

# 1 Generalized Radiograph Representation Learning 2 via Cross-supervision between Images and 3 Free-text Radiology Reports

4 Hong-Yu Zhou<sup>1,†</sup>, Xiaoyu Chen<sup>2,†</sup>, Yinghao Zhang<sup>2,†</sup>, Ruibang Luo<sup>1</sup>,  
Liansheng Wang<sup>2,\*</sup>, Yizhou Yu<sup>1,\*</sup>

<sup>1</sup> Department of Computer Science, The University of Hong Kong,  
Pokfulam, Hong Kong

<sup>2</sup> Department of Computer Science, Xiamen University,  
Xiamen, China

<sup>†</sup> These authors contributed equally

\* Corresponding authors: L.Wang (lswang@xmu.edu.cn) and  
Y.Yu (yizhouy@acm.org)

## 5 Abstract

6 Pre-training lays the foundation for recent successes in radiograph analysis sup-  
7 ported by deep learning. It learns transferable image representations by conducting  
8 large-scale fully-supervised or self-supervised learning on a source domain. However,  
9 supervised pre-training requires a complex and labor intensive two-stage human-assisted  
10 annotation process while self-supervised learning cannot compete with the supervised  
11 paradigm. To tackle these issues, we propose a cross-supervised methodology named  
12 REviewing FreE-text Reports for Supervision (REFERS), which acquires free supervi-  
13 sion signals from original radiology reports accompanying the radiographs. The pro-  
14 posed approach employs a vision transformer and is designed to learn joint representa-  
15 tions from multiple views within every patient study. REFERS outperforms its transfer  
16 learning and self-supervised learning counterparts on 4 well-known X-ray datasets under  
17 extremely limited supervision. Moreover, REFERS even surpasses methods based on  
18 a source domain of radiographs with human-assisted structured labels. Thus REFERS  
19 has the potential to replace canonical pre-training methodologies.

## 20 1 Introduction

21 Medical image analysis has achieved tremendous progress in recent years, thanks to the  
22 development of deep convolutional neural networks (DCNNs) [1, 2, 3, 4, 5]. At the core of  
23 DCNNs is visual representation learning [6], where pre-training has been widely adopted  
24 and become the most dominant approach to obtain transferable representations. Typically,

25 a large-scale dataset, also called the source domain, is first used for model pre-training.  
26 Transferable representations from the pre-trained model are further fine-tuned on other  
27 smaller downstream datasets, called target domains.

28 As one of the most general forms of medical images, radiographs have a great po-  
29 tential to be used in widespread applications [7, 8, 9]. In order to achieve (or at least  
30 approximate) radiologist-level diagnosis performance in these applications, it is common  
31 to transfer learned representations from natural images to radiographs [10, 11], and Im-  
32 geNet [12] based pre-training is most widely adopted in this context. On the other hand,  
33 self-supervised learning [13, 14, 15, 16] has attracted much attention in the community be-  
34 cause it is capable of learning transferable radiograph representations without any human  
35 annotations. Both methodologies have been proven to be effective in solving medical image  
36 analysis tasks, especially when the amount of labeled data in the target domain is quite  
37 limited. However, in the first approach, there is an inevitable problem, which is the exis-  
38 tence of domain shifts between medical and natural images. For instance, it is possible to  
39 introduce harmful noises from natural images as radiographs have a different pixel intensity  
40 distribution. As for self-supervised learning, to the best of our knowledge, there still exist  
41 clear performance gaps between radiograph representations learned through self-supervised  
42 and label-supervised pre-training. To avoid these problems, building large-scale annotated  
43 radiograph datasets for label-supervised pre-training becomes an essential and urgent issue  
44 in radiograph analysis.

45 Recently, radiologists and computer scientists have managed to build medical datasets  
46 for label-supervised pre-training at the size of hundreds of thousands of images, such as  
47 ChestX-ray [11], MIMIC [17] and CheXpert [18]. To acquire accurate labels for radiographs,  
48 these datasets often rely on a two-stage human intervention process. A radiology report is  
49 first prepared by radiologists for every patient study as part of the clinical routine. In the  
50 second stage, human annotators extract and confirm structured labels from these reports  
51 using artificial rules and existing natural language processing (NLP) tools. However, there  
52 are two major limitations of this label extraction workflow. First, it is still complex and  
53 labor intensive. For example, human annotators have to define a list of alternate spellings,  
54 synonyms, and abbreviations for every target label. Consequently, the final accuracy of  
55 extracted labels heavily depends on the quality of human assistance and various NLP tools.

56 A small mistake in a single step or a single tool may give rise to disastrous annotation  
57 results. Second, those human-defined rules are often severely restricted to application-  
58 oriented tasks instead of general-purpose tasks. It is difficult for DCNNs to learn universal  
59 representations from such application-oriented tasks.

60 In this paper, we propose **RE**viewing **FreE**-text **R**eports for **S**upervision (REFERS)  
61 to directly learn radiograph representations from accompanying free-text radiology reports.  
62 We believe abstract and complex logic reasoning sentences in radiology reports provide  
63 sufficient information for learning well-transferable visual features. As shown in Figure  
64 1a, REFERS is realized using a set of transformers, where the most important part is  
65 a radiograph transformer serving as the backbone. The main reason why we choose the  
66 transformer as the backbone in REFERS is that it not only exhibits the advantages of  
67 DCNNs, but also has been shown to be more effective [19] because of the self-attention  
68 mechanism [20]. Moreover, we have found that, in comparison to features generated from  
69 DCNNs, features from transformers are more compatible with textual tasks.

70 Different from aforementioned representation learning methodologies, REFERS per-  
71 forms cross-supervised learning and does not need structured labels during the pre-training  
72 stage. Instead, supervision signals are defined by automatically cross-checking the two  
73 different data modalities, radiographs and free-text reports. Considering in daily clinical  
74 routine, there is typically a free-text report associated with every patient study, which  
75 usually involves more than one radiographs. To fully utilize the study-level information in  
76 each report, we design a view fusion module based on an attention mechanism to process all  
77 radiographs in a patient study simultaneously, and fuse the resulting multiple features. In  
78 this way, the learned representations are able to preserve both study-level and image-level  
79 information. In contrast, only image-level information is addressed in traditional represen-  
80 tation learning paradigms [11, 13, 14, 15, 16] that use a single image as input. On top  
81 of the view fusion module, we conduct two tasks, i.e., report generation and study-report  
82 representation consistency reinforcement, to extract study-level supervision signals from  
83 free-text reports. To carry out the first task, we apply a decoder, called report transformer,  
84 to the fused feature with the goal to reproduce the radiology report associated with the  
85 study. For the second task, we apply our radiograph transformer and an NLP transformer  
86 to a study-report pair. These transformers produce a pair of feature representations for the

87 patient study and radiology report in the pair, respectively. The consistency between such a  
88 pair of feature representations within every study-report pair is reinforced via a contrastive  
89 loss function. Some previous works [21, 22] tried to learn joint text-image representations  
90 for single-domain medical image analysis tasks. Compared to them, REFERS focuses on  
91 learning well-transferable image features from study-level free-text reports on a large-scale  
92 source domain and fine-tuning them on one or more target domains.

93 On four well-known X-ray datasets, REFERS outperforms self-supervised learning and  
94 transfer learning on natural source images in producing more transferable representations,  
95 often bringing impressive improvements (more than 5%) under limited supervision from  
96 target domains. This capability can be extremely important in real-world applications as  
97 medical data is scarce and their annotations are usually hard to acquire. More surprisingly,  
98 we found that REFERS clearly surpasses those methods that employ a source domain with  
99 a large collection of medical images with structured labels. In terms of specific abnormal-  
100 ities and diseases, REFERS is quite effective under extremely limited supervision (< 1k  
101 annotated radiographs during fine-tuning). For instance, REFERS brings about 9-percent  
102 improvements on pneumothorax. Meanwhile, over 7-percent improvements are achieved on  
103 two common lung diseases (atelectasis and emphysema).

