# CXR-CLIP: Toward Large Scale Chest X-ray Language-Image Pre-training

Kihyun You<sup>1</sup>, Jawook Gu<sup>1</sup>, Jiyeon Ham<sup>1</sup>, Beomhee Park<sup>1</sup>, Jiho Kim<sup>1</sup>, Eun K. Hong<sup>1</sup>, Woonhyuk Baek<sup>1</sup>, and Byungseok Roh<sup>1</sup>

Kakaobrain, Seongnam, Republic of Korea {ukihyun, jawook.gu, jiyeon.ham, brook.park, tyler.md, amy.hong, wbaek, peter.roh}@kakaobrain.com

Abstract. A large-scale image-text pair dataset has greatly contributed to the development of vision-language pre-training (VLP) models, which enable zero-shot or few-shot classification without costly annotation. However, in the medical domain, the scarcity of data remains a significant challenge for developing a powerful VLP model. In this paper, we tackle the lack of image-text data in chest X-ray by expanding image-label pair as image-text pair via general prompt and utilizing multiple images and multiple sections in a radiologic report. We also design two contrastive losses, named ICL and TCL, for learning study-level characteristics of medical images and reports, respectively. Our model outperforms the state-of-the-art models trained under the same conditions. Also, enlarged dataset improve the discriminative power of our pre-trained model for classification, while sacrificing marginal retrieval performance. Code is available at https://github.com/kakaobrain/cxr-clip.

**Keywords:** Chest X-ray · Vision-Language Pre-training · Contrastive Learning

### 1 Introduction

Chest X-ray (CXR) plays a vital role in screening and diagnosis of thoracic diseases [20]. The effectiveness of deep-learning based computer-aided diagnosis has been demonstrated in disease detection [22]. However, one of the major challenges in training deep learning models for medical purposes is the need for extensive, high-quality clinical annotation, which is time-consuming and costly.

Recently, CLIP [23] and ALIGN [11] have shown the ability to perform vision tasks without any supervision. However, vision-language pre-training (VLP) in the CXR domain still lacks sufficient image-text datasets because many public datasets consist of image-label pairs with different class compositions. Med-CLIP [27] attempted to a rule-based labler to use both image-text data and image-label data. However, it relies on the performance of the rule-based labeler and is not scalable to other diseases that the labeler cannot address.

In this paper, we propose a training method, *CXR-CLIP*, that integrates image-text data with image-label data using class-specific prompts made by

radiologists. Our method does not depend on a rule-based labeler and can be applied to any image-label data. Also, inspired by DeCLIP [14], we used Multi-View Supervision (MVS) utilizing multiple images and texts in a CXR study to make more image-text pairs for efficient learning. In addition, we introduce two contrastive loss functions, named image contrastive loss (ICL) and text contrastive loss (TCL), to learn study-level characteristics of the CXR images and reports respectively.

The main contributions of this paper are summarized as follows. 1) We tackle the lack of data for VLP in CXR by generating image-text pairs from image-label datasets using prompt templates designed by radiologists and utilizing multiple images and texts in a study. 2) Two additional contrastive losses are introduced to learn discriminate features of image and text, improving image-text retrieval performances. 3) Performance of our model is validated on diverse datasets with zero-shot and few-shot settings.

# 2 Related Work

**Data Efficient VLP** Recent studies [14,18] have proposed data-efficient VLP via joint learning with self-supervision. DeCLIP [14] suggested MVS that utilizes image and text augmentation to leverage positive pairs along with other self-supervisions. In CXR domain, GloRIA [8] aligned words in reports and sub-regions in an image for label efficiency, and BioVIL [2] combined self-supervision for label efficiency. We modify MVS as two distinct images and texts from a study and present self-supervised loss functions, ICL and TCL for efficient learning.

Self-supervision within CXR study A CXR study could include several images in different views and two report sections: 'findings' and 'impression'. The impression section includes the differential diagnosis inferred from the findings section. BioVIL [2] enhanced the text encoder by matching two sections during language pre-training. MedAug [25] shows that self-supervised learning by matching images in a study is better than differently augmented images. We utilize both of multiple images and texts from a single study in VLP in an end-to-end fashion.

Leveraging image-label data in VLP MedCLIP [27] integrated unpaired images, texts, and labels using rule-based labeler [9], which is less capable of retrieving the exact report for a given image due to the effect of decoupling image-text pairs. UniCL [29] suggested using prompts to leverage image-label dataset [4], considering the samples from the same label to be a positive pair. To our knowledge, this is the first work to utilize prompting for training in CXR domain.

## 3 Method

CXR-CLIP samples image-text pairs from not only image-text data but also image-label data, and learns study-level characteristics with two images and two texts per study. The overview of the proposed method is illustrated in Fig 1.





