50         | 1.142 |

Table 6: RSNA Pneumonia zero-shot and fine-tuned classification. We compare to GLORIA scores reported in [30] which outperforms ConVIRT [84] (see [30]). Training size: GLORIA ( $N = 186k$ , private dataset), BioViL ( $N = 146.7k$  of MIMIC-CXR).



| Method      | Type       | Text model   | Loss           | % of labels | Acc.  | F1    | AUROC |
|-------------|------------|--------------|----------------|-------------|-------|-------|-------|
| SimCLR [6]  | Image only | -            | Global         | 1%          | 0.545 | 0.522 | 0.701 |
|             |            |              |                | 10%         | 0.760 | 0.639 | 0.802 |
|             |            |              |                | 100%        | 0.788 | 0.675 | 0.849 |
| GLORIA [30] | Joint      | ClinicalBERT | Global & local | Zero-shot   | 0.70  | 0.58  | -     |
|             |            |              |                | 1%          | 0.72  | 0.63  | 0.861 |
|             |            |              |                | 10%         | 0.78  | 0.63  | 0.880 |
|             |            |              |                | 100%        | 0.79  | 0.65  | 0.886 |
| Baseline    | Joint      | ClinicalBERT | Global         | Zero-shot   | 0.719 | 0.614 | 0.812 |
| BioViL      | Joint      | CXR-BERT     | Global         | Zero-shot   | 0.732 | 0.665 | 0.831 |
|             |            |              |                | 1%          | 0.805 | 0.723 | 0.881 |
|             |            |              |                | 10%         | 0.812 | 0.727 | 0.884 |
|             |            |              |                | 100%        | 0.822 | 0.733 | 0.891 |

### 4.3 Local Alignment Evaluation – Phrase Grounding

We perform a phrase grounding evaluation of the pretrained BioViL model on the MS-CXR dataset. For each image–phrase pair, the image is passed to the CNN image encoder and projected to obtain a grid of image representations  $\mathbf{V}$  in the joint space. Similarly, the phrase is embedded via the text encoder and projected to the joint space to obtain  $\mathbf{t}$ . Cosine similarity between  $\mathbf{t}$  and elements of  $\mathbf{V}$  produces a similarity grid, which is evaluated against the ground-truth bounding boxes. Table 5 shows the superior phrase grounding results achieved by BioViL across radiological findings. We also create BioViL-L by adding a local loss term as in [30], which further improves phrase grounding performance for almost all findings. Moreover, the ablation in Table 4 demonstrates that there are clear gains to be had in visual grounding performance by improving the text model.

### 4.4 Global Alignment Evaluation – Zero-shot and Fine-tuned Classification

To measure the quality of the global alignment, the joint models are also benchmarked on zero-/few-shot binary pneumonia classification problems (image-level) using the external RSNA dataset [65]. Fine-tuning is done via linear probing, i.e. only a last linear layer is trained. The evaluation is conducted on  $\mathcal{D}_{\text{test}} = 9006$  images as in [30] (30% eval. / 70% train.) using the dataset’s ground-truth labels. We define two simple text prompts for BioViL, representing presence/absence of pneumonia, namely “Findings suggesting pneumonia” and “No evidence of pneumonia”. The image encoders are utilised and fine-tuned as described in Section 2.2.

The zero-shot and fine-tuned results in Table 6 show that our focus on better text modelling results in improved joint modelling of shared latent information between text-image pairs. Note that, to achieve its superior performance here and in Section 4.5, BioViL does not require extensive human expert text-prompt engineering as for example conducted in GLORIA [30], where variations over severity and/or location were created (see Supp. A.1 for a text-prompt sensitivity analysis on BioViL).

### 4.5 Local Alignment Evaluation – Semantic Segmentation

We evaluate models on an RSNA pneumonia segmentation task, using grid-level image representations in the joint latent space. We use the same text prompts as in the previous section for all models, and evaluate against ground-truth bounding boxes of the RSNA pneumonia dataset ( $|\mathcal{D}_{\text{train}}| = 6634$  and  $|\mathcal{D}_{\text{test}}| = 2907$ ).



Table 7 shows that BioViL significantly reduces the need for dense annotations as compared to similar multi-modal and image-only pretraining approaches, outperforming them when using the same number of labelled data points. Note that our proposed modelling framework BioViL (Fig. 1), uses neither a local loss term [30, 55], nor a separate object detection [62] or segmentation network [64]. Further, while Table 7 shows results using two simple queries, we find that BioViL continues to outperform related work even when more prompts are used for all models as in [30]. Dice and IoU are computed using the same threshold value (0.6) on predictions scaled between [0, 1].

Table 7: RSNA pneumonia segmentation. Related work is reproduced in the same experimental setup except for LoVT [55]. *Zero-shot* and *linear probing* results demonstrate the effectiveness of learning and pretraining with free-text data.

| Method       | % of Labels | Supervision | IoU   | Dice  | CNR   |
|--------------|-------------|-------------|-------|-------|-------|
| LoVT [55]    | 100%        | Lin. prob.  | -     | 0.518 | -     |
| ConVIRT [84] | -           | Zero-shot   | 0.228 | 0.348 | 0.849 |
| GLoRIA [30]  | -           | Zero-shot   | 0.245 | 0.366 | 1.052 |
| BioViL       | -           | Zero-shot   | 0.355 | 0.496 | 1.477 |
| SimCLR [6]   | 5%          | Lin. prob.  | 0.382 | 0.525 | 1.722 |
| SimCLR [6]   | 100%        | Lin. prob.  | 0.427 | 0.570 | 1.922 |
| BioViL       | 5%          | Lin. prob.  | 0.446 | 0.592 | 2.077 |
| BioViL       | 100%        | Lin. prob.  | 0.469 | 0.614 | 2.178 |

## 5 Related Work

We refer the reader to Supp. D for a more detailed review of related work.

**Biomedical Vision-Language Processing.** Multiple studies explore joint representation learning for paired image and text radiology data [29, 30, 44, 55, 84]. [84] follow a contrastive learning formulation for instance-level representation learning, while [30, 55] introduce approaches that combine instance-level image-report learning with local terms for radiology data. An alternative, local-only objective is explored by [44], approximating the mutual information between local image features and sentence-level text features. While most related approaches use no ground truth, [5] study a semi-supervised edema severity classification setting, and [27] assume sets of seen and unseen labels towards CXR zero-shot classification.

Related medical VLP work commonly uses publicly available contextual word embedding models including BioBERT [38], ClinicalBERT [2], BioClinicalBERT [2], or PubMedBERT [25]. The models are either trained from scratch or fine-tuned via continual pretraining using an MLM objective. Additional objectives such as adversarial losses [46] are added infrequently. The specialised corpora these models use include PubMed abstracts and PubMed Central full texts (see [2]), as well as MIMIC-III [34] clinical notes.

**Local Alignment Datasets.** Presently, no datasets exist that allow for phrase grounding of radiology findings, but some enable different forms of local image evaluations. VinDr [56], RSNA Pneumonia [65], and the NIH Chest X-ray Dataset [75] provide bounding-box annotations, but lack free-text descriptions. REFLACX [36] provides gaze locations (ellipses) captured with an eye tracker, dictated reports, and some ground truth annotations for gaze locations, but no full phrase matches to image regions. Phrase annotations for MIMIC-CXR data released in [70] are of small size (350 studies), only contain two abnormalities, and for some samples have shortened phrases that were adapted to simplify the task. The ground-truth set of ImaGenome [78] only contains 500 studies, bounding-box regions annotate anatomical regions rather than radiological findings, and its sentence annotations are not curated for grounding evaluation.

## 6 Conclusion

In this article, we show that careful attention to text modelling can lead to large benefits for all learned models and representations in self-supervised vision language processing frameworks for medical image-text applications. We introduce a novel pretraining procedure and publicly release a chest X-ray (CXR) domain-specific language model: CXR-BERT. It has an improved vocabulary, increased masked token prediction performance on CXR data, achieves superior performance on a radiology natural language inference benchmark, and contributes to improved downstream performance for all aspects of CXR VLP approaches.

We also present BioViL, as a simple yet effective baseline for self-supervised multi-modal learning for paired image-text radiology data, with a focus on improved text modelling. The approach displays state-of-the-art performance on a large number of downstream tasks evaluating global and local aspects of the image



model, text model, and joint latent space. On zero-shot tasks, the model does not require extensive text-prompt engineering compared to prior work. Notably, it outperforms related work on segmentation despite not using a local loss term, and without requiring an additional vision model to produce region proposals. We do not advocate against local losses. In fact, adding a local loss term improves phrase grounding (Table 5). But our study highlights that careful text modelling enables even global alignment to learn local aspects, providing a strong baseline to compare against.

To support the research community in evaluating fine-grained image-text understanding in the radiology domain, we also publicly release a chest X-ray phrase grounding dataset called **MS-CXR**. It presents a more challenging benchmark for joint image-text understanding compared to existing datasets, requiring reasoning over real-world radiology language to ground findings in the correct image locations.