## 104 **2 Results**

105 All self-supervised learning (SSL) and label-supervised pre-training (LSP) baselines as well  
106 as our REFERS are first pre-trained on a source domain of medical images (i.e., MIMIC-  
107 CXR-JPG [23]). Then, pre-trained models are fine-tuned on each of four well-established  
108 datasets (target domains with labels), including NIH ChestX-ray [11], VinBigData Chest  
109 X-ray Abnormalities Detection [24], Shenzhen Tuberculosis [25] and COVID-19 Image Data  
110 Collection [26]. During the fine-tuning stage, we always perform fully-supervised learning  
111 on the target domain, which only consists of radiographs with structured labels. Further-  
112 more, we verify model performance by varying the percentage of actually used training  
113 images (sampled from the predefined whole training set) in the target domain, and this  
114 percentage is called *label ratio*. When the label ratio is 100%, we use the whole training set  
115 in the target domain for fine-tuning.

116

117 **NIH ChestX-ray.** Table 1, Supplementary Figures 1a and 2a present experimental results  
118 from our REFERS and other approaches under different label ratios. As shown in Table 1  
119 and Supplementary Figure 1a, our approach significantly outperforms self-supervised base-  
120 lines and transfer learning on natural source images. To be specific, REFERS achieves  
121 the highest AUC on all 14 classes using different amounts of training data during the  
122 fine-tuning stage. Moreover, REFERS shows the largest performance improvements with  
123 respect to these baselines when only 0.8k training images (1% label ratio) in the target  
124 domain are utilized. For example, REFERS surpasses the widely adopted ImageNet-based  
125 pre-training [11] by about 7 percents on average. Even when compared to LSP, our REFERS  
126 still gives quite competitive results. In Table 2, it is easy to find out that the average perfor-  
127 mance of REFERS actually surpasses LSP, and consistently maintains an advantage of at  
128 least 2 percents. Compared to self-supervised baselines [13, 14, 15, 16] and ImageNet-based  
129 pre-training [11], REFERS achieves the largest improvements on emphysema (7 percents)  
130 and cardiomegaly ( $> 10$  percents), especially under limited supervision. When compared  
131 to LSP, our method achieves consistent improvements on mass ( $> 4$  percents).

132

133 **VinBigData Chest X-ray Abnormalities Detection.** Our REFERS exhibits more  
134 advantage on this target domain dataset than it does on NIH ChestX-ray as VinBigData  
135 comprises a much smaller number of annotated radiographs (about  $\frac{1}{8}$  of the NIH dataset).  
136 This phenomenon again demonstrates the ability of REFERS in dealing with limited su-  
137 pervision. REFERS consistently maintains large advantages over other methods under dif-  
138 ferent conditions (see Tables 1, 2, Supplementary Figures 1b and 2b). For instance, when  
139 we only have 105 annotated radiographs (1% label ratio) as fine-tuning data, REFERS sur-  
140 passes C2L [16], the best performing self-supervised method, by over 7 percents in AUC.  
141 The performance of REFERS once again surpasses LSP with human-assisted structured  
142 labels even when all annotated training data (100% label ratio) in the target domain is  
143 used. When we check specific abnormalities and diseases, we found REFERS consistently  
144 improves the diagnosis of atelectasis, lung opacity and pneumothorax in comparison to LSP.

145

146 **COVID-19 and Shenzhen Tuberculosis Image Collections** Both datasets serve as

147 target domains and comprise a small number of labeled images (fewer than 1k X-rays),  
148 which are employed to test the transferability of the representation learned on the source  
149 domain. This is because few training images in such small target domains are not capable of  
150 training powerful models themselves. Thus, the performance of the trained models is more  
151 dependent on the quality of the learned representation. In Table 1, although separating  
152 tuberculosis from normal cases is not a hard task, our method still achieves 2.5% improve-  
153 ments over C2L [16] in AUC. When looking at COVID-19 Image Data Collection which  
154 includes two harder tasks, we can find that the relative performance improvements over  
155 self-supervised baselines [13, 14, 15, 16] and transfer learning on natural source images [11]  
156 become quite clear. For instance, on the “Viral vs. Bacterial” task, REFERS outperforms  
157 C2L [16] by 7 percents in AUC, demonstrating the effectiveness of REFERS in helping  
158 achieve better performance over small-scale target datasets. Even if we compare REFERS  
159 against LSP, the performance advantage is still maintained at more than 1 percent.

160

### 161 **3 Discussion**

162 **REFERS outperforms self-supervised learning and transfer learning on natural**  
163 **source images by substantial and significant margins.** This is the most promi-  
164 nent observation obtained from our experimental results, which holds on different datasets  
165 and with different amounts of annotated training data during fine-tuning. Among self-  
166 supervised baselines [13, 14, 15, 16], C2L [16] and TransVW [15] are the two best per-  
167 forming methods. Our REFERS outperforms C2L and TransVW by at least 4 percents  
168 when very limited annotated training data (at most 10% label ratio) from NIH ChestX-  
169 ray and VinBigData datasets is used. Somewhat interestingly, as the label ratio increases,  
170 ImageNet-based pre-training [11] gradually narrows its gap with self-supervised learning.  
171 Nonetheless, our REFERS still surpasses it by a large margin (4 percents at least). Similar  
172 results can also be observed on Shenzhen Tuberculosis and COVID Image Collection. Since  
173 our REFERS employs a cross-supervised learning manner, it does not require structured  
174 labels as conventional fully-supervised learning approaches. As radiographs and radiology  
175 reports are readily available medical data, we believe our approach is as practical as self-

176 supervised learning methodologies in real-world scenarios.

177

178 **REFERS consistently surpasses label-supervised pre-training with human-assisted**

179 **structured labels.** This is another clear observation obtained from our experimental

180 results. Even though our approach does not use any structured labels in the source do-

181 main, over all four target domain datasets, our pre-trained model exhibits clear advantages.

182 Specifically, REFERS outperforms the most competitive LSP method, LSP (Transformer),

183 which is based on Transformer and human-assisted structured labels in the source domain.

184 In particular, our method shows more advantages at small label ratios. For instance, when

185 NIH ChestX-ray and VinBigData are used as target domain datasets, REFERS achieves

186 about 2.5% improvements when the number of training images is smaller than 10k. Sim-

187 ilarly, on Shenzhen Tuberculosis and COVID-19 Data Collection, REFERS consistently

188 surpasses LSP by significant margins. It is worth mentioning that when a classification

189 problem is difficult to solve and has limited supervision, REFERS becomes more advan-

190 tageous and achieves impressive improvements. For example, on the “Viral vs. Bacterial”

191 task (Table 2), REFERS surpasses label-supervised pre-training methods based on two-

192 stage human intervention by approximately 4 percents. These improvements demonstrate

193 that raw radiology reports contain more useful information than human-assisted structured

194 labels. In other words, the advantages exhibited by our approach on small-scale target

195 domain training data can be attributed to the rich information carried by radiology re-

196 ports in the source domain. Such information provides additional supervision to help learn

197 transferable representations for radiographs while the supervision signals from structured

198 labels have less information. We believe this is an important step towards directly using

199 natural language descriptions as supervision signals for image representation learning. As

200 an example, our REFERS can be used to learn natural image representations from text

201 descriptions at corresponding websites.