**Fig. 1.** Overview of the proposed method with a training batch sampling n studies, where each study has a pair of images  $(x^1, x^2)$  and a pair of text  $(t^1, t^2)$ . If a study has one image or one text, data augmentation is conducted to make second examples. For the image-label data, two different prompts are generated from class labels as  $(t^1, t^2)$ . Using sampled pairs, the encoders are trained with three kinds of contrastive losses (MVS, ICL, and TCL).

### 3.1 Data Sampling

We define a CXR study as  $s = \{X, T\}$ , where X is a set of images, and T is a set of "findings" and "impression" sections. The study of image-label dataset has a set of image labels Y instead of T. For the image-label dataset, we make promptbased texts  $T = Concat(\{p \sim P(y)\}_{y \in Y})$ , where p is a sampled prompt sentence, P(y) is a set of prompts given the class name and value y, and  $Concat(\cdot)$  means concatenating texts. The set of prompts is used to generate sentences such as actual clinical reports, taking into account class labels and their values (positive, negative, etc.), unlike the previous prompt [8] for evaluation which randomly combines a level of severity, location, and sub-type of disease. Our prompts are available in Appendix.

We sample two images  $(x^1, x^2)$  in X if there are multiple images. Otherwise, we use augmented image  $A_i(x^1)$  as  $x^2$ , where  $A_i$  is image augmentation. To leverage various information from different views in CXR (AP, PA, or lateral), we sample images from two distinct views as possible. Similarly, we sample two texts  $(t^1, t^2)$  in T if there are both "findings" and "impression". Otherwise, we use augmented text  $A_t(t^1)$  as  $t^2$ , where  $A_t$  is text augmentation. For the imagelabel data, we sample two prompt sentences as  $t^1$  and  $t^2$  from the constructed  $T = Concat(\{p \sim P(y)\}_{y \in Y}).$ 

#### 3.2 Model Architecture

We construct image encoder  $E^i$  and text encoder  $E^t$  to obtain global representations of image and text, and a projection layer  $f^i$  and  $f^t$  to match the size of final embedding vectors.

**Image Encoder** We have tested two different image encoders; ResNet-50 [7] and Swin-Tiny [15] as follow [8,27]. We extract global visual features from the global average pooled output of the image encoder. A linear layer is adopted to project the embeddings into the same size as text embeddings. The normalized visual embedding v is obtained by  $v = f^i(E^i(x)) / ||f^i(E^i(x))||$ . We denote a batch of the visual embeddings as  $V = \{v\}_{i=1}^n$ , where n is a batch size.

**Text Encoder** We use BioClinicalBERT [1] model, which is the same architecture as BERT [5] but pre-trained with medical texts [12] as follow [8,27]. We use **[EOS]** token's final output as the global textual representation. Also, a linear projection layer is adopted the same as the image encoder. The normalized text embedding u is denoted as  $u = f^t(E^t(t)) / ||f^t(E^t(t))||$ . We denote a batch of the text embedding as  $U = \{u\}_{i=1}^n$  and  $(v_i, u_i)$  are paired.

#### 3.3 Loss Function

In this section, we first describe CLIP loss [23] and then describe our losses (MVS, ICL, TCL) in terms of CLIP loss. The goal of CLIP loss is to pull image embedding and corresponding text embedding closer and to push unpaired image and text farther in the embedding space. The InfoNCE loss is generally adopted as a type of contrastive loss, and CLIP uses the average of two InfoNCE losses; image-to-text and text-to-image. The formula for CLIP loss is given by

$$L_{CLIP}(U,V) = -\frac{1}{2n} \left( \sum_{u_i \in U} \log \frac{\exp(v_i^T u_i/\tau)}{\sum_{v_j \in V} \exp(u_i^T v_j/\tau)} + \sum_{v_i \in V} \log \frac{\exp(u_i^T v_i/\tau)}{\sum_{u_j \in U} \exp(v_i^T u_j/\tau)} \right)$$
(1)

, where  $\tau$  is a learnable temperature to scale logits.