Limitations of the proposed joint approach include that it does not explicitly deal with false negatives in the contrastive losses. Furthermore, co-occurrence of multiple abnormalities could enable contrastive methods to focus only on a subset to match pairs, e.g. pneumothorax and chest tubes commonly occur together [24]. Amongst its failure cases (see Supp. A.2 for more), we have seen that the approach struggles with very small structures, likely due to image resolution limits. Future work will explore the presented ideas in other domains, expand the evaluated radiological findings, and explore using larger image resolution.

## References

- [1] Hassan Akbari, Svebor Karaman, Surabhi Bhargava, Brian Chen, Carl Vondrick, and Shih-Fu Chang. Multi-level multimodal common semantic space for image-phrase grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 12476–12486. Computer Vision Foundation / IEEE, 2019.
- [2] Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA, 2019. Association for Computational Linguistics.
- [3] Aurelia Bustos, Antonio Pertusa, Jose-Maria Salinas, and Maria de la Iglesia-Vayá. PadChest: A large chest X-ray image dataset with multi-label annotated reports. *Medical image analysis*, 66:101797, 2020.
- [4] Arlene Casey, Emma Davidson, Michael Poon, Hang Dong, Daniel Duma, Andreas Grivas, Claire Grover, Víctor Suárez-Paniagua, Richard Tobin, William Whiteley, et al. A systematic review of natural language processing applied to radiology reports. *BMC medical informatics and decision making*, 21(1):1–18, 2021.
- [5] Geeticka Chauhan, Ruizhi Liao, William Wells, Jacob Andreas, Xin Wang, Seth Berkowitz, Steven Horng, Peter Szolovits, and Polina Golland. Joint modeling of chest radiographs and radiology reports for pulmonary edema assessment. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 529–539. Springer, 2020.
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020.
- [7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020.
- [8] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, November 2020.
- [9] Sasank Chilamkurthy, Rohit Ghosh, Swetha Tanamala, Mustafa Biviji, Norbert G Campeau, Vasantha Kumar Venugopal, Vidur Mahajan, Pooja Rao, and Prashant Warier. Deep learning algorithms for

- detection of critical findings in head CT scans: a retrospective study. *The Lancet*, 392(10162):2388–2396, 2018.
- [10] Songtai Dai, Quan Wang, Yajuan Lyu, and Yong Zhu. BDKG at MEDIQA 2021: System report for the radiology report summarization task. In *Proceedings of the 20th Workshop on Biomedical Language Processing*, pages 103–111. Association for Computational Linguistics, 2021.
  - [11] Samyak Datta, Karan Sikka, Anirban Roy, Karuna Ahuja, Devi Parikh, and Ajay Divakaran. Align2Ground: Weakly supervised phrase grounding guided by image-caption alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 2601–2610. IEEE, 2019.
  - [12] Surabhi Datta and Kirk Roberts. A hybrid deep learning approach for spatial trigger extraction from radiology reports. In *Proceedings of the Third International Workshop on Spatial Language Understanding*, volume 2020, pages 50–55. Association for Computational Linguistics, 2020.
  - [13] Surabhi Datta, Yuqi Si, Laritza Rodriguez, Sonya E Shooshan, Dina Demner-Fushman, and Kirk Roberts. Understanding spatial language in radiology: Representation framework, annotation, and spatial relation extraction from chest X-ray reports using deep learning. *Journal of biomedical informatics*, 108:103473, 2020.
  - [14] Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 2016.
  - [15] Karan Desai and Justin Johnson. VirTex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11162–11173, 2021.
  - [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, 2019. Association for Computational Linguistics.
  - [17] Dmitriy Dligach, Steven Bethard, Lee Becker, Timothy Miller, and Guergana K Savova. Discovering body site and severity modifiers in clinical texts. *Journal of the American Medical Informatics Association*, 21(3):448–454, 2014.
  - [18] Jared A. Dunnmon, Alexander J. Ratner, Khaled Saab, Nishith Khandwala, Matthew Markert, Hersh Sagreiya, Roger Goldman, Christopher Lee-Messer, Matthew P. Lungren, Daniel L. Rubin, and Christopher Re. Cross-modal data programming enables rapid medical machine learning. *Patterns*, 1(2):100019, 2020.
  - [19] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017.
  - [20] Sabri Eyuboglu, Geoffrey Angus, Bhavik N Patel, Anuj Pareek, Guido Davidzon, Jin Long, Jared Dunnmon, and Matthew P Lungren. Multi-task weak supervision enables anatomically-resolved abnormality detection in whole-body FDG-PET/CT. *Nature communications*, 12(1):1–15, 2021.
  - [21] Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C Platt, et al. From captions to visual concepts and back. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1473–1482. IEEE Computer Society, 2015.

- [22] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, 2021.
- [23] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23):e215–e220, 2000.
- [24] Benedikt Graf, Arkadiusz Sitek, Amin Katouzian, Yen-Fu Lu, Arun Krishnan, Justin Rafael, Kirstin Small, and Yiting Xie. Pneumothorax and chest tube classification on chest X-rays for detection of missed pneumothorax. *Machine Learning for Health (ML4H) NeurIPS Workshop: Extended Abstract*, 2020.
- [25] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.
- [26] Tanmay Gupta, Arash Vahdat, Gal Chechik, Xiaodong Yang, Jan Kautz, and Derek Hoiem. Contrastive learning for weakly supervised phrase grounding. In *16th European Conference on Computer Vision, ECCV 2020*, pages 752–768. Springer, 2020.
- [27] Nasir Hayat, Hazem Lashen, and Farah E Shamout. Multi-label generalized zero shot learning for the classification of disease in chest radiographs. In *Machine Learning for Healthcare Conference*, pages 461–477. PMLR, 2021.
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778. IEEE Computer Society, 2016.
- [29] Tzu-Ming Harry Hsu, Wei-Hung Weng, Willie Boag, Matthew McDermott, and Peter Szolovits. Unsupervised multimodal representation learning across medical images and reports. *Machine Learning for Health (ML4H) NeurIPS Workshop*, 2018.
- [30] Shih-Cheng Huang, Liyue Shen, Matthew P Lungren, and Serena Yeung. GLoRIA: A multimodal global-local representation learning framework for label-efficient medical image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3942–3951, 2021.
- [31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [32] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn Ball, Katie Shpanskaya, et al. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Thirty-Third AAAI Conference on Artificial Intelligence*, pages 590–597. AAAI Press, 2019.
- [33] A Johnson, T Pollard, SJ Berkowitz, R Mark, and S Horng. MIMIC-CXR database (version 2.0.0). PhysioNet, 2019.
- [34] 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. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [35] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *European Conference on Computer Vision*, pages 67–84. Springer, 2016.

- [36] Ricardo Bigolin Lanfredi, Mingyuan Zhang, William F Auffermann, Jessica Chan, Phuong-Anh T Duong, Vivek Srikumar, Trafton Drew, Joyce D Schroeder, and Tolga Tasdizen. REFLACX, a dataset of reports and eye-tracking data for localization of abnormalities in chest x-rays. *arXiv preprint arXiv:2109.14187*, 2021.
- [37] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [38] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.
- [39] Ang Li, Allan Jabri, Armand Joulin, and Laurens Van Der Maaten. Learning visual n-grams from web data. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 4183–4192. IEEE Computer Society, 2017.
- [40] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-VL: A universal encoder for vision and language by cross-modal pre-training. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, volume 34(7), pages 11336–11344. AAAI Press, 2020.
- [41] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. VisualBERT: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [42] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. *arXiv preprint arXiv:2110.05208*, 2021.
- [43] Yikuan Li, Hanyin Wang, and Yuan Luo. A comparison of pre-trained vision-and-language models for multimodal representation learning across medical images and reports. In *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1999–2004. IEEE, 2020.
- [44] Ruizhi Liao, Daniel Moyer, Miriam Cha, Keegan Quigley, Seth Berkowitz, Steven Horng, Polina Golland, and William M Wells. Multimodal representation learning via maximization of local mutual information. *International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI)*, 2021.
- [45] Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew McDermott, Willie Boag, Wei-Hung Weng, Peter Szolovits, and Marzyeh Ghassemi. Clinically accurate chest X-ray report generation. In *Machine Learning for Healthcare Conference*, pages 249–269. PMLR, 2019.
- [46] Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. Adversarial training for large neural language models. *arXiv preprint arXiv:2004.08994*, 2020.
- [47] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [48] Yongfei Liu, Bo Wan, Lin Ma, and Xuming He. Relation-aware instance refinement for weakly supervised visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5612–5621, 2021.
- [49] Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence representations. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- [50] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.