202

203 **REFERS significantly reduces the need of annotated data in target domains.**

204 Figures 2a and 2b present the performance of our approach under various label ratios. On

205 NIH ChestX-ray, REFERS needs 90% fewer annotated target domain data (10% label ratio)

206 to deliver a performance comparable to those of Model Genesis [14] and ImageNet-based

207 pre-training [11]. Similarly, on VinBigData, our method only needs 10% annotated train-  
208 ing data to achieve much better results than those of Model Genesis and ImageNet-based  
209 pre-training under 100% label ratio. This phenomenon shows the potential of REFERS in  
210 providing high-quality pre-trained representations for downstream fine-tuning tasks with  
211 limited annotations. Due to the difficulty to acquire reliable annotations for medical image  
212 analysis, the ability to achieve good performance with limited annotations means much to  
213 the community.

214

215 **Improvements on specific abnormalities and diseases.** In Supplementary Figures 1  
216 and 2, REFERS brings 5-percent performance gains on emphysema and mass even when  
217 compared to LSP with limited supervision in the target domain ( $< 10k$  training images).  
218 Since both abnormalities have a dispersed spatial distribution in the lung area, the consider-  
219 able improvements demonstrate that REFERS is able to handle elusive chest abnormalities  
220 in radiographs well. When the amount of supervision in the target domain becomes ex-  
221 tremely limited, such as using 105 training images from VinBigData, REFERS becomes  
222 more advantageous. For instance, REFERS outperforms LSP on atelectasis and pneumotho-  
223 rax by over 7 and 9 percents, respectively. Different from emphysema, mass and atelectasis,  
224 pneumothorax maintains a concentrated spatial distribution and is often located around  
225 the pleura. These successes imply that REFERS can deal with the diagnosis of both elusive  
226 and regular abnormalities and diseases well using a small number of training radiographs  
227 in the target domain. A similar phenomenon can be observed when REFERS is used for  
228 distinguishing viral pneumonia cases from bacterial ones in Tables 1 and 2.

229

230 **Transformer is more effective under limited supervision.** In Tables 1 and 2, we  
231 observe a trend of CNNs (i.e., ResNet series [4]): LSP (ConvNet) shows mediocre perfor-  
232 mance when a relatively small number of training images in the target domain are used.  
233 However, when all training data (100% label ratio) is used, ConvNet shows competitive  
234 results. It seems that LSP (ConvNet) cannot well handle little amount of supervision. In  
235 contrast, LSP (Transformer) exhibits much better performance at small label ratios. This  
236 comparison demonstrates that pre-trained transformers generate more transferable repre-  
237 sentations than pre-trained CNNs. The underlying reason might be that the self-attention



238 mechanism in transformers makes the learned representations more transferable due to cap-  
239 tured long-distance dependencies.

240

241 **REFERS provides reliable evidences for clinical decisions.** Figure 3 presents ran-  
242 domly chosen radiographs and their corresponding class activation maps (CAMs) [27]. We  
243 can find that REFERS generates reliable attention regions, on top of which we can apply  
244 a fixed confidence threshold to further identify the location of different types of lesions  
245 (green boxes in Figure 3). The overall IoUs (Intersection over Unions) between green and  
246 red boxes (drawn by radiologists) are mostly higher than 0.5, indicating that the generated  
247 attention regions can well match radiologists’ diagnoses. When lesions have a large size  
248 (such as the fifth image from NIH ChestX-ray), our method captures well-aligned lesion  
249 areas. Even when lesions are quite small and thus hard to detect (such as the last image  
250 from NIH ChestX-ray and the first image from VinBigData), REFERS can still identify  
251 the right locations.

252

253 **Replication of experimental results and their statistical significance.** There are a  
254 number of factors that influence pre-training results exhibit a certain level of randomness.  
255 These factors include, but are not limited to network initialization, training strategy (e.g.,  
256 how to randomly crop images and perform mini-batch gradient descent) and even non-  
257 deterministic characteristics in computational tools (e.g., cuDNN [28] would choose differ-  
258 ent algorithms in different runs due to benchmarking noise and hardware configuration). A  
259 good pre-training methodology should be able to produce relatively stable pre-trained rep-  
260 resentations when randomness in these factors is controlled within an acceptable limit. To  
261 take into account the influence of such randomness on experimental results, when REFERS  
262 and baseline pre-trained models are fine-tuned, we independently repeat each experiment  
263 three times and report their average results in Tables 1 and 2. Then, we calculate p-values  
264 between mean class AUCs of our REFERS and the best performing baseline model ac-  
265 cording to their fine-tuned performance using independent two-sample t-test. According to  
266 Tables 1 and 2, nearly all p-values are much smaller than 0.01, indicating that our REFERS  
267 is significantly better than its counterparts when various amounts of labeled training data  
268 in the target domain is used. In contrast, making the number of times (repeating each

269 experiment) smaller than three would give rise to less stable mean AUCs while simply  
270 repeating more times would produce meaninglessly smaller p-values.

## 271 **4 Methods**

272 **Dataset for pre-training (source domain).** MIMIC-CXR-JPG [23] contains over 370k  
273 radiographs organized into patient studies, each of which may have one or more radiographs  
274 taken from different views or at different times for the same patient. Each patient study  
275 has one free-text radiology report, and each radiograph is associated with a set of abnor-  
276 mality/disease labels obtained from two-stage human-assisted intervention as mentioned  
277 above. There are two major sections in each report: Findings and Impressions. The Find-  
278 ings section includes detailed descriptions of important aspects in the radiographs while  
279 the Impressions section summarizes most immediately relevant findings.

280 To acquire human-assisted structured labels for radiographs (i.e., two-stage human in-  
281 tervention), annotators need to first define a list of labels for abnormalities and diseases,  
282 including alternate spellings, synonyms, and abbreviations. On the basis of local contexts  
283 and existing NLP tools, mentions of labels in reports are classified as positive, uncertain,  
284 or negative. An aggregation procedure is further applied to aggregate multiple mentions of  
285 a single label. Uncertain labels need to be double-checked by radiologists.

286 As radiology reports were originally prepared by radiologists as part of the daily clinical  
287 routine, they can be regarded as freely available information that does not require extra  
288 human efforts in contrast to structured labels. In practice, we only keep the Findings and  
289 Impressions sections in the reports. Also, we remove all study-report pairs, where the text  
290 section has less than 3 tokens (words and phrases), from the dataset. This screening pro-  
291 cedure produces 217k patient studies.

292

293 **Datasets for fine-tuning (target domains).** We do not require these datasets adopted  
294 for fine-tuning to have radiology reports. Instead, only human-assisted annotations are  
295 used during the fine-tuning stage. We follow the official split of NIH ChestX-ray, where the  
296 percentages of training, validation and testing sets are 70%, 10% and 20%, respectively. The  
297 same set of ratios are also employed for VinBigData Chest X-ray, Shenzhen Tuberculosis

298 and COVID-19 Image Data Collection to build randomly split training, validation and  
299 testing sets.