In DeCLIP [14], MVS uses four  $L_{CLIP}$  loss with all possible pairs augmented views;  $(x, t), (x, A_t(t)), (A_i(x), t)$  and  $(A_i(x), A_t(t))$ . We modify DeCLIP's MVS to fit the CXR domain by the composition of the second example. DeCLIP only utilizes an augmented view of the original sample, but we sample a pair of the

Data	]	Evaluation					
Split	MIMIC-CXR	CheXpert	ChestX-ray14	VinDR	$\operatorname{RSNA}$	$\operatorname{SIIM}$	Open-I
Train	222,628	$216,\!478$	89,696	12,000	$18,\!678$	8,422	
Valid	1,808	233	22,423	$3,\!000$	4,003	1,808	
Test	3,264	1,000		3,000	4,003	1,807	3,788

Table 1. The number of studies for each dataset and split in this paper

second image and text as described in 3.1. We denote the first and the second sets of image embeddings as  $U^1$ ,  $U^2$ , and text embeddings as  $V^1$ ,  $V^2$ .

$$L_{MVS} = \frac{1}{4} (L_{CLIP}(U^1, V^1) + L_{CLIP}(U^2, V^1) + L_{CLIP}(U^1, V^2) + L_{CLIP}(U^2, V^2))$$
(2)

The goal of ICL and TCL is to learn modality-specific characteristics in terms of image and text respectively. We design ICL and TCL as same as CLIP loss, but the input embeddings are different. ICL only uses image embeddings;  $L_{ICL} = L_{CLIP}(V^1, V^2)$  and TCL only uses text embeddings;  $L_{TCL} = L_{CLIP}(U^1, U^2)$ . ICL pulls image embeddings from the same study and pushes image embeddings from the different studies, so that, the image encoder can learn study-level diversity. Similarly, TCL pulls embeddings of "findings" and "impression" in the same study or diverse expressions of prompts from the same label and pushes the other studies' text embeddings, so that the text encoder can match diverse clinical expressions on the same diagnosis. Thereby, the final training objective consists of three contrastive losses balanced each component by  $\lambda_I$  and  $\lambda_T$ , formulated by  $L = L_{MVS} + \lambda_I L_{ICL} + \lambda_T L_{TCL}$ .

## 4 Experiment

#### 4.1 Datasets

We used three pre-trained datasets and tested with various external datasets to test the generalizability of models. The statistics of the datasets used are summarized in Table 1.

**MIMIC-CXR** [13] consists of CXR studies, each with one or more images and free-form reports. We extracted "findings" and "impression" from the reports. We used the training split for pre-training and the test split for image-to-text retrieval.

**CheXpert** [9] is an image-label data with 14 classes, obtained from the impression section by its rule-based labeler, and each class is labeled as positive, negative, uncertain, or none (not mentioned). We used the training split for pre-training with class-specific prompts. **CheXpert5x200** is a subset of CheXpert for 5-way classification, which has 200 exclusively positive images for each class. Note that only the reports of CheXpert5x200 are publicly available, but the reports of CheXpert are not. Following the previous works [8,27], we excluded CheXpert5x200 from the training set and used it for test.

**ChestX-ray14** [26] consists of frontal images with binary labels for 14 diseases. Prompts are generated by sampling 3 negative classes per study. We used 20% of the original training set for validation, and the remaining 80% for pre-training.

**RSNA pneumonia** [24] is binary-labeled data as pneumonia or normal. We split train/valid/test set 70%, 15%, 15% of the dataset following [8] for the external classification task.

**SIIM Pneumothorax**<sup>1</sup> is also binary labeled as pneumothorax or normal. We split the train/valid/test set same ratio as RSNA pneumonia following [8] and used it for the classification task.

VinDR-CXR [19] contains 22 local labels and 6 global labels of disease, which were obtained by experienced radiologists. We split the validation set from the original training set. Of 28 classes, "other diseases" and "other lesions" classes were excluded. Then, only 18 classes having 10 or more samples within the test set were evaluated for the binary classification of each class as follow [10].

**Open-I** [3] is an image-text dataset. From each study, one of the report sections and one frontal-view image were sampled and used for image-to-text retrieval.

#### 4.2 Implementation Details

We used augmentations  $A_i$  and  $A_t$  to fit medical images and reports. For  $A_i$ , we resize and crop with scale [0.8, 1.1], randomly adapt CLAHE [21], and random color jittering; brightness, hue ratios from [0.9, 1.1] and contrast, saturation [0.8, 1.2]. For  $A_t$ , to preserve clinical meaning, sentence swap and back-translation<sup>2</sup> from Italian to English is used. The image size and final-embedding size are set to 224 and 512 respectively as in previous work [27]. We set  $\lambda_I$  and  $\lambda_T$ to 1.0, 0.5 for balancing total loss. Two encoders were trained for 15 epochs in a mixed-precision manner, early stopped by validation loss, and optimized by AdamW [17] with an initial learning rate 5e-5 and a weight decay 1e-4. We used cosine-annealing learning-rate scheduler [16] with warm-up for 1 epoch. A training batch consists of 128 studies with 256 image-text pairs. We implemented all experiments on PyTorch with 4 NVIDIA V100 GPUs.