- [51] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13–23, 2019.
- [52] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 11–20. IEEE Computer Society, 2016.
- [53] Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. 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. Association for Computational Linguistics, 2021.
- [54] Zongshen Mu, Siliang Tang, Jie Tan, Qiang Yu, and Yueting Zhuang. Disentangled motif-aware graph learning for phrase grounding. *AAAI*, 2021.
- [55] Philip Müller, Georgios Kaissis, Congyu Zou, and Daniel Rückert. Joint learning of localized representations from medical images and reports. *arXiv preprint arXiv:2112.02889*, 2021.
- [56] Ha Q Nguyen, Khanh Lam, Linh T Le, Hieu H Pham, Dat Q Tran, Dung B Nguyen, Dung D Le, Chi M Pham, Hang TT Tong, Diep H Dinh, et al. VinDr-CXR: An open dataset of chest X-rays with radiologist’s annotations. *arXiv preprint arXiv:2012.15029*, 2020.
- [57] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [58] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015.
- [59] Konpat Preechakul, Chawan Piansaddhayanon, Burin Naowarat, Tirasan Khandhawit, Sira Sriswasdi, and Ekapol Chuangsuwanich. Set prediction in the latent space. *Advances in Neural Information Processing Systems*, 34, 2021.
- [60] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [61] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. DenseCLIP: Language-guided dense prediction with context-aware prompting. *arXiv preprint arXiv:2112.01518*, 2021.
- [62] Joseph Redmon and Ali Farhadi. YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [63] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, 28:91–99, 2015.
- [64] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [65] George Shih, Carol C Wu, Safwan S Halabi, Marc D Kohli, Luciano M Prevedello, Tessa S Cook, Arjun Sharma, Judith K Amorosa, Veronica Arteaga, Maya Galperin-Aizenberg, et al. Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. *Radiology: Artificial Intelligence*, 1(1):e180041, 2019.

- [66] Chaitanya Shivade. MedNLI - A natural language inference dataset for the clinical domain. PhysioNet, October 2019.
- [67] PY Simard, D Steinkraus, and JC Platt. Best practices for convolutional neural networks applied to visual document analysis. In *Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings.*, pages 958–963. IEEE, 2003.
- [68] Akshay Smit, Saahil Jain, Pranav Rajpurkar, Anuj Pareek, Andrew Y Ng, and Matthew Lungren. Combining automatic labelers and expert annotations for accurate radiology report labeling using BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1500–1519. Association for Computational Linguistics, 2020.
- [69] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VL-BERT: pre-training of generic visual-linguistic representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2019.
- [70] L.K. Tam, X. Wang, E. Turkbey, K. Lu, Y. Wen, and D. Xu. Weakly supervised one-stage vision and language disease detection using large scale pneumonia and pneumothorax studies. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2020*, March 2020.
- [71] Joseph J Titano, Marcus Badgeley, Javin Schefflein, Margaret Pain, Andres Su, Michael Cai, Nathaniel Swinburne, John Zech, Jun Kim, Joshua Bederson, et al. Automated deep-neural-network surveillance of cranial images for acute neurologic events. *Nature medicine*, 24(9):1337–1341, 2018.
- [72] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008, 2017.
- [73] A Wallis and P McCoubrie. The radiology report—are we getting the message across? *Clinical radiology*, 66(11):1015–1022, 2011.
- [74] Qinxin Wang, Hao Tan, Sheng Shen, Michael Mahoney, and Zhewei Yao. MAF: Multimodal alignment framework for weakly-supervised phrase grounding. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2030–2038, Online, 2020. Association for Computational Linguistics.
- [75] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. ChestX-Ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2097–2106. IEEE Computer Society, 2017.
- [76] John R Wilcox. The written radiology report. *Applied Radiology*, 35(7):33, 2006.
- [77] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [78] Joy T Wu, Nkechinyere Nneka Agu, Ismini Lourentzou, Arjun Sharma, Joseph Alexander Paguio, Jasper Seth Yao, Edward Christopher Dee, William G Mitchell, Satyananda Kashyap, Andrea Giovannini, et al. Chest imagenome dataset for clinical reasoning. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [79] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [80] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.



- [81] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [82] Tianyu Yu, Tianrui Hui, Zhihao Yu, Yue Liao, Sansi Yu, Faxi Zhang, and Si Liu. Cross-modal omni interaction modeling for phrase grounding. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*, pages 1725–1734, 2020.
- [83] Yuhao Zhang, Daisy Yi Ding, Tianpei Qian, Christopher D Manning, and Curtis P Langlotz. Learning to summarize radiology findings. In *Proceedings of the Ninth International Workshop on Health Text Mining and Information Analysis*, pages 204–213. Association for Computational Linguistics, 2018.
- [84] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020.
- [85] Zhu Zhang, Zhou Zhao, Zhijie Lin, Xiuqiang He, et al. Counterfactual contrastive learning for weakly-supervised vision-language grounding. *Advances in Neural Information Processing Systems*, 33:18123–18134, 2020.

## A Additional Experiments

### A.1 Zero-shot Text-prompt Sensitivity Analysis

Vision-language pretraining aligns image and text data in a joint representation space, which enables impressive zero-shot downstream image classification performance via input text prompts. However, some recent work [30, 84] has shown that downstream task performance can heavily depend on the choice of text prompts. Constructing good text prompts (prompt engineering) may require expert domain knowledge and can be costly and time-consuming. In Table A.1, we study RSNA pneumonia zero-shot classification performance using different text prompt combinations. Compared to the baseline, BioViL demonstrates much lower sensitivity to prompt choices selected from the data distribution. BioViL maintains its high performance even when faced with relatively long queries, which is not the case for the baseline model. These observations suggest that our improved text encoder CXR-BERT is more robust to prompt variations, and makes prompt engineering easier and less of a requirement to achieve high zero-shot classification performance.

Table A.1: Text prompt sensitivity analysis on the RSNA pneumonia zero-shot classification task. Image-text models trained without the proposed text modelling improvements (Table 4) show higher sensitivity to different input text prompts as the latent text embeddings are inconsistent for synonym phrases. For this reason, baseline methods often require post-hoc text prompt engineering heuristics (e.g. [30]).

| Method       | Pos. Query                                   | Neg. Query                                | F1 Score | ROC-AUC | $ \Delta AUC $ |
|--------------|--|---|----------|---------|----------------|
| BioViL       | “Findings suggesting pneumonia”              | “There is no evidence of acute pneumonia” | 0.657    | 0.822   | -              |
| ClinicalBert | “Findings suggesting pneumonia”              | “There is no evidence of acute pneumonia” | 0.581    | 0.731   | -              |
| BioViL       | “Findings suggesting pneumonia”              | “No evidence of pneumonia”                | 0.665    | 0.831   | -              |
| BioViL       | “Consistent with the diagnosis of pneumonia” | “There is no evidence of acute pneumonia” | 0.669    | 0.839   | 0.008          |
| ClinicalBert | “Findings suggesting pneumonia”              | “No evidence of pneumonia”                | 0.614    | 0.815   | -              |
| ClinicalBert | “Consistent with the diagnosis of pneumonia” | “There is no evidence of acute pneumonia” | 0.621    | 0.694   | 0.121          |
| BioViL       | “Findings consistent with pneumonia”         | “No evidence of pneumonia”                | 0.672    | 0.838   | -              |
| BioViL       | “Findings consistent with pneumonia”         | “There is no pneumonia”                   | 0.679    | 0.847   | 0.009          |
| ClinicalBert | “Findings consistent with pneumonia”         | “No evidence of pneumonia”                | 0.640    | 0.782   | -              |
| ClinicalBert | “Findings consistent with pneumonia”         | “There is no pneumonia”                   | 0.586    | 0.724   | 0.058          |

### A.2 Qualitative Results – Phrase Grounding

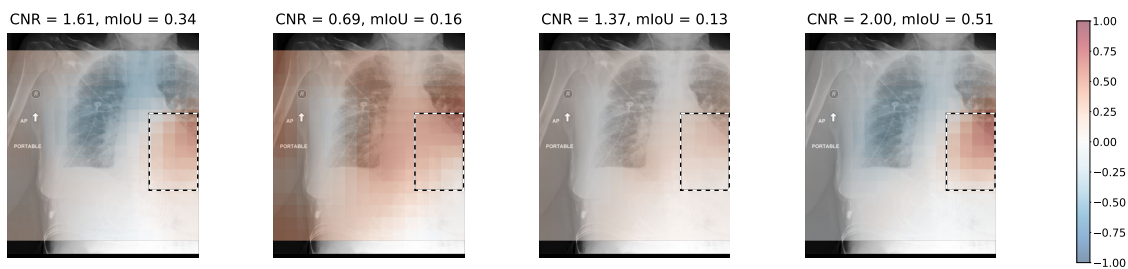
In Fig. A.1, we show and describe some phrase grounding examples obtained with different models on the MS-CXR dataset. From left to right, the figure shows the ClinicalBERT baseline, ConVIRT, GLoRIA, and BioViL similarity maps. While the figure only illustrates a few examples, the results demonstrate that phrase grounding performance can be significantly enhanced by leveraging improved text modelling (BioViL). The examples include clinical findings that differ in size, type, and anatomical location.