- 300 • *NIH ChestX-ray* is a dataset for multi-label classification of 14 chest abnormalities  
301 (i.e., Atelectasis, Cardiomegaly, Consolidation, Edema, Effusion, Emphysema, Fibro-  
302 sis, Hernia, Infiltration, Mass, Nodule, Pleural Thickening, Pneumonia and Pneu-  
303 mothorax). There are over 100k frontal-view X-ray images of about 32k patients in  
304 NIH ChestX-ray, where labels of radiographs were extracted from associated reports  
305 following a similar procedure as that for MIMIC-CXR-JPG.
- 306 • *VinBigData Chest X-ray* provides labels of 14 chest diseases (i.e., Aortic enlargement,  
307 Atelectasis, Pneumothorax, Lung Opacity, Pleural thickening, ILD, Pulmonary fibro-  
308 sis, Calcification, Pleural effusion, Consolidation, Cardiomegaly, Other lesion, Nodule-  
309 Mass and Infiltration), and consists of 15k postero-anterior chest X-ray images. Here  
310 we did not use the test set in Kaggle, which does not provide any annotations. All  
311 images were labeled by a panel of experienced radiologists.
- 312 • *Shenzhen Tuberculosis* is a small dataset containing 662 frontal chest X-ray images  
313 primarily from hospital clinical routine. 336 abnormal X-rays show various manifes-  
314 tations of tuberculosis, and the remaining 326 images are normal. We simply perform  
315 binary classification on this dataset.
- 316 • *COVID-19 Image Data Collection* is a dataset involving more than 900 pneumonia  
317 cases with chest X-rays, which was built to improve the identification of COVID-19.  
318 We conduct experiments on two tasks, which are a) distinguishing COVID-19 from  
319 the rest and b) separating viral pneumonia cases from bacterial ones.

320 **Baselines and label-supervised pre-training.** Since our method does not need struc-  
321 tured labels required by traditional fully-supervised learning, we compare it against four re-  
322 cent self-supervised learning methods [13, 14, 15, 16] and ImageNet-based pre-training [11]:

- 323 • *Context Restoration* [13] repeats the operation of swapping two randomly chosen small  
324 X-ray patches for a fixed number of times, and the neural network is asked to restore  
325 each altered image back to its original version.

- 326 • *Model Genesis* [14] applies multiple types of distortions to the input X-ray, including  
327 local shuffling, non-linear transformation, in- and out-painting. Similar to Context  
328 Restoration, Model Genesis asks the model to reconstruct the original image from the  
329 distorted one.
- 330 • *TransVW* [15] contrasts local X-ray patches to exploit the semantics of anatomical  
331 patterns while restoring distorted image contents.
- 332 • *C2L* [16] proposes to construct homogeneous and heterogeneous data pairs by mixing  
333 both images and features on top of MoCo [29]. C2L outperforms MoCo by observable  
334 margins on multiple X-ray benchmarks.
- 335 • *ImageNet-based pre-training* [11] is taken as a representative method that sets a large-  
336 scale dataset of annotated natural images as the source domain.

337 Note that all above baselines are implemented using the same transformer-based network  
338 architecture as our REFERS (i.e, a ViT architecture plus the proposed recurrent concate-  
339 nation module). Such an implementation arrangement is meant to rule out the influence of  
340 network architectures on final performance and maintain fairness in experimental compar-  
341 isons.

342 Finally, our approach is compared against label-supervised pre-training (LSP) that di-  
343 rectly sets a large collection of X-ray images with human-assisted structured labels as  
344 the source domain. For better comparison, we implement LSP on top of both CNN and  
345 Transformer based backbone networks. Specifically, LSP (Transformer) adopts the same  
346 Transformer based network architecture as REFERS and the aforementioned self-supervised  
347 and ImageNet-based pre-training baselines. LSP (ConvNet) stands for the best performing  
348 residual network among ResNet-18, ResNet-50 and ResNet-101 [4].

349

350 **Data augmentation and image resizing.** During the *pre-training* stage, we resize each  
351 radiograph in the source domain to  $256 \times 256$  pixels, and then apply random cropping to  
352 produce  $224 \times 224$  images. Random horizontal flip, random rotation (-10 to 10 degrees)  
353 and random grayscale (brightness and contrast) are also applied to generate augmented  
354 images. When using random horizontal flip, we change the words ‘left’ and ‘right’ in the  
355 accompanying radiology report accordingly. During the *fine-tuning* stage, we apply the

356 same set of data augmentation strategies, which are random cropping, random rotation,  
357 random grayscale and random horizontal flip, to all four target domain datasets. As in  
358 the pre-training stage, we resize each radiograph in a target domain to  $256 \times 256$ , and then  
359 generate  $224 \times 224$  cropped and augmented radiographs as input images.

360

361 **Algorithm Overview.** REFERS performs cross-supervised learning on top of a trans-  
362 former based backbone, called radiograph transformer. Given a patient study, we first  
363 forward its views to the radiograph transformer for extracting view-dependent feature rep-  
364 resentations. Next, we perform cross-supervised learning that acquires study-level supervi-  
365 sion signals from free-text radiology reports. To this aim, it is necessary and essential to  
366 use view fusion to obtain a unified visual representation for an entire patient study because  
367 each radiology report is associated with a patient study but not individual radiographs  
368 within the patient study. Such fused representations are then used in two tasks during  
369 the pre-training stage: report generation and study-report representation consistency rein-  
370 forcement. The first task takes the free texts in original radiology reports to supervise the  
371 training process of the radiograph transformer. The second task reinforces the consistency  
372 between the visual representations of patient studies and the textual representations of  
373 their corresponding reports.

## 374 **4.1 Radiograph Transformer**

375 The radiograph transformer accepts image patches as inputs. We divide each image into a  
376 grid of  $14 \times 14$  cells, each of which has  $16 \times 16$  pixels. We then flatten each image patch to  
377 form a 1D vector of pixels, and feed it to the transformer. At the beginning of the trans-  
378 former, a patch embedding layer linearly transforms each 1D pixel vector into a feature  
379 vector. This vector is concatenated with a position feature produced from a learnable posi-  
380 tion embedding to help clarify the relative location of each patch in the whole input patch  
381 sequence. The concatenated feature is then passed through another linear transformation  
382 layer to make its dimensionality the same as that of the final radiograph feature. At the  
383 core part of the radiograph transformer, we stack twelve self-attention blocks, which have  
384 the same architecture but independent parameters (Figure 1b). We first follow the practice  
385 in [20] to build a single self-attention block and then repeat its operations multiple times.

386 In each block, we apply layer normalization [30] before the multi-head attention and per-  
387 ceptron layers, after which residual connections are added to stabilize the training process.  
388 In the perceptron layer, we employ a two-layer perceptron with the Rectified Linear Unit  
389 (ReLU) [31] as the activation function. Moreover, we add an aggregation embedding, which  
390 is responsible for gathering the information from different input features. As shown in Fig-  
391 ure 1b, in the last layer, recurrent concatenation is performed to repeatedly concatenates  
392 the learned aggregation embedding with the learned representation of every patch. This is  
393 different from the operation in vision transformer (ViT) [19], which only concatenates the  
394 aggregation embedding with patch features once.

## 395 **4.2 Cross-supervised Learning**

396 There are two major components in cross-supervised learning: the view fusion module for  
397 producing study-level representations and two report-related tasks exploiting study-level  
398 information from associated free-text reports.

399 As aforementioned, we forward all radiographs in a patient study through the radio-  
400 graph transformer simultaneously to obtain their individual representations. We further  
401 employ an attention mechanism to fuse these individual representations to obtain an over-  
402 all representation of the given study. Supposing a study has three radiographs (i.e., views),  
403 as shown in Figure 1c. We first concatenate the features of all views, and then feed the  
404 concatenated features to a multi-layer perceptron to compute an attention value for each  
405 view. Next, we apply the softmax function to normalize these attention values, which are  
406 used as weights to produce a weighted version of the individual representations. Finally,  
407 these weighted representations are concatenated to form a unified visual feature for de-  
408 scribing the whole study. Note that for studies that contain few than three radiographs,  
409 we randomly select one of the radiographs, and then repeat it once or twice to have a total  
410 of three views. For studies that contain more than three radiographs, we randomly select  
411 three of them from each study as input views.