#### 4.3 Comparison with State-of-the-arts

**Zero-shot and few-shot classification** Table 2 shows performance on classification tasks of our models and state-of-the-art models. To evaluate zero-shot classification fairly, we used evaluation prompts suggested from previous works [2,8,10]. The evaluation prompts are available in Appendix. We evaluate binary classification computed by Area Under ROC (AUC) and multi-class classification computed by accuracy (ACC). Our ResNet model trained with MIMIC-CXR outperforms GloRIA [8] except for CheXpert5x200, as GloRIA

<sup>&</sup>lt;sup>1</sup> https://siim.org/page/pneumothorax challenge

<sup>&</sup>lt;sup>2</sup> https://huggingface.co/Helsinki-NLP

**Table 2.** Comparison with state-of-the-art for zero-shot(ZS) or few-shot(10%) classification tasks. M, C, and C14 mean MIMIC-CXR, CheXpert, and ChestX-ray14, respectively. C<sup>\*</sup> means CheXpert with reports, which are not publicly available. ResNet50 (R50) and SwinTiny (SwinT) mean the image encoder used for each model.

Madal Nama	Pre-train	Vir	DR-0	CXR		RSNA	ł		SIIM		C5x200
Model Name	Dataset	ZS	10%	100%	ZS	10%	100%	ZS	10%	100%	ZS-ACC
$GloRIA_{R50}$	C*	78.0	73.0	73.1	80.6	88.2	88.5	84.0	91.5	91.9	$62.4^{*}$
$\text{CXR-CLIP}_{R50}$	Μ	78.8	82.1	82.2	83.3	88.5	89.2	85.2	88.3	90.5	56.2
$CXR-CLIP_{SwinT}$	Μ	78.3	84.9	85.4	81.3	88.0	88.4	85.5	86.9	88.3	54.3
$MedCLIP_{SwinT}$	M,C	82.4	84.9	85.1	81.9	88.9	89.0	89.0	90.4	90.8	59.2
$\text{CXR-CLIP}_{R50}$	$^{\rm M,C}$	83.0	81.4	82.1	81.7	88.5	88.9	86.4	88.4	90.7	61.7
$\text{CXR-CLIP}_{SwinT}$	M,C	82.7	86.1	86.7	84.5	88.1	88.8	87.9	89.6	91.2	60.1
$\text{CXR-CLIP}_{R50}$	M,C,C14	78.1	80.2	81.0	81.8	88.7	89.3	85.2	91.5	92.8	60.3
$CXR-CLIP_{SwinT}$	M,C,C14	78.9	88.0	89.0	80.1	89.2	89.8	91.4	92.9	94.0	62.8

 Table 3. Comparison with state-of-the-arts for image-to-text retrieval. The notations of datasets and models are same to Table 2.

Model Name	Pre-Train	Chel	Xpert	5x200	MI	MIC-0	CXR		Open	-I	Total
model manie	Dataset	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	RSUM
$GloRIA_{R50}$	C*	17.8	38.8	<b>49.9</b>	7.2	20.6	30.3	1.5	4.4	6.5	177.0
$\text{CXR-CLIP}_{R50}$	Μ	9.4	23.0	32.6	21.4	46.0	59.2	<b>3.8</b>	8.2	12.3	216.9
$\operatorname{CXR-CLIP}_{SwinT}$	Μ	8.4	21.5	30.2	21.6	<b>48.9</b>	60.2	3.6	8.3	11.5	214.2
$MedCLIP_{SwinT}$	M,C	2.6	3.0	3.6	1.1	1.4	5.5	0.1	0.4	0.7	18.4
$\text{CXR-CLIP}_{R50}$	$^{\rm M,C}$	5.5	19.2	27.4	20.2	45.9	58.2	3.5	8.2	12.0	200.1
$\operatorname{CXR-CLIP}_{SwinT}$	$^{\rm M,C}$	8.5	23.0	31.6	19.6	44.2	57.1	3.1	8.3	11.6	207.0
$\text{CXR-CLIP}_{R50}$	M,C,C14	5.7	18.0	28.3	19.7	44.4	56.4	2.3	6.7	10.1	191.6
$\operatorname{CXR-CLIP}_{SwinT}$	M,C,C14	7.0	20.1	29.7	20.9	46.2	58.8	2.4	6.6	9.4	201.1

trained with image-text pair in CheXpert. Our SwinTiny model trained with MIMIC-CXR and CheXpert outperforms MedCLIP [27], which is the same architecture trained with the same datasets, in most of the metrics. Adding more pre-training datasets by prompting image-label datasets tends to improve performance for classifications, while the SwinTiny CXR-CLIP pre-trained with three datasets, performs the best for most of the metrics. More comparison with self-supervised models is available in Appendix.