Additionally, in Fig. A.2, we show and describe some failure cases of BioViL on the MS-CXR dataset to motivate any further research on this topic. In particular, the models show limitations in grounding the descriptions relating to smaller structures (e.g., rib fracture, pneumothorax), and in a few cases the location modifier is not disassociated from the entities corresponding to abnormalities, see (a) in Fig. A.2.

### A.3 Additional Experimental Results

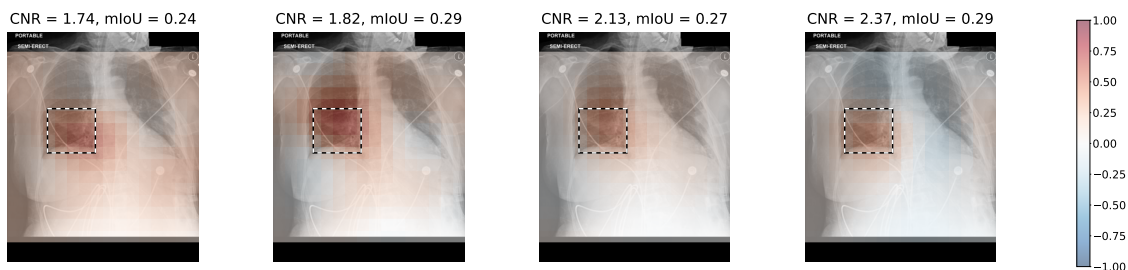
In Table A.2, an extension of Table 6 is provided to include the sensitivity and specificity metrics for the zero-shot and fine-tuned classification experiments presented in Section 4.4. The classification thresholds are set to maximise the F1 scores for each method. Further, in Table A.3 we provide mean IoU scores for the phrase grounding experiments presented in Section 4.3, which evaluates the pretrained BioViL model on the MS-CXR dataset. We observed that the distribution of similarity scores is different for GLoRIA and BioViL-L due to the different temperature parameter used in the local loss term in [30]. To provide a fair comparison, we adjust the similarity scores via min-max scaling to the full  $[-1, 1]$  range. The same scaling strategy is utilised in the implementation of the baseline method [30]. Note that the CNR scores are not affected by this linear re-scaling.

**Query:** "Retrocardiac opacification is again seen with air bronchograms suggesting pneumonia"



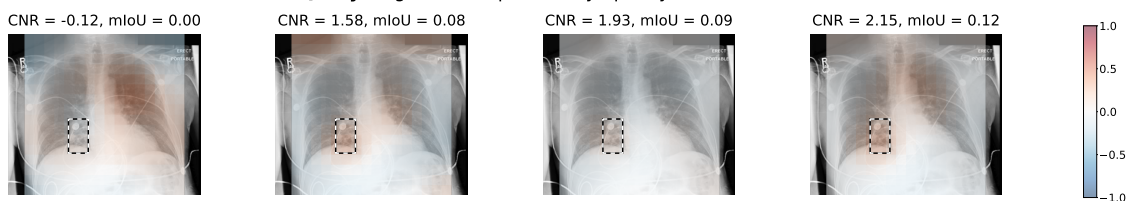
(a) Relatively long and complex query

**Query:** "area of opacification has developed in the right juxta hilar region"



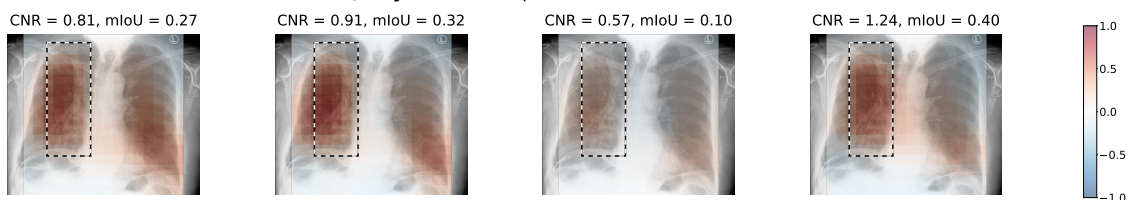
(b) Complex anatomical location specification

**Query:** "right basilar pulmonary opacity"



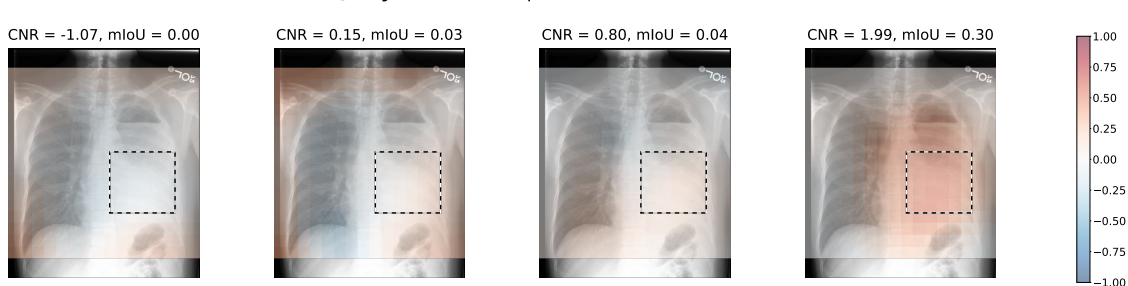
(c) Small ground-truth bounding box

**Query:** "multifocal pneumonia"



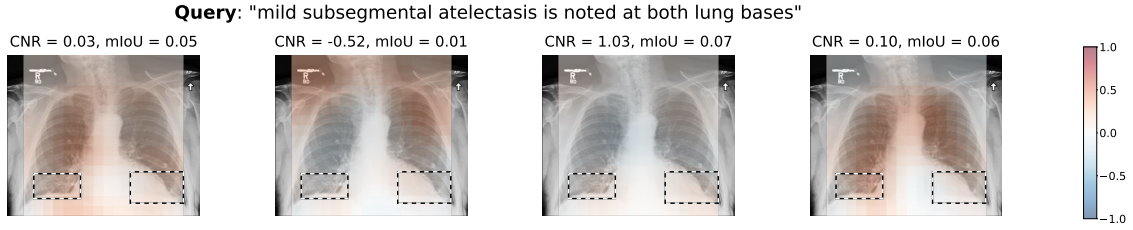
(d) Multifocal pneumonia example which is localized in the right lobe

**Query:** "left basilar opacification"

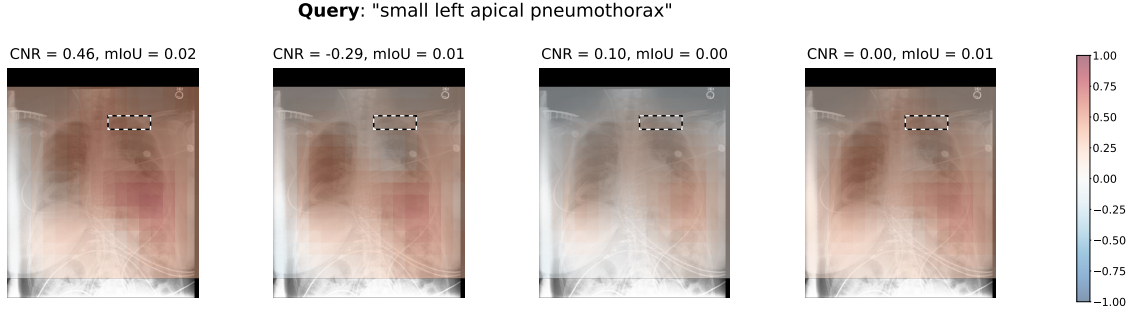


(e) Location modifier "left basilar"

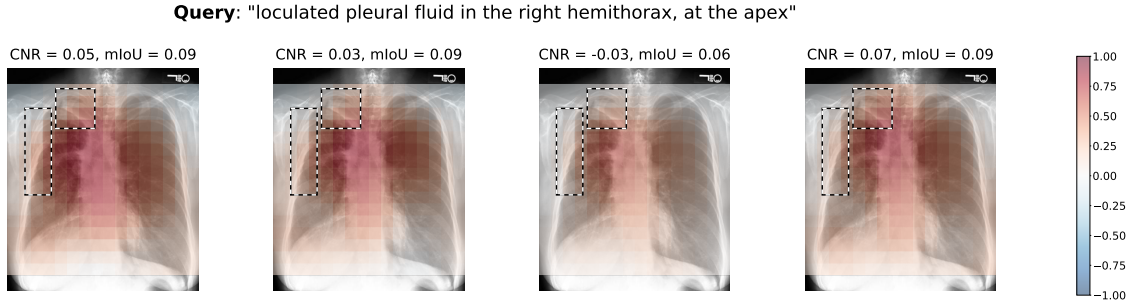
Figure A.1: Qualitative examples from MS-CXR phrase grounding benchmark. Model outputs (latent vector similarity) are compared (from left, ClinicalBERT baseline, ConVIRT, GLoRIA, and BioViL)