412 We design two report-related tasks that acquire cross-supervision signals from free-text  
413 reports: report generation and study-report representation consistency reinforcement. In  
414 practice, these two tasks exploit study-level free-text information for better training study-  
415 level visual representations produced from the view fusion module. The first task applies

416 a decoder, called report transformer, to the unified visual feature  $\mathbf{v}^k$  of the  $k$ -th patient  
417 study to reproduce its associated radiology report denoted as  $c_{1:T}^k$ . Here,  $c_1^k$  represents  
418 the start-of-sequence token and  $c_T^k$  the end-of-sequence token. As a result, the report  
419 transformer generates a sequence of token-level predictions,  $\hat{c}_{1:T}^k$ , for the  $k$ -th patient study.  
420 The prediction of the  $t$ -th token in this sequence depends on the predicted subsequence  $\hat{c}_{1:t-1}^k$   
421 and the visual feature  $\mathbf{v}^k$ . The network architecture of the report transformer follows the  
422 architecture of the decoder in [20]. We wish the predicted token sequence ( $\hat{c}_{1:T}^k$ ) resembles  
423 the sequence ( $c_{1:T}^k$ ) representing the original report of the  $k$ -th patient study. Therefore, as  
424 shown in Figure 1d, we apply a language modeling loss to both  $\hat{c}_{1:T}^k$  and  $c_{1:T}^k$  to maximize  
425 the following log-likelihood of the tokens in the original report.

$$\mathcal{L}_{\text{language}}^k = \sum_{t=2}^T \log P \left( c_t^k \mid \hat{c}_{1:t-1}^k, \mathbf{v}^k; \phi_v, \phi_t \right), \quad (1)$$

426 where  $\hat{c}_1^k$  is a special symbol indicating the start of the predicted sequence,  $\phi_v$  and  $\phi_t$  stand  
427 for the parameters of the radiograph transformer and report transformer, respectively.

428 For the second task on study-report representation consistency reinforcement, we employ  
429 a contrastive loss [32] to align cross-modal representations. Here, we use  $\mathbf{t}^k$  to stand for the  
430 textual feature vector of the  $k$ -th radiology report. In practice, we obtain  $\mathbf{t}_k$  by forwarding  
431 the sequence of tokens in the  $k$ -th report (i.e.,  $c_{1:T}^k$ ) to a BERT (i.e., Bidirectional Encoder  
432 Representations from Transformer) model [33]. BERT is built on top of the encoder in  
433 [20] using large-scale pre-training on a great number of corpus resources. Thus, BERT can  
434 help produce a generalized textual representation for the input report. Suppose we have  $B$   
435 patient studies in each training mini-batch, as shown in Figure 1d. The contrastive loss for  
436 the  $k$ -th study can be formulated as

$$\mathcal{L}_{\text{contrast}}^k = -\log \frac{e^{\cos(\mathbf{v}^k, \mathbf{t}^k)/\tau}}{\sum_{i=1}^B e^{\cos(\mathbf{v}^k, \mathbf{t}^i)/\tau}}, \quad (2)$$

437 where  $\cos(\cdot, \cdot)$  means the cosine similarity,  $\cos(\mathbf{v}^k, \mathbf{t}^k) = \frac{(\mathbf{v}^k)^\top \mathbf{t}^k}{\|\mathbf{v}^k\| \|\mathbf{t}^k\|}$ ,  $\top$  denotes the transpose  
438 operation,  $\|\cdot\|$  stands for L2 normalization, and  $\tau$  is the temperature factor. Finally, for  
439 each patient study, we simply sum up  $\mathcal{L}_{\text{contrast}}^k$  and  $\mathcal{L}_{\text{language}}^k$  as the overall loss. During the  
440 *fine-tuning* stage, we typically use the cross entropy loss for model tuning.

441

442 **Training and testing methodologies.** We first pre-train the radiograph transformer

443 on the source domain and then fine-tune it on downstream target domain datasets to  
444 verify the quality of pre-training. During the *pre-training* stage, we sample 4.6k studies to  
445 form a held-out validation set according to the official division of the MIMIC-CXR-JPG  
446 dataset [23]. We train the entire network using stochastic gradient descent (SGD) while  
447 setting the momentum value to 0.9 [34] and the weight decay to 1e-4. Following [33], we  
448 do not apply weight decay to layer normalization and the bias terms in all layers. We use  
449 a fixed batch size of 32 for 300k iterations (about 45 epochs). We calculate the validation  
450 loss after each epoch and save the checkpoint that achieves the lowest validation loss. We  
451 adopt the linear learning rate warm-up strategy [35] for the first 10k iterations, and then  
452 switch to cosine decay [36] until the end. Empirically, we found that training the radiograph  
453 transformer requires a large learning rate for fast convergence. Thus, its learning rate is  
454 set to 3e-3 while the learning rate for the report transformer and BERT is set to 3e-4.  
455 We initialize the aggregation embedding to all zeros while randomly initializing all position  
456 embeddings. We use PyTorch [37] and NVIDIA Apex for mixed-precision training [38].  
457 The complete pre-training process on the MIMIC-CXR dataset takes about 2 days on a  
458 single RTX 3090 GPU.

459 During the *fine-tuning* stage, we fine-tune all transformer based models (including trans-  
460 former based baselines) using SGD with the momentum set to 0.9 and the initial learning  
461 rate set to 3e-3 for all datasets. We fine-tune ResNet models using Adam [39] instead of  
462 SGD, and set the initial learning rate to 1e-4. All downstream models use the same learning  
463 rate decay strategy as that used in the pre-training stage, and are trained with a batch size  
464 of 128.

### 465 **4.3 Ablation Study**

466 We conduct a thorough ablation study of REFERS by removing or replacing individual  
467 modules, and the results are shown in Table 3.

468 First, we investigate the impact of replacing the radiograph transformer (rows 1-2 in  
469 Table 3). If we replace the radiograph transformer with ResNet-101 [4] (row 1), the overall  
470 performance of REFERS on COVID-19 Image Data Collection would drop by about 7 per-  
471 cents (compared to row 0). This comparison demonstrates that the radiograph transformer  
472 is more effective in dealing with limited annotations, which is also verified with results



473 in Tables 1 and 2. Next, when we replace the radiograph transformer with the original  
474 ViT architecture (row 2), which does not have the recurrent concatenation operator, the  
475 overall performance would drop by 3.3 percents. This result verifies the helpfulness of re-  
476 currently concatenating the learned aggregation embedding with patch representations. We  
477 also note that there exists a 3.8-percent performance difference between ResNet and ViT  
478 based architectures (rows 1&2), showing the advantage of a transformer-like architecture.

479 In addition to the radiograph transformer, we also investigate the impact of cross-  
480 supervised learning. First of all, we remove the view fusion module so that different radio-  
481 graphs within a patient study become associated with the same study-level radiology report  
482 (row 3). Such an operation is counter-intuitive as each individual radiograph alone cannot  
483 provide enough information to produce a study-level report. By comparing row 3 with row  
484 0, we found that dropping the view fusion module would reduce the performance by nearly  
485 2 percents on COVID-19 Image Data Collection. This result implies that learning study-  
486 level pre-trained representation is better than image-level pre-training as the former includes  
487 more patient-level information. Next, we completely replace cross-supervised learning with  
488 label-supervised learning (row 4), and REFERS deteriorates into LSP (Transformer) in  
489 Table 2. We found that dropping the two report-related tasks would adversely affect the  
490 performance by 2 percents. Last but not the least, we study the two report-related learning  
491 tasks individually. By comparing row 0 with row 5 and row 6, respectively, we observed  
492 that dropping either of them would not affect the overall performance too much (about 1  
493 percent). This result implies that the effects of both tasks may partially overlap to some  
494 extent. Nonetheless, either of them along with the view fusion module can still outperform  
495 LSP (Transformer) (row 4). In addition, we found that although both of them improve the  
496 overall performance, reinforcing the consistency between representations of each patient  
497 study and its associated report (i.e., the second task) is more crucial than report genera-  
498 tion (i.e., the first task). We believe the reason behind is that the representation learned  
499 in the second task can be regarded as a summary of each report, and thus provides more  
500 global information than token-level predictions in the first task. Such advantages make  
501 it more beneficial for the second task to include more study-level information for learning  
502 better study-level radiograph features.