Image-to-text retrieval We evaluated image-to-text retrieval computed by R@K, the recall of the exact report in the top K retrieved reports for a given image. (Table 3) While GloRIA [8] uses image-text pairs in CheXpert(C<sup>\*</sup>) which is not available in public, CXR-CLIP uses image-text in MIMIC-CXR. So we adapt an external image-text dataset Open-I [3] for a fair comparison. GloRIA has the best performance on CheXpert but our model trained with MIMIC-CXR, which has similar amounts of studies to CheXpert, outperforms on Open-I. MedCLIP almost lost the ability to retrieve image-text due to decoupling pairs of image and text during pre-training. In CXR-CLIP, adding more image-label datasets such as CheXpert and ChestX-ray14 degrades the image-text retrieval

**Table 4.** Ablations and comparison with CLIP [23] and DeCLIP [14]. Our augmentations effectively preserves clinical meaning than EDA. Our full methodology (CXR-CLIP) outperforms DeCLIP.

Mathad	CheXpert 5x200			MIMIC-CXR			Total	
Method	ACC	R@1	R@5	R@10	R@1	R@5	R@10	RSUM
Vanila CLIP	58.9	4.4	14.4	22.6	17.3	41.2	52.6	152.5
+ Study Level Sampling	58.7	4.6	15.1	23.2	17.8	42.5	54.2	157.4
+ Augmentations	60.6	5.7	17.0	24.9	16.1	40.2	51.5	155.4
+ MVS	61.2	5.4	17.1	24.7	16.3	40.6	53.3	157.4
+ ICL	61.6	6.8	<b>20.3</b>	28.6	17.5	41.6	53.2	168.0
+ TCL (CXR-CLIP)	61.7	6.2	18.2	29.1	19.6	<b>44.8</b>	56.6	174.5
MVS of DeCLIP (EDA)	59.5	3.2	15.5	22.9	15.8	39.1	51.5	148.0
MVS of DeCLIP (Our aug)	59.4	6.0	17.0	24.4	15.1	38.8	51.8	153.1
DeCLIP (Our aug)	59.4	5.7	16.1	24.6	18.1	44.0	55.3	163.8

performance, possibly because the contribution of the text in original reports was diluted.

#### 4.4 Ablations

For the ablation study, models with ResNet-50 [7] backbone were trained on MIMIC-CXR and CheXpert datasets and tested on zero-shot classification and image-to-text retrieval tasks with MIMIC-CXR and CheXpert5x200 datasets.

We conducted two ablations shown in Table 4. First, we analyzed the effect of each component of CXR-CLIP by adding the components to vanilla CLIP [23] one by one. To validate our data sampling closer, we divided the sampling method into three parts 1) study-level sampling 2) data augmentations 3) Multi-view and Multi-text sampling (MVS). Our study-level sampling strategy improves performance compared to vanilla CLIP, which uses a naive sampling method bringing an image and corresponding report. Additionally, the modified data augmentation to fit the CXR domain contributes to performance increment of classification, the similar performance on retrieval. MVS slightly improves performances in both classification and image-text retrieval. Adding more supervision (ICL and TCL) improves performance by utilizing better multi-views and multi-text inputs. However, TCL drops the performance of recalls in CheXpert5x200, TCL could be hard to optimize variation of the radiologic report and prompt not diverse as images.

In the second ablation study, CXR-CLIP was compared to DeCLIP [14] to confirm that our MVS using two image-text pairs per study is better than the MVS of DeCLIP which uses naively augmented images and texts. We show that our text augmentation outperforms DeCLIP's text augmentation named EDA [28] in terms of image-to-text recall, which implies our text augmentation preserves clinical meaning. The superiority of our MVS over DeCLIP's MVS confirms that using multiple images and texts from one study is better than using images and texts from augmented examples. Also, our full methodology (CXR-CLIP) outperforms DeCLIP, suggesting that our method efficiently learns in the CXR domain more than DeCLIP.

## 5 Conclusion

We presented a framework enlarging training image-text pair by using imagelabel datasets as image-text pair with prompts and utilizing multiple images and report sections in a study. Adding image-label datasets achieved performance gain in classification tasks including zero-shot and few-shot settings, on the other hand, lost the performance of retrieval tasks. We also proposed loss functions ICL and TCL to enhance the discriminating power within each modality, which effectively increases image-text retrieval performance. Our additional loss functions are designed to efficiently learn CXR domain knowledge along with image-text contrastive learning.