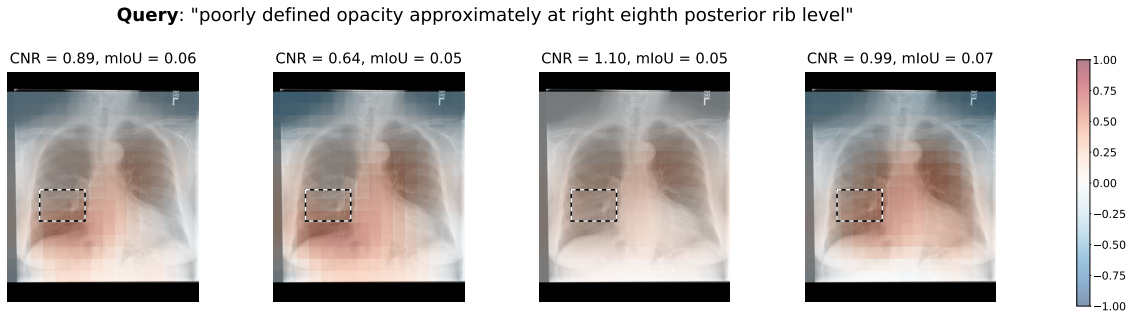
(a) Failed to recognise atelectasis despite having “both lung bases” location specification



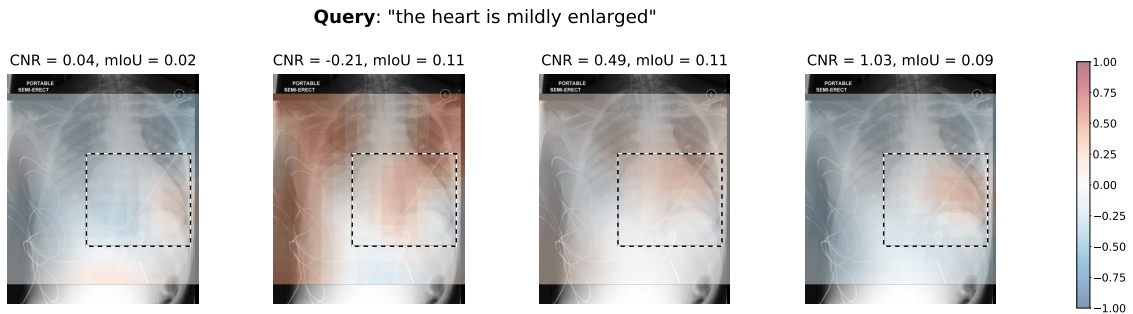
(b) Failed to recognise small pneumothorax despite having “apical” location modifier



(c) Failed to recognise loculated pleural fluid despite having “apical” and “right hemithorax”



(d) Failed to recognise the rib position



(e) Mismatch between bounding box and salient region: Models attend to the salient region (enlarged area) to identify the abnormality instead of the entire heart.

Figure A.2: Failure cases from MS-CXR phrase grounding benchmark. Model outputs (latent vector similarity) are compared (from left, ClinicalBERT baseline, ConVIRT, QoRIA, and BioViL)

Table A.2: An extension of Table 6 to include Sensitivity and Specificity for the RSNA Pneumonia zero-shot and fine-tuned classification. We compare to GLoRIA scores reported in [30] which outperforms ConVIRT [84] (see [30]). Training size: GLoRIA ( $N = 186k$ , private dataset), BioViL ( $N = 146.7k$  of MIMIC-CXR).

| Method      | Type       | Text Model   | Loss           | % of labels | Acc.  | Sens. | Spec. | F1    | AUROC |
|-------------|------------|--------------|----------------|-------------|-------|-------|-------|-------|-------|
| SimCLR [6]  | Image only | -            | Global         | 1%          | 0.545 | 0.776 | 0.436 | 0.522 | 0.701 |
|             |            |              |                | 10%         | 0.760 | 0.663 | 0.806 | 0.639 | 0.802 |
|             |            |              |                | 100%        | 0.788 | 0.685 | 0.837 | 0.675 | 0.849 |
| GLoRIA [30] | Joint      | ClinicalBERT | Global & local | Zero-shot   | 0.70  | 0.89  | 0.65  | 0.58  | -     |
|             |            |              |                | 1%          | 0.72  | 0.82  | 0.69  | 0.63  | 0.861 |
|             |            |              |                | 10%         | 0.78  | 0.78  | 0.79  | 0.63  | 0.880 |
|             |            |              |                | 100%        | 0.79  | 0.87  | 0.76  | 0.65  | 0.886 |
| Baseline    | Joint      | ClinicalBERT | Global         | Zero-shot   | 0.719 | 0.648 | 0.781 | 0.614 | 0.812 |
| BioViL      | Joint      | CXR-BERT     | Global         | Zero-shot   | 0.732 | 0.831 | 0.685 | 0.665 | 0.831 |
|             |            |              |                | 1%          | 0.805 | 0.791 | 0.812 | 0.723 | 0.881 |
|             |            |              |                | 10%         | 0.812 | 0.781 | 0.826 | 0.727 | 0.884 |
|             |            |              |                | 100%        | 0.822 | 0.755 | 0.856 | 0.733 | 0.891 |

Table A.3: Mean IoU scores obtained on the newly released MS-CXR dataset, averaged over four runs with different seeds. The results are collected using different text encoder and training objectives (G&L: Global and local loss).

| Method       | Objective | Text encoder | Atelectasis | Cardiomegaly | Consolidation | Lung opacity | Edema | Pneumonia | Pneumothorax | Pl. effusion | Avg.  |
|--------------|-----------|--------------|-------------|--------------|---------------|--------------|-------|-----------|--------------|--------------|-------|
| Baseline     | Global    | ClinicalBERT | 0.228       | 0.269        | 0.293         | 0.173        | 0.268 | 0.249     | 0.084        | 0.232        | 0.224 |
| Baseline     | Global    | PubMedBERT   | 0.225       | 0.293        | 0.297         | 0.167        | 0.266 | 0.286     | 0.077        | 0.222        | 0.225 |
| ConVIRT [84] | Global    | ClinicalBERT | 0.257       | 0.281        | 0.313         | 0.177        | 0.272 | 0.238     | 0.091        | 0.227        | 0.238 |
| GLoRIA [30]  | G&L       | ClinicalBERT | 0.261       | 0.273        | 0.324         | 0.198        | 0.251 | 0.246     | 0.100        | 0.254        | 0.246 |
| BioViL       | Global    | CXR-BERT     | 0.296       | 0.292        | 0.338         | 0.202        | 0.281 | 0.323     | 0.109        | 0.290        | 0.266 |
| BioViL-L     | G&L       | CXR-BERT     | 0.302       | 0.375        | 0.346         | 0.209        | 0.275 | 0.315     | 0.135        | 0.315        | 0.284 |

## B Background in Chest Radiology

Chest X-rays are the most commonly performed diagnostic X-ray examination, and a typical text report for such an exam consists of three sections: a “Background” section describing the reason for examination and the exam type, a “Findings” section describing abnormalities as well as normal clinical findings in the scan, and an “Impression” section which summarises the findings and offers interpretation with possible recommendations. Multiple large Chest X-ray datasets have been released to the public (see [70] for an overview of CXR image datasets), including multi-modal ones of images and text such as MIMIC-CXR [33], some also accompanied by small sets of expert-verified ground-truth annotations of various nature, making the application a popular candidate for exploring self-supervised VLP on biomedical data.

The application area also possesses a strong clinical motivation. Globally, there is a shortage of qualified trained radiologists and a constantly increasing number of examinations in healthcare systems, workflows are hampered by issues such as a lack of standardisation in report writing, and fatigue-based errors occur too frequently. Thus, decision-support systems that can analyse incoming images or image-report pairs in order to provide real-time feedback to radiologists are a promising avenue towards improving workflow efficiency and the quality of medical image readings. In practice, the existing radiology workflow can for example be augmented via machine learning models by providing feedback on any incorrect or missing information in reports, and by standardising the reports’ structure and terminology.

### B.1 Key NLP and Dataset Challenges in Radiology

In this work, we focus on developing text and image models to enable clinical decision-support systems for biomedical applications via self-supervised VLP, without ground-truth annotations, and we conduct experiments in CXR applications. Image and text understanding in the biomedical domain is distinct from general-domain applications and requires careful consideration. Medical images are elaborately structured, which is reflected in the corresponding notes. To be able to harness the dense information captured in text notes for free-text natural language supervision, it becomes imperative to obtain finely tuned text models.

**Complex Sentence Structure.** Linguistic characteristics in radiology reports, many shared with related clinical text settings, decidedly differ from general domain text and thus require carefully tuned text models to acquire the best possible free-text natural language supervision in self-supervised VLP. For one, negations



are frequently used to indicate the absence of findings, in particular to make references as to how a patient’s health has evolved, e.g. “there are no new areas of consolidation to suggest the presence of pneumonia”. This sentence is for example falsely captured as positive for pneumonia by the automated CheXpert labeller [32]. Furthermore, as exemplified in this example, long-range dependencies are common, which makes understanding of relations within sentences challenging.