## 503 **Code Availability**

504 All codes are available at <https://github.com/funnyzhou/REFERS> [40].

## 505 **Data Availability**

506 **MIMIC-CXR-JPG:** <https://physionet.org/content/mimic-cxr-jpg/2.0.0/>.

507

508 **NIH ChestX-ray:** <https://nihcc.app.box.com/v/ChestXray-NIHCC/folder/36938765345>.

509

510 **VinBigData Chest X-ray Abnormalities Detection:** [https://www.kaggle.com/c/](https://www.kaggle.com/c/vinbigdata-chest-xray-abnormalities-detection)

511 [vinbigdata-chest-xray-abnormalities-detection](https://www.kaggle.com/c/vinbigdata-chest-xray-abnormalities-detection).

512

513 **Shenzhen Tuberculosis:** <https://www.kaggle.com/raddar/tuberculosis-chest-xrays-shenzhen>.

514

515 **COVID-19 Image Data Collection:** <https://github.com/ieee8023/covid-chestxray-dataset>.

## 516 **Acknowledgements**

517 This work was supported in part by the Fundamental Research Funds for the Central  
518 Universities (Grant No. 20720190012, 20720210121).

## 519 **Author Contributions Statement**

520 H.Z. and Y.Y. conceived the idea and designed the experiments. H.Z., X.C. and Y.Z.  
521 implemented and performed the experiments. H.Z. and Y.Y. wrote the manuscript. All  
522 authors analyzed the data and experimental results, commented on the manuscript.

## 523 **Competing Interests Statement**

524 The authors declare no competing interests.

## 525 Tables

	NIH	NIH	NIH	VBD	VBD	VBD	SZ	C-T1	C-T2
Method	0.8k (1%)	8k (10%)	80k (100%)	0.1k (1%)	1k (10%)	10k (100%)	All	All	All
Our REFERS	<b>76.7</b>	<b>80.9</b>	<b>84.7</b>	<b>83.0</b>	<b>88.2</b>	<b>90.1</b>	<b>98.0</b>	<b>82.1</b>	<b>80.4</b>
Model Genesis	70.3	75.7	81.0	70.7	82.7	85.8	94.9	76.0	71.8
C2L	71.0	76.6	82.2	75.3	83.3	85.9	95.5	77.8	73.0
Context Restoration	67.8	73.9	78.7	67.9	82.4	83.8	92.7	74.6	69.8
TransVW	71.2	74.3	81.7	73.6	83.8	86.2	94.2	76.1	71.5
ImageNet Pre-training	69.8	74.4	80.0	69.7	82.9	84.5	94.5	74.1	70.3
p-value	8.35e-4	8.72e-4	1.94e-3	8.72e-5	4.34e-4	9.33e-4	1.73e-3	5.88e-4	3.59e-4

Table 1: Comparison with self-supervised learning and transfer learning baselines. **NIH**, **VBD** and **SZ** stand for NIH ChestX-ray, VinBigData Chest X-ray Abnormalities Detection and Shenzhen Tuberculosis datasets, respectively. **C-T1** and **C-T2** denote the two tasks in COVID-19 Image Data Collection, where one task is to distinguish COVID-19 from the rest (C-T1) and the other task is to separate viral pneumonia cases from bacterial ones (C-T2). Note that for the sake of fairness, all baselines use the same transformer-based backbone as the radiograph transformer of REFERS (i.e., a ViT-like architecture plus the recurrent concatenation operator). Each p-value is calculated between our REFERS and the best performing baseline. The evaluation metric is Area under the ROC Curve (AUC). Best results are bolded.

	NIH	NIH	NIH	VBD	VBD	VBD	SZ	C-T1	C-T2
Method	0.8k (1%)	8k (10%)	80k (100%)	0.1k (1%)	1k (10%)	10k (100%)	All	All	All
Our REFERS	<b>76.7</b>	<b>80.9</b>	<b>84.7</b>	<b>83.0</b>	<b>88.2</b>	<b>90.1</b>	<b>98.0</b>	<b>82.1</b>	<b>80.4</b>
LSP (Transformer)	74.2	78.2	82.1	78.5	85.8	87.6	96.4	80.2	76.6
LSP (ConvNet)	65.8	74.5	81.9	76.0	85.2	87.2	96.7	80.1	76.2
p-value	3.25e-3	2.89e-3	5.23e-3	3.56e-4	8.69e-4	1.05e-3	9.65e-3	7.61e-3	1.47e-3

Table 2: Comparison with methods using human-assisted structured labels. **NIH**, **VBD** and **SZ** stand for NIH ChestX-ray, VinBigData Chest X-ray Abnormalities Detection and Shenzhen Tuberculosis datasets, respectively. **C-T1** and **C-T2** denote the two tasks in COVID-19 Image Data Collection, where one task is to distinguish COVID-19 from the rest (C-T1) and the other task is to separate viral pneumonia cases from bacterial ones (C-T2). Note that for fairness, both LSP (Transformer) and REFERS share the same transformer-based backbone (i.e., the ViT architecture plus the recurrent concatenation operator). Each p-value is calculated between the results from our REFERS and LSP (Transformer). The evaluation metric is Area under the ROC Curve (AUC). Best results are bolded.

Row	ViT	RecConcat	View Fusion	Task1	Task2	Viral vs. Bacterial
0	✓	✓	✓	✓	✓	<b>80.4</b>
1			✓	✓	✓	73.3
2	✓		✓	✓	✓	77.1
3	✓	✓		✓	✓	78.6
4	✓	✓				76.6
5	✓	✓	✓	✓		79.1
6	✓	✓	✓		✓	79.3

Table 3: An ablation study of REFERS by removing or replacing individual modules. **RecConcat** stands for the recurrent concatenation operation in the radiograph transformer. **Task1** and **Task2** refer to the two tasks in cross-supervised learning. Row 1 corresponds to the result of a convolutional neural network while row 4 corresponds to LSP (Transformer).