### References

- Alsentzer, E., Murphy, J.R., Boag, W., Weng, W., Jin, D., Naumann, T., McDermott, M.B.A.: Publicly available clinical BERT embeddings. CoRR abs/1904.03323 (2019), http://arxiv.org/abs/1904.03323
- Boecking, B., Usuyama, N., Bannur, S., Castro, D.C., Schwaighofer, A., Hyland, S., Wetscherek, M., Naumann, T., Nori, A., Alvarez-Valle, J., Poon, H., Oktay, O.: Making the most of text semantics to improve biomedical vision-language processing. In: Lecture Notes in Computer Science, pp. 1–21. Springer Nature Switzerland (2022). https://doi.org/10.1007/978-3-031-20059-5\_1
- Demner-Fushman, D., Kohli, M.D., Rosenman, M.B., Shooshan, S.E., Rodriguez, L., Antani, S., Thoma, G.R., McDonald, C.J.: Preparing a collection of radiology examinations for distribution and retrieval. J. Am. Med. Inform. Assoc. 23(2), 304–310 (Mar 2016)
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A largescale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
- Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. CoRR abs/1810.04805 (2018), http://arxiv.org/abs/1810.04805
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. CoRR abs/2010.11929 (2020), https://arxiv.org/abs/2010.11929
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. CoRR abs/1512.03385 (2015), http://arxiv.org/abs/1512.03385
- Huang, S.C., Shen, L., Lungren, M.P., Yeung, S.: 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 (ICCV). pp. 3942–3951 (October 2021)
- Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R.L., Shpanskaya, K.S., Seekins, J., Mong, D.A., Halabi,

S.S., Sandberg, J.K., Jones, R., Larson, D.B., Langlotz, C.P., Patel, B.N., Lungren, M.P., Ng, A.Y.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. CoRR abs/1901.07031 (2019), http://arxiv.org/abs/1901.07031

- Jang, J., Kyung, D., Kim, S.H., Lee, H., Bae, K., Choi, E.: Significantly improving zero-shot x-ray pathology classification via fine-tuning pre-trained image-text encoders (2022), https://arxiv.org/abs/2212.07050
- Jia, C., Yang, Y., Xia, Y., Chen, Y., Parekh, Z., Pham, H., Le, Q.V., Sung, Y., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. CoRR abs/2102.05918 (2021), https://arxiv.org/abs/ 2102.05918
- 12. Johnson, A., Pollard, T., Mark, R.: MIMIC-III clinical database (2020)
- Johnson, A.E.W., Pollard, T., Mark, R., Berkowitz, S., Horng, S.: The MIMIC-CXR database (2019)
- Li, Y., Liang, F., Zhao, L., Cui, Y., Ouyang, W., Shao, J., Yu, F., Yan, J.: Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. CoRR abs/2110.05208 (2021), https://arxiv.org/abs/2110.05208
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. CoRR abs/2103.14030 (2021), https://arxiv.org/abs/2103.14030
- Loshchilov, I., Hutter, F.: SGDR: stochastic gradient descent with restarts. CoRR abs/1608.03983 (2016), http://arxiv.org/abs/1608.03983
- 17. Loshchilov, I., Hutter, F.: Fixing weight decay regularization in adam. CoRR abs/1711.05101 (2017), http://arxiv.org/abs/1711.05101
- Mu, N., Kirillov, A., Wagner, D.A., Xie, S.: SLIP: self-supervision meets languageimage pre-training. CoRR abs/2112.12750 (2021), https://arxiv.org/abs/ 2112.12750
- Nguyen, H.Q., Lam, K., Le, L.T., Pham, H.H., Tran, D.Q., Nguyen, D.B., Le, D.D., Pham, C.M., Tong, H.T.T., Dinh, D.H., Do, C.D., Doan, L.T., Nguyen, C.N., Nguyen, B.T., Nguyen, Q.V., Hoang, A.D., Phan, H.N., Nguyen, A.T., Ho, P.H., Ngo, D.T., Nguyen, N.T., Nguyen, N.T., Dao, M., Vu, V.: Vindr-cxr: An open dataset of chest x-rays with radiologist's annotations. Scientific Data 9(1), 429 (Jul 2022), https://doi.org/10.1038/s41597-022-01498-w
- 20. Organization, W.H., et al.: Communicating radiation risks in paediatric imaging: information to support health care discussions about benefit and risk (2016)
- Pisano, E.D., Zong, S., Hemminger, B.M., DeLuca, M., Johnston, R.E., Muller, K., Braeuning, M.P., Pizer, S.M.: Contrast limited adaptive histogram equalization image processing to improve the detection of simulated spiculations in dense mammograms. Journal of Digital Imaging 11(4), 193 (Nov 1998), https: //doi.org/10.1007/BF03178082
- Qin, C., Yao, D., Shi, Y., Song, Z.: Computer-aided detection in chest radiography based on artificial intelligence: a survey. BioMedical Engineering OnLine 17(1), 113 (Aug 2018), https://doi.org/10.1186/s12938-018-0544-y
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. CoRR abs/2103.00020 (2021), https://arxiv.org/abs/2103.00020
- 24. Shih, G., wu, C., Halabi, S., Kohli, M., Prevedello, L., Cook, T., Sharma, A., Amorosa, J., Arteaga, V., Galperin-Aizenberg, M., Gill, R., Godoy, M., Hobbs, S., Jeudy, J., Laroia, A., Shah, P., Vummidi, D., Yaddanapudi, K., Stein, A.:

Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia. Radiology: Artificial Intelligence 1, e180041 (01 2019). https://doi.org/10.1148/ryai.2019180041

- 25. Vu, Y.N.T., Wang, R., Balachandar, N., Liu, C., Ng, A.Y., Rajpurkar, P.: Medaug: Contrastive learning leveraging patient metadata improves representations for chest x-ray interpretation. In: Jung, K., Yeung, S., Sendak, M., Sjoding, M., Ranganath, R. (eds.) Proceedings of the 6th Machine Learning for Healthcare Conference. Proceedings of Machine Learning Research, vol. 149, pp. 755–769. PMLR (06–07 Aug 2021), https://proceedings.mlr.press/v149/vu21a.html
- 26. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
- 27. Wang, Z., Wu, Z., Agarwal, D., Sun, J.: Medclip: Contrastive learning from unpaired medical images and text (2022), https://arxiv.org/abs/2210.10163
- Wei, J.W., Zou, K.: EDA: easy data augmentation techniques for boosting performance on text classification tasks. CoRR abs/1901.11196 (2019), http://arxiv.org/abs/1901.11196
- Yang, J., Li, C., Zhang, P., Xiao, B., Liu, C., Yuan, L., Gao, J.: Unified contrastive learning in image-text-label space (2022), https://arxiv.org/abs/2204.03610
- Zhou, H.Y., Chen, X., Zhang, Y., Luo, R., Wang, L., Yu, Y.: Generalized radiograph representation learning via cross-supervision between images and free-text radiology reports. Nature Machine Intelligence 4(1), 32-40 (jan 2022). https://doi.org/10.1038/s42256-021-00425-9, https://doi.org/10. 1038%2Fs42256-021-00425-9
- Zhou, H.Y., Lian, C., Wang, L., Yu, Y.: Advancing radiograph representation learning with masked record modeling. In: The Eleventh International Conference on Learning Representations (2023), https://openreview.net/forum?id= w-x7U26GM7j

# Supplementary Materials, CXR-CLIP: Toward Large Scale Chest X-ray Language-Image Pre-training

Kihyun You<sup>1</sup>, Jawook Gu<sup>1</sup>, Jiyeon Ham<sup>1</sup>, Beomhee Park<sup>1</sup>, Jiho Kim<sup>1</sup>, Eun K. Hong<sup>1</sup>, Woonhyuk Baek<sup>1</sup>, and Byungseok Roh<sup>1</sup>

Kakaobrain, Seongnam, Republic of Korea

#### 

**Table 1.** Comparison between our and GloRIA [8] prompt on cheXpert5x200 [9]. Performance gain of the model not trained with prompt (MIMIC-CXR [13]) suggests that our prompts also worth to evaluate. Training with our prompts further improved performance.

Madal Nama	Pre-train	GloRIA prompt	Our prompt
Model Name	Dataset	C5x200 ACC	C5x200 ACC
$\text{CXR-CLIP}_{SwinT}$	М	54.3	56.1
$\text{CXR-CLIP}_{SwinT}$	M,C	60.1	64.2
$\text{CXR-CLIP}_{SwinT}$	M,C,C14	62.8	65.7

**Table 2.** Comparison with self-supervised models (REFERS [30] and MRM [31]) in terms of classification tasks. We compared two-settings linear-probing and fine-tune whole visual backbone. All the models has ViT-base [6] backbone and are trained on MIMIC-CXR

Madal Nama	VinDR-CXR		R	SNA	SIIM	
Model Name	Linaer	Fine-tune	Linaer	Fine-tune	Linaer	Fine-tune
$\operatorname{REFERS}_{ViT-B}$	83.6	90.1	86.7	87.9	81.3	89.5
$MRM_{ViT-B}$	77.0	91.3	86.7	89.9	86.0	93.3
$ CXR-CLIP_{ViT-B} $	89.3	91.6	89.6	90.3	90.2	92.7

Table 3. Evaluation prompts for zero-shot classification. For VinDR-CXR and SIIM, we use simple prompt in Jang et el. [10], and we use prompt in BioVIL [2] for RSNA.

Dataset	Positive	Negative
VinDR-CXR, SIIM	{classname}	No {classname}
RNSA-pneumonia	Findings suggesting pneumonia.	No evidence of pneumonia.