**Use of Modifiers.** Another characteristic is the use of highly specialised spatial language in radiology, which is crucial for correct diagnosis, often describing the positioning of radiographic findings or medical devices with respect to anatomical structures, see e.g. [12, 13]. The use of words like “medial”, “apical”, “bilateral” or “basilar” as spatial modifiers is unlikely to appear in the general domain but very common in CXR radiology reports. In addition to spatial modifiers, severity modifiers such as “mild”, “moderate” or “severe” are also commonly attached to an identified disorder or abnormality [17].

**Expressions of Uncertainty.** Another interesting difference to most general domain VLP applications and datasets such as Internet image captions, are expressions of uncertainty that one frequently encounters in radiology reports. We rarely expect to find an image caption to read “We see a person petting an animal, it is likely a dog but it could also be a cat”. In contrast, consider the following real radiology example: “New abnormality in the right lower chest could be either consolidation in the lower lobe due to rapid pneumonia or collapse, and/or moderate right pleural effusion, more likely abnormality in the lung because of absent contralateral mediastinal shift.” It is an extremely long description expressing uncertainty and containing long range dependencies.

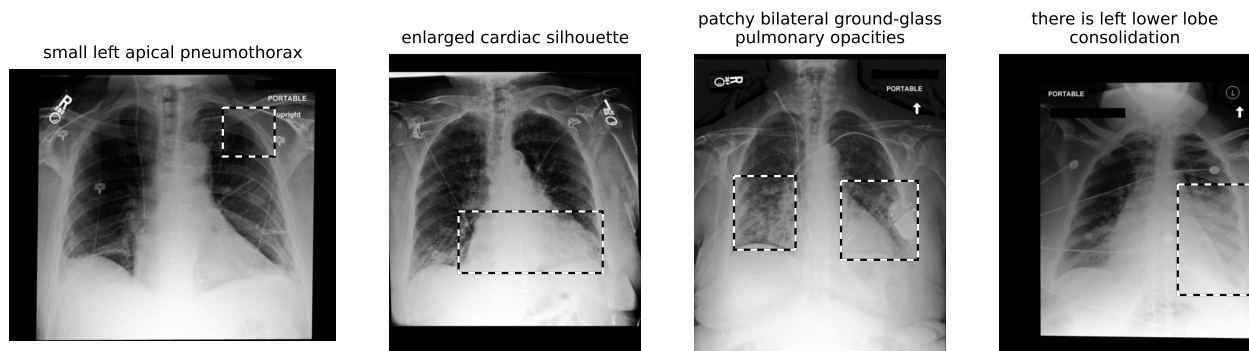
**Class Imbalance.** Finally, a challenge for many domain-specific VLP applications that is far less pronounced in the general domain setting is that of imbalanced latent entities. An example of such entities are the normal and anomalous findings in radiology images that doctors will describe in their report. In the CXR application, reports can roughly be divided into normal and abnormal scans, where abnormal ones reveal signs or findings observed during the exam [10]. Normal scans that do not show any signs of disease are far more common than any other findings, which leads to a larger number of false negatives in contrastive objectives compared to the general domain. An important detail is that normal scans tend to be expressed in specific forms and doctors frequently use templates to produce reports with no abnormalities.

## C MS-CXR Dataset Details

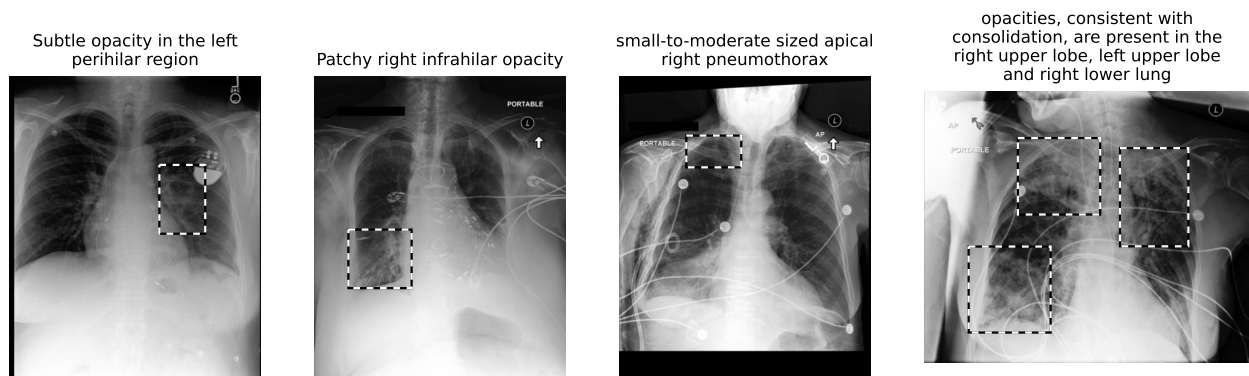
**General Overview.** With this new benchmark dataset, we provide bounding box and sentence pair annotations describing clinical findings visible in a given chest X-ray image. Each sentence describes a single pathology present in the image, and there could be multiple manually annotated bounding boxes corresponding to the description of the single radiological finding. Additionally, an image may have more than one pathology present, and we provide separate sets of bounding boxes for each phrase describing a unique pathology associated with an image. The annotations were collected on a subset of MIMIC-CXR images, which additionally contains labels across eight different pathologies: atelectasis, cardiomegaly, consolidation, edema, lung opacity, pleural effusion, pneumonia and pneumothorax. These pathologies were chosen based on the overlap between pathology classes present in the existing datasets and the CheXbert classifier [68]. In Fig. C.1 and Table C.1, we show some representative image and text examples from MS-CXR. Additionally, the distribution of samples across the pathology classes is shown in Table C.2 together with demographics across subjects in MS-CXR.

**Differences to Existing Annotations.** The proposed benchmark builds on top of publicly available bounding-box/ellipse annotations in MIMIC-CXR-Annotations [70] and REFLACX [36], where the former also contains simplified text phrases for pneumonia and pneumothorax. MS-CXR extends and curates these annotation sets by (I) reviewing their clinical correctness and suitability for the grounding task (see Section 3.1), (II) creating, verifying, and correcting bounding boxes where necessary, (III) pairing them up with real clinical descriptions extracted from MIMIC-CXR reports if none were present, and (IV) covering a wider range of clinical findings and pathologies. Most importantly, the textual descriptions paired with

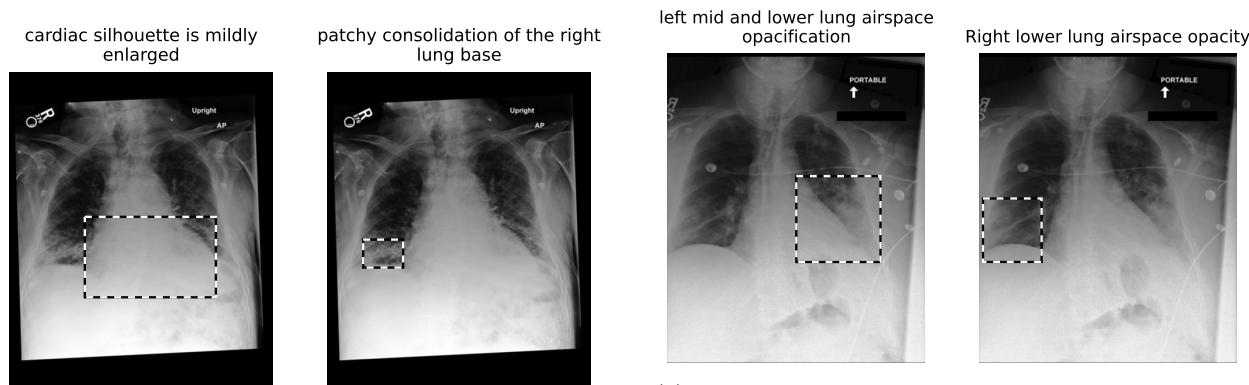




(a) Spatial extent of abnormalities ranging from highly localised to large and diffuse



(b) Complex spatial modifiers commonly seen in radiology reports



(c) Multiple pathologies reported for the same study

(d) Findings with multiple spatial locations reported separately

Figure C.1: We here provide some examples illustrating important axes of variability present in the MS-CXR dataset. Text descriptions include clinical findings of varying spatial extent (a) and a range of different spatial modifiers (b). Additionally, a subset of studies contain multiple bounding-box and sentence annotations per image (c-d).