526 **Figures**

Overall workflow of REFERS

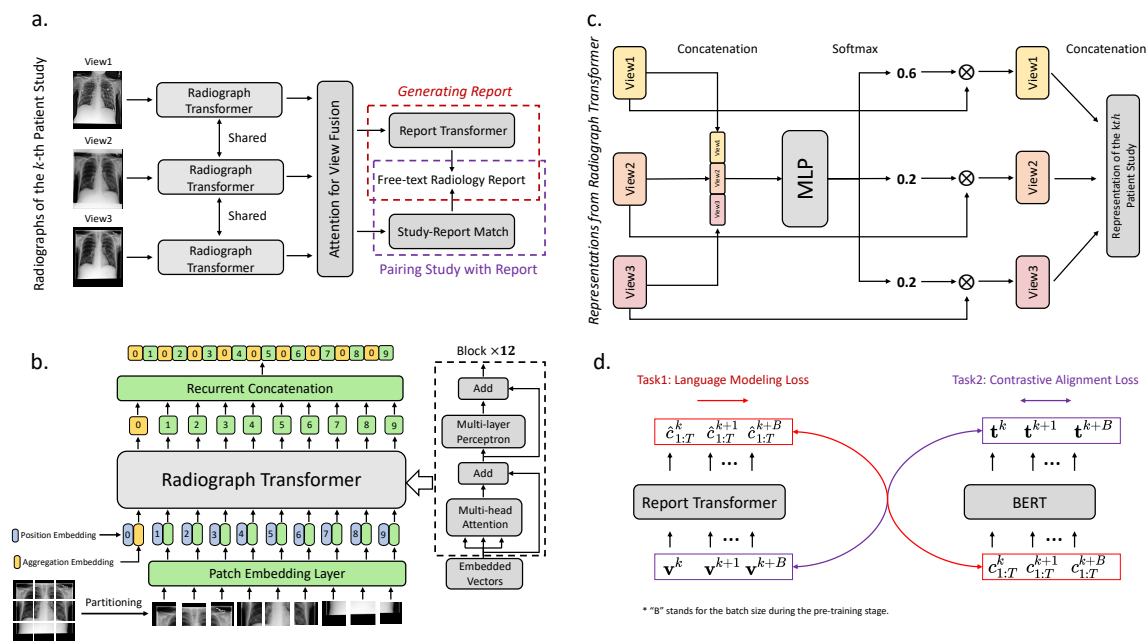


Figure 1: Workflow of REFERS: forwarding radiographs of the  $k$ -th patient study through the radiograph transformer, fusing representations of different views using an attention mechanism, and utilizing report generation and study-report representation consistency reinforcement to exploit the information in radiology reports. Part a provides an overview of the whole pipeline. Part b shows the architecture of the radiograph transformer. Attention for view fusion is elaborated in Part c. Part d presents two supervision tasks, report generation and study-report representation consistency reinforcement. In Part d,  $\mathbf{v}^k$  and  $\mathbf{t}^k$  denote the visual and textual features of the  $k$ -th patient study, respectively.  $\hat{c}_{1:T}^k$  and  $c_{1:T}^k$  stand for the token-level prediction and ground truth of the  $k$ -th radiology report whose length is  $T$ .

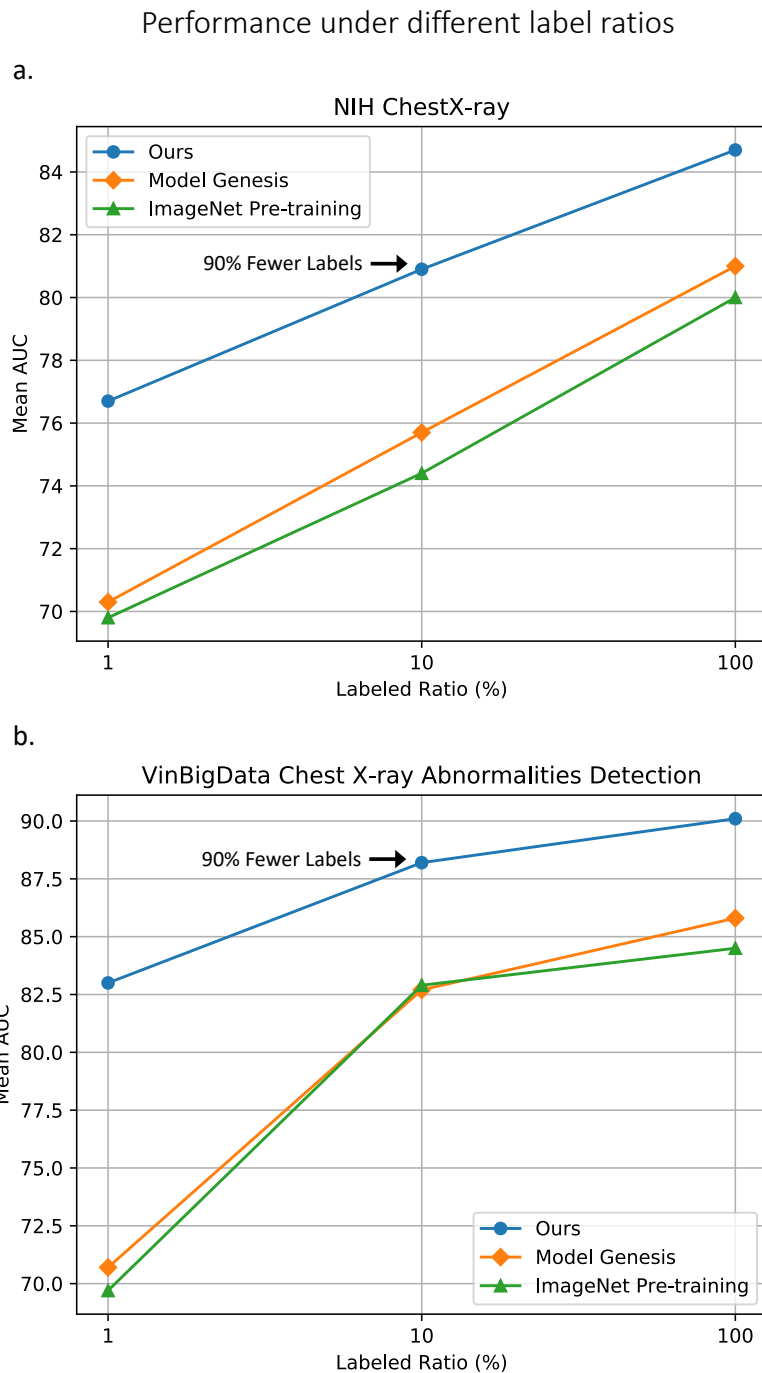


Figure 2: Performance obtained with different amounts of annotated training data in the target domain (a. NIH ChestX-ray and b. VinBigData Chest X-ray Abnormalities Detection). We also denote the percentage of annotated training data in the target domain that our REFERS requires to achieve comparable results with those of Model Genesis and ImageNet pre-training. Note that all three methods share the same transformer-based backbone.