13

**Table 4.** To compare BioVIL [2], we train our ResNet models with image resolution 512, denoted CXR-CLIP<sup>+</sup><sub>R50</sub>. RSUM is sum of recall@k, where  $k = \{1, 5, 10\}$ . Our models trained MIMIC outperforms BioVIL and CXR-CLIP<sup>+</sup><sub>R50</sub> generally outperforms CXR-CLIP<sup>+</sup><sub>R50</sub>

Model Name	Pre-train	VinDR	RSNA	SIIM	CheXper	t5x200	MIMIC	OpenI
Model Malle	Dataset	ZS(AUC)	ZS(AUC)	ZS(AUC)	ZS(ACC)	RSUM	RSUM	RSUM
$\operatorname{BioVIL}_{R50}$	М	-	83.1	-	-	-	-	-
$\text{CXR-CLIP}_{R50}$	Μ	78.8	83.3	85.2	54.0	65.0	126.6	25.3
$\text{CXR-CLIP}^+_{R50}$	Μ	82.2	84.8	85.2	56.8	64.6	133.7	27.5
$CXR-CLIP_{R50}$	M,C	83.0	85.0	86.4	61.7	52.1	124.3	23.7
$CXR-CLIP^+_{R50}$	$^{\rm M,C}$	87.5	85.3	89.0	57.2	50.1	117.6	22.6
CXR-CLIP <sub>R50</sub>	M,C,C14	78.1	81.8	85.2	60.3	52.0	120.5	19.1
$CXR-CLIP^+_{R50}$	M,C,C14	84.6	86.7	87.3	62.0	<b>59.6</b>	120.7	25.1

**Table 5.** Default positive and negative templates for suggested prompts, as well as class-specific templates. E is expressions for each class, + means text concatenation,  $[\cdot]$  means random selection from the given list, and () is blank text.

	Positive templates	Negative templates
Default	[ $\{E\}$ ., There is $\{E\}$ ., $\{E\}$ is [present, seen, noted]., the presence of $\{E\}$ is [seen, noted].]	<pre>[[There is, ()] + [no {E}., no radiographic evidence for {E}., no [visible, definite, obvious, appreciable, evident] {E}.</pre>
Edema Pneumonia	[Default Positive templates, Findings are + [suggesting, compatible with, suggestive of, representing] + $\{E\}$ .]	no [convincing, definite, ( )] evidence of $\{E\}$ ., no convincing signs of $\{E\}$ .], No $\{E\}$ is [visible, present, noted]. ]
Cardiomegaly	[heart size, cardiac size, cardiac silhouette, cardiac shadow, cardiac contour] + [is, appears] + [enlarged, increased].	[heart size, cardiac size, cardiac silhouette, cardiac shadow, cardiac contour] + [is, appears] + [normal, within normal limits, unremarkable].
Enlarged Cardio- mediastinum	[[cardiomediastinal, mediastinal] silhouette, [cardiomediastinum, mediastinum], mediastinal contour] + [is, appears] + [enlarged, widened].	[[cardiomediastinal, mediastinal] silhouette, [cardiomediastinum, mediastinum], mediastinal contour] + [is, appears] + [normal, within normal limits, unremarkable].
No Finding	[the lungs, both lungs, the lung fields, both lung fields] + [are clear, appear clear].	-

**Table 6.** Various expressions E for classes using default templates. +,  $[\cdot]$  and ( ) are same as Table 4.

Class Name	Expressions $E$
Atelectasis	[Atelectasis]
Consolidation	[Consolidation]
Edema	[Pulmonary edema]
Emphysema	[Emphysema, Emphysematous change]
Fibrogia	[[( ), pulmonary] + [( ), fibrotic] + [scar, scarring],
I'IDIOSIS	parenchymal + [scar, scarring], fibrotic change]
Fracture	[Fracture, Acute fracture]
Hernia	[Hernia, Herniation, Hiatal Hernia]
Infiltration	[[( ), pulmonary] + infiltration, infiltrate,
	infiltrative + [density, opacity, process]]
Lung Lesion	Pos: [lung lesion]
Lung Lesion	Neg: $[[lung, pulmonary] + [nodule, mass, lesions, nodules or masses]]$
Lung Opacity	[pulmonary opacity]
Mass	[[pulmonary, lung] + mass]
Nodule	[[(), pulmonary] + [nodule, nodular opacity, nodular density]]
Pleural Effusion	[Pleural Effusion]
Pleural Other	[Pleural Abnormality]
Pleural Thickening	[Pleural Thickening, Thickened pleura]
Pneumonia	[Pneumonia]
Pneumothorax	[Pneumothorax]
Support Devices	[Support Devices]