Table C.1: Example findings in MS-CXR with complex syntactic structures. Please note how radiological sentences are most often not just a simple statement of the form “[class1, class2, ...]” that can be parsed with a simple bag-of-words approach, as in typical natural image captioning benchmarks (e.g., “A couple getting married” retrieved from Flickr30k [58]).

| Sentence  | Difficulty                   | Class            |
|---|------------------------------|------------------|
| “Abnormal opacity in the basilar right hemithorax is likely atelectasis involving the right lower and middle lobes”   | Complex syntactic structure  | Atelectasis      |
| “Multisegmental lower lobe opacities are present, consistent with areas of consolidated and atelectatic lung”   | Complex syntactic structure  | Atelectasis      |
| “Parenchymal opacification in the mid and lower lung”   | Less common expression       | Pneumonia        |
| “Air bronchograms extending from the left hilum throughout the left lung which has the appearance of infection”   | Complex location description | Pneumonia        |
| “Persistent focal bibasilar opacities, most consistent with infection”  | Domain-specific modifier     | Pneumonia        |
| “Widespread infection, less severe on the left”   | Location partially specified | Pneumonia        |
| “Airspace consolidation in the right upper, right middle and lower lobes”   | Multiple locations           | Pneumonia        |
| “Subsegmental-sized opacities are present in the bilateral infrahilar lungs”  | Domain specific modifiers    | Lung opacity     |
| “There continues to be a diffuse bilateral predominantly interstitial abnormality in the lungs with more focal vague opacity in the left upper peripheral lung” | Complex syntactic structure  | Lung opacity     |
| “Left apical pneumothorax”  | Domain-specific modifier     | Pneumothorax     |
| “Fluid level posteriorly, which represents a loculated hydropneumothorax”   | Domain-specific language     | Pneumothorax     |
| “Mild-to-moderate left pneumothorax”  | Severity modifier            | Pneumothorax     |
| “There is no pulmonary edema or pneumothorax, but small pleural effusions are still present”  | Negated disease entities     | Pleural effusion |
| “Pleural effusions are presumed but impossible to quantify, except say they are not large”  | Complex sentence structure   | Pleural effusion |

Table C.2: Distribution of the annotation pairs (image bounding-box and sentence) across different clinical findings. The demographic statistics (e.g., gender, age) of the subjects are collected from MIMIC-IV dataset for MS-CXR and all MIMIC-CXR.

| Findings                   | # of annotation pairs | # of subjects | Gender - F (%)   | Avg Age (std) |
|----------------------------|-----------------------|---------------|------------------|---------------|
| Atelectasis                | 61                    | 61            | 28 (45.90%)      | 64.52 (15.95) |
| Cardiomegaly               | 333                   | 282           | 135 (47.87%)     | 68.10 (14.81) |
| Consolidation              | 117                   | 109           | 40 (36.70%)      | 60.08 (17.67) |
| Edema                      | 46                    | 42            | 18 (42.86%)      | 68.79 (14.04) |
| Lung opacity               | 82                    | 82            | 33 (40.24%)      | 62.07 (17.20) |
| Pleural effusion           | 96                    | 95            | 41 (43.16%)      | 66.36 (15.29) |
| Pneumonia                  | 182                   | 146           | 65 (44.52%)      | 64.32 (17.17) |
| Pneumothorax               | 245                   | 151           | 66 (43.71%)      | 60.71 (18.04) |
| Total                      | 1162                  | 851           | 382 (44.89%)     | 64.37 (16.61) |
| Background (all MIMIC-CXR) | -                     | 65379         | 34134.0 (52.39%) | 56.85 (19.47) |

dense image region annotations are sampled from the original distribution of word tokens, which capture dense text semantics and are better aligned with real-world clinical applications that build on good local alignment.

## C.1 Label Collection and Review

We first parse original MIMIC reports and REFLACX [36] radiology transcripts by extracting sentences to form a large pool of text descriptions of pathologies. These candidates are later filtered by deploying the CheXbert [68] text classifier, in order to only keep phrases associated with the target pathologies whilst ensuring the following two criteria: (I) For a given study, there is only one sentence describing the target pathology, and (II) the sentence does not mention more than one findings that are irrelevant to each other. After extracting the text descriptions, they are paired with image annotations on a study level. At the final stage, a review process is conducted with two board certified radiologists mainly to verify the match between the text and bounding box candidates. Moreover, in this review process, we also assessed the suitability of the annotation pairs for the grounding task whilst ensuring clinical accuracy.

In detail, the phrase-image samples are filtered out if at least one of following conditions is met:

1. Text describing a finding not present in the image.
2. Phrase/sentence does not describe a clinical finding or describes multiple unrelated abnormalities that appear in different lung regions.
3. There is a mismatch between the bounding box and phrase, such as image annotations are placed incorrectly or do not capture the true extent of the abnormality.
4. High uncertainty is expressed regarding reported findings, e.g. “there is questionable right lower lobe opacity”.
5. Chest X-ray is not suitable for assessment of the finding or has poor image quality.
6. Text contains differential diagnosis or longitudinal information that prohibits correct grounding via the single paired image.
7. Sentences longer than 30 tokens, which often contain patient meta-information that is not shared between the two modalities (e.g., de-identified tokens).

Note that we only filter out phrases containing multiple findings, not images with multiple findings. For instance, if an image contains both pneumonia and atelectasis, with separate descriptions for each in the report, then we create two instances of phrase-bounding box pairs.

To further increase the size of our dataset, and to balance samples across classes, additional CXR studies are sampled at random, conditioned on the underrepresented pathologies. The following procedure is applied to create the pairs of image and text annotations for these selected studies: Text descriptions are extracted

using the same methodology outlined above, using MIMIC-CXR and ImaGenome datasets [78], where the latter provides sentence extracts from a subset of MIMIC-CXR dataset for clinical findings. However, differently from the initial step, the corresponding bounding box annotations (either one or more per sentence) are created from scratch by radiologists for the finding described in the text, and the same filtering as above is applied by the annotator to discard candidates if the image and/or sentence is found unsuitable for the grounding task.

**Patient Demographics.** As shown in Table C.2, the average age of subjects in MS-CXR is higher than the average for all subjects in MIMIC-CXR. We explain this observation with the fact that we do not sample studies from healthy subjects that do not display any anomalous findings and who are statistically likely to be younger. Similarly, we do not expect gender bias to be present due to our sampling as none of the pathologies we sample are gender-specific. Overall MS-CXR does not deviate far from the MIMIC-CXR distribution.

## D Related Work

Here we provide a more detailed overview of related work to complement the brief review provided in the main article.

**Joint Image-Text Representation Learning.** A variety of self-supervised VLP approaches have been proposed towards jointly learning visual and textual representations of paired data without supervision, such as frameworks using contrastive objectives [26, 42, 60], approaches based on joint transformer architectures [40, 41, 51, 69], self-supervised VLP with word-region alignment and language grounding [7], and text prediction tasks to learn image features [15]. For example, [60] use a contrastive loss over embeddings of text and image pairs to train a model on large data collected from the internet ( $\sim 400$ M pairs) enabling zero-shot transfer of the model to downstream tasks. Some of the proposed approaches utilise a single architecture, usually a transformer, to learn a representation, following encoders for the individual modalities [7, 41, 69]. Another common theme is the use of cross-modal attention mechanisms to improve the aggregation of image regions in convolutional architectures [1, 11, 26].

A number of different objectives have been explored for representation learning in VLP, including the prediction of words in image captions [35], predicting phrase n-grams [39], predicting of entire captions [15], *global* contrastive objectives defined on the embeddings of the entire image and text instances [84], and combinations of global and *local* contrastive terms [30, 55], where local means that objectives are defined over text fragments (words or phrases) and image regions.

A task closely related to instance representation learning in VLP is *phrase grounding*, also known as visual grounding, phrase localisation, local alignment, or word-region alignment. The goal here is to connect natural language descriptions to local *image regions*. In a supervised learning setting such as in [52, 54], this problem requires expensive manual annotation for region-phrase correspondence. Thus, settings for visual grounding have been explored in which cross-modal pairs are the only form of supervision that is available [7, 11, 21, 26, 48, 74], i.e. the supervision signal is the knowledge of which caption belongs to which image. This setting of paired images and text has also been referred to as weakly supervision. Much of the general domain prior work on phrase grounding relies on off-the-shelf object-detection networks [7, 11, 26, 74, 82, 85] such as Faster R-CNN [63] which are pretrained on large labelled datasets to extract region candidates from images. This considerably simplifies the problem of matching regions to phrases as the set of possible regions to match can be assumed to be known, a luxury that is often unavailable in domain specific contexts.

**Biomedical VLP Representation Learning.** Several studies [29, 30, 44, 55, 84] have explored joint representation learning for paired image and text data in the medical domain. Contrastive VISual Representation Learning from Text (ConVIRT) [84] uses a contrastive learning formulation for instance-level representation learning from paired medical images and text. The authors uniformly sample sentences and maximise their similarity to true augmented paired images via the InfoNCE contrastive loss [57], while reducing similarity between negative pairs in the same batch. [30, 55] both introduce approaches that combine instance-level image-report contrastive learning with local contrastive learning for medical data. In contrast, [44] use a local-only objective in an approach that approximates the mutual information between grid-like local features of images and sentence-level text features of medical data. The formulation learns image and

Table D.1: Example findings in ImaGenome which would make grounding of phrases difficult.

| Sentence   | Difficulty  | Annotated Finding |
|--|---|-------------------|
| “Even though Mediastinal veins are more distended, previous pulmonary vascular congestion has improved slightly, but there is more peribronchial opacification and consolidation in both lower lobes which could be atelectasis or alternatively results of recent aspiration, possibly progressing to pneumonia.” | Multiple findings, uncertainty, different sub-parts of lung | Pneumonia         |
| “Moderate right pleural effusion and bilateral heterogenous airpace opacities, concerning for pneumonia.”  | Multiple findings, differing laterality                     | Pneumonia         |
| “It could be an early infection”   | Region unclear  | Pneumonia         |
| “There is also a new small left-sided pleural effusion.”   | Differential diagnosis, there could be another effusion     | Effusion          |

text encoders as well as a discriminator trained to distinguish positive and negative pairs. While most related approaches use no ground truth, [5] study a semi-supervised edema severity classification setting, and [27] assume sets of seen and unseen labels towards zero-shot classification on CXR data. [43] evaluate pretrained joint embedding models—general domain VLP representation learning models that use a transformer to learn a joint embedding—by fine-tuning the models on CXR data.

Multiple CXR datasets exist that enable a partial evaluation of phrase grounding, but all come with some limitations we hope to mitigate with our MS-CXR dataset (see Section 3.1). VinDr [56], RSNA Pneumonia [65], and the NIH Chest X-ray Dataset [75] are datasets that provide bounding-box image annotations, but lack accompanying free-text descriptions. REFLACX [36] provides gaze locations captured with an eye tracker, dictated reports and some ground truth annotations for gaze locations, but no full phrase matches to image regions. Phrase annotations for MIMIC-CXR data released in [70] are of small size (350 studies), only contain two abnormalities, and for some samples have shortened phrases that were adapted to simplify the task. ImaGenome [78] provides a large number of weak local labels for CXR images and reports, with a focus on anatomical regions. However, its ground-truth set is smaller (500 studies), bounding-box regions annotate anatomical regions rather than radiological findings. Furthermore, ImaGenome sentence annotations are not curated, see Table D.1 for some examples. Sentences often contain multiple diseases as well as uncertain findings, making an accurate, largely noiseless grounding evaluation difficult. Some sentences also contain differential diagnosis and temporal change information, which cannot be grounded without access to prior scans.

**Language Modelling in Radiology.** Most recent general domain VLP work relies on transformer based contextual word embedding models, in particular BERT [16], pretrained on general domain data from newswire and web domains such as Wikipedia. But specific domains often exhibit differences in linguistic characteristics from general text and even related domains, such as between clinical and non-clinical biomedical text as noted in [2], motivating the use of more specialised language models in most related work with a focus on the medical domain. Here, related multi-modal work commonly uses publicly available models including BioBERT [38], ClinicalBERT [2], BioClinicalBERT [2], or PubMedBERT [25], which are either trained from scratch or fine-tuned via continual pretraining using a Masked Language Modelling (MLM) objective. Sometimes additional objectives are added such as adversarial losses [46] or Next Sentence Prediction. [25] provide evidence that training language models from scratch for specialised domains with abundant amounts of unlabelled text can result in substantial gains over continual pretraining of models first fit to general domain text. The specialised corpora these biomedical and clinical domain models use include PubMed abstracts and PubMed Central full texts, and de-identified clinical notes from MIMIC-III [34]. All the aforementioned language models have a pre-specified vocabulary size consisting of words and subwords, usually 30,000 words in standard BERT. The in-domain vocabulary plays a particularly important role in representative power for a specialised domain. A vocabulary that is not adapted will break up more words into subwords and additionally contain word pieces that have no specific relevance in the specialised domain, hindering downstream learning (see e.g. [25]). As [25] highlight, BERT models that use continual pretraining are stuck with the original vocabulary from the general-domain corpora.

Table E.1: Hyper-parameter values used for image data augmentations.

|   | Image-Text Pretraining | Image-only Pretraining | Fine-tuning for Downstream Tasks |
|---|------------------------|------------------------|----------------------------------|
| Affine transform – shear                    | 15°                    | 40°                    | 25°                              |
| Affine transform – angle                    | 30°                    | 180°                   | 45°                              |
| Colour jitter – brightness                  | 0.2                    | 0.2                    | 0.2                              |
| Colour jitter – contrast                    | 0.2                    | 0.2                    | 0.2                              |
| Horizontal flip probability                 | -                      | 0.5                    | 0.5                              |
| Random crop scale                           | -                      | (0.75, 1.0)            | -                                |
| Occlusion scale                             | -                      | (0.15, 0.4)            | -                                |
| Occlusion ratio                             | -                      | (0.33, 0.3)            | -                                |
| Elastic transform ( $\sigma, \alpha$ ) [67] | -                      | (4, 34)                | -                                |
| Elastic transform probability               | -                      | 0.4                    | -                                |
| Gaussian noise                              | -                      | 0.05                   | -                                |

Other closely related tasks in the CXR domain that share similar NLP challenges include report summarisation [10, 83], automatic report generation [8, 45, 53], and natural language inference for radiology reports [53]. Finally, while the name implies close similarity to our CXR-BERT, CheXbert [68] is a BERT based sentence classification model developed for improving the CheXpert [32] labeller, and the model does not have a domain-specific vocabulary like ours or PubMedBERT.

We note that most related work on self-supervised multi-modal learning on CXR data neither explores text augmentation nor maintains text losses such as MLM during multi-modal training. An exception is found in [55], who use the Findings and Impression/Assessment sections of radiology reports, and randomly change the sentence order by swapping pairs of them.

## E Model Details

### E.1 CXR-BERT Pretraining Details

Our CXR-BERT text encoder is based on the BERT (base size) architecture [72]. We adopt an implementation available via the Huggingface transformers library [77]. The model weights are randomly initialised and pretrained from scratch. As described in Section 2.1, CXR-BERT is pretrained in three phases before the joint pretraining phase. For Phase (I), we use the Huggingface tokeniser library<sup>4</sup> to generate our custom WordPiece vocabulary of 30k tokens. For Phase (II), we use the AdamW [50] optimiser with a batch size of 2048 sequences and a linear learning rate schedule over 250k training steps with a 5% warm up period. We set a base learning rate of 4e-4. Following RoBERTa [47], we pack multiple sentences into one input sequence of up to 512 tokens and use dynamic whole-word masking. In Phase (III), we continue pretraining the model using only MIMIC-CXR text reports. In addition to the MLM loss, we add our RSM loss to pretrain the projection layer. The projection layer  $P_{\text{txt}}$  is used to project the 768-dimensional feature vector  $\mathbf{t}$  to a 128-dimensional report representation  $\mathbf{t}$ . We use the AdamW optimiser with a batch size of 256 sequences and a linear learning rate schedule over 100 epochs with a 3% warm up period. We set the base learning rate to 2e-5.

### E.2 Image Encoder

**Pretraining Details.** For the image encoder, we adopt the ResNet50 [28] architecture. The 2048-dimensional feature maps  $\tilde{\mathbf{V}}$  of the ResNet50 are projected to 128-dimensional feature maps  $\mathbf{V}$  using a two-layer perceptron  $P_{\text{img}}$  implemented with  $1 \times 1$  convolutional layers and batch-normalisation [31]. The global image representation  $\mathbf{v}$  is obtained by average-pooling the projected local features  $\mathbf{V}$ . Prior to image-text joint training, the model weights are randomly initialised and pretrained on MIMIC-CXR images using SimCLR [6] — an image-only self-supervised learning approach. We use a large-batch optimisation (LARS) technique [80] on top of ADAM with a batch size of 256 and a linear learning rate scheduler over 100 epochs with a 3% warm up period. We set the base learning rate to 1e-3.

<sup>4</sup><https://github.com/huggingface/tokenizers>



**Augmentations.** For each training stage, we apply a different set of image augmentations to have a better control over the learnt feature invariances (e.g., laterality). During the image-text joint pretraining stage, we use affine transformations (random rotation and shearing) and contrast and brightness colour jitter. Unlike ConVIRT [84] and GLoRIA [30], we do not apply horizontal flips during the joint training to preserve location information (e.g. “pneumonia in the left lung”). During the image-only SSL (SimCLR) pretraining phase, we use additional image augmentations including random occlusion, additive Gaussian noise, and elastic spatial transforms [67]. We use the implementations available in the torchvision library<sup>5</sup>. The image augmentation parameters and their corresponding values are listed in Table E.1. Before applying these transformations, we normalise the input image intensities by re-scaling each colour channel values to the  $[0, 255]$  range. During inference, we only apply centre cropping and resizing.

---

<sup>5</sup><https://pytorch.org/vision/stable/transforms.html>