### Visualization of samples from NIH and VBD

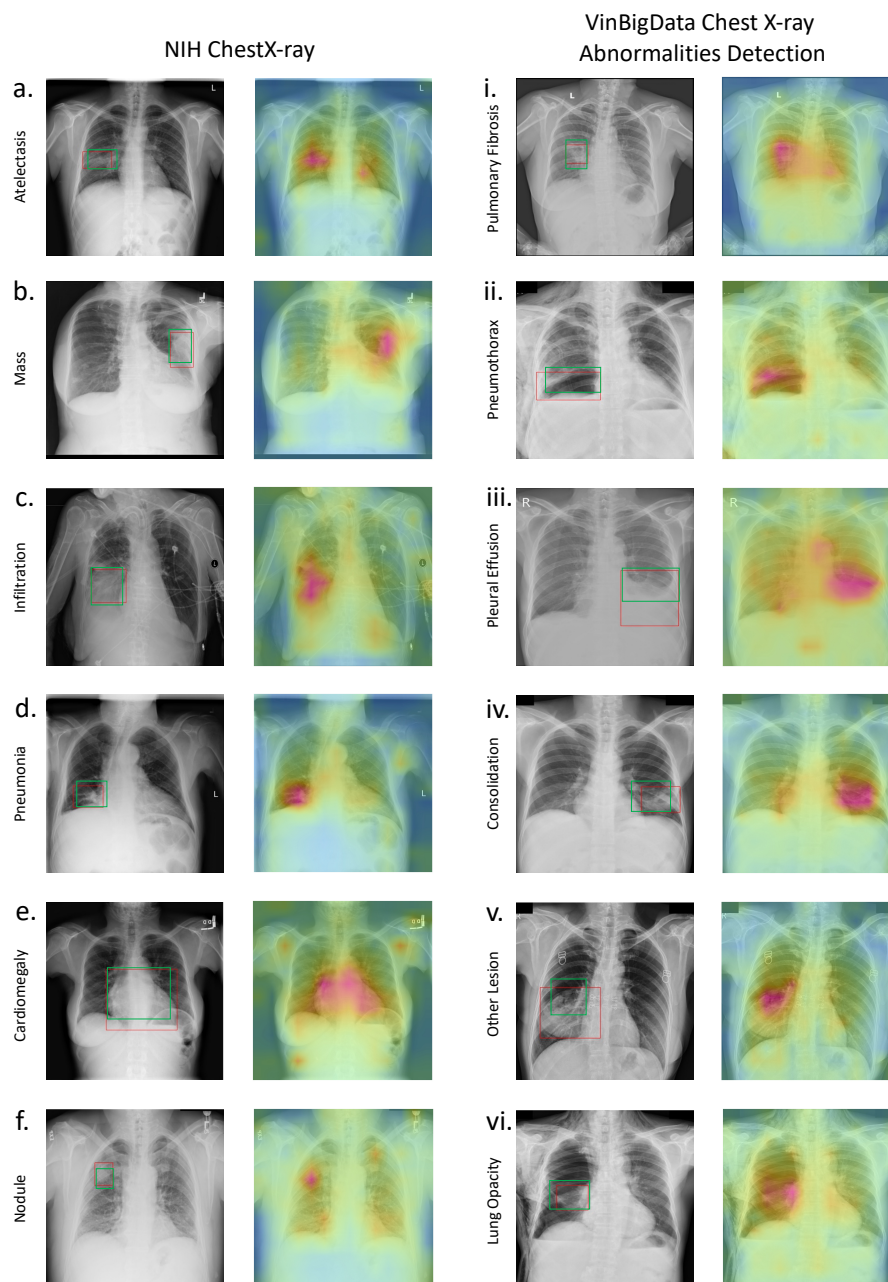


Figure 3: Visualization of twelve randomly chosen samples from NIH ChestX-ray (a-f) and VinBigData (i-vi) (fine-tuned with all annotated training data). For each sample, we present both the original image (left) and an attention map generated from REFERS. In each original image, red boxes denote lesion areas annotated by radiologists. In attention maps, fuchsia color stands for attention values generated from REFERS. The darker the fuchsia color, the higher the confidence of a specific disease. Green boxes in original images are our predicted lesion areas generated by applying a fixed confidence threshold to attention maps.



## 527 **References**

- 528 [1] Krizhevsky, A., Sutskever, I. & Hinton, G.E. ImageNet classification with deep convo-  
529 lutional neural networks. In *Proc. Advances in Neural Information Processing Systems*,  
530 1097-1105 (2012).
- 531 [2] Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image  
532 recognition. In *International Conference on Learning Representations* (2014).
- 533 [3] Szegedy, C. et al. Going deeper with convolutions. In *Proc. IEEE Conference on Com-*  
534 *puter Vision and Pattern Recognition*, 1-9 (IEEE, 2015).
- 535 [4] He, K.M., Zhang, X.Y., Ren, S.Q. & Sun, J. Deep residual learning for image recogni-  
536 tion. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 770-778  
537 (IEEE, 2016).
- 538 [5] Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K.Q. Densely connected convo-  
539 lutional networks. In *Proc. IEEE Conference on Computer Vision and Pattern Recog-*  
540 *nition*, 4700-4708 (IEEE, 2017).
- 541 [6] Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new  
542 perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**, 1798-1828 (IEEE, 2013).
- 543 [7] Phillips, N.A. et al. CheXphoto: 10,000+ photos and transformations of chest x-rays  
544 for benchmarking deep learning robustness. In *Proc. Machine Learning for Health*,  
545 318-327 (PMLR, 2020).
- 546 [8] Taylor, A.G., Mielke, C. & Mongan, J. Automated detection of moderate and large  
547 pneumothorax on frontal chest x-rays using deep convolutional neural networks: a  
548 retrospective study. *PLoS medicine* **15**, e1002697 (Public Library of Science San Fran-  
549 cisco, 2018).
- 550 [9] Carlile, M. et al. Deployment of artificial intelligence for radiographic diagnosis of  
551 COVID-19 pneumonia in the emergency department. *Jour. of the Amer. Coll. of Emer.*  
552 *Phys. Open* **1**, 1459-1464 (Wiley Online Library, 2018).

- 553 [10] Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in deep  
554 neural networks? In *Advances in Neural Information Processing Systems*, 3320-3328  
555 (2014).
- 556 [11] Wang, X.S. et al. ChestX-ray8: Hospital-scale chest x-ray database and benchmarks on  
557 weakly-supervised classification and localization of common thorax diseases. In *Proc.*  
558 *IEEE Conference on Computer Vision and Pattern Recognition*, 2097-2106 (IEEE,  
559 2017).
- 560 [12] Deng, J. et al. ImageNet: A large-scale hierarchical image database. In *Proc. IEEE*  
561 *Conference on Computer Vision and Pattern Recognition*, 248-255 (IEEE, 2009).
- 562 [13] Chen, L. et al. Self-supervised learning for medical image analysis using image context  
563 restoration. *Med. Image Anal.* **58**, 101539 (Elsevier, 2019).
- 564 [14] Zhou, Z.W., Sodha, V., Pang, J.X., Gotway, M.B. & Liang, J.M. Model genesis. *Med.*  
565 *Image Anal.* **67**, 101840 (Elsevier, 2021).
- 566 [15] Haghighi, F., Taher, M.R.H., Zhou, Z.W., Gotway, M.B. & Liang, J.M. Transfer-  
567 able visual words: Exploiting the semantics of anatomical patterns for self-supervised  
568 learning. *IEEE Trans. Med. Imag.*, early access (IEEE, 2021).
- 569 [16] Zhou, H.-Y. et al. Comparing to learn: Surpassing ImageNet pretraining on radio-  
570 graphs by comparing image representations. In *Proc. International Conference on Med-*  
571 *ical Image Computing and Computer-Assisted Intervention*, 398-407 (Springer, 2020).
- 572 [17] Johnson, A.E.W. et al. MIMIC-CXR, a de-identified publicly available database of  
573 chest radiographs with free-text reports. *Sci. Data* **6**, 1-8 (NPG, 2019).
- 574 [18] Irvin, J. et al. CheXpert: A large chest radiograph dataset with uncertainty labels and  
575 expert comparison. In *Proc. the AAAI Conference on Artificial Intelligence*, 590-597  
576 (AAAI, 2019).
- 577 [19] Dosovitskiy, A. et al. An image is worth 16x16 words: Transformers for image recog-  
578 nition at scale. In *International Conference on Learning Representations* (2021).
- 579 [20] Vaswani, A. et al. Attention is all you need. In *Proc. Advances in Neural Information*  
580 *Processing Systems*, 5998-6008 (2017).

- 581 [21] Shin, H.-C. et al. Interleaved text/image deep mining on a very large-scale radiology  
582 database. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*,  
583 1090-1099 (IEEE, 2015).
- 584 [22] Wang, X.S., Peng, Y.F., Lu, L., Lu, Z.Y & Summers, R.M. Tienet: Text-image embed-  
585 ding network for common thorax disease classification and reporting in chest x-rays.  
586 In *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 9049-9058  
587 (IEEE, 2018).
- 588 [23] Johnson, A.E.W. et al. MIMIC-CXR-JPG, a large publicly available database of la-  
589 beled chest radiographs. Preprint at <https://arxiv.org/abs/1901.07042> (2019).
- 590 [24] Nguyen, H.Q. et al. VinDr-CXR: An open dataset of chest x-rays with radiologist's  
591 annotations. Preprint at <https://arxiv.org/abs/2012.15029> (2021).
- 592 [25] Jaeger, S. et al. Two public chest x-ray datasets for computer-aided screening of pul-  
593 monary diseases. *Quantitative Imaging in Medicine and Surgery* **4**, 475 (AME Publi-  
594 cations, 2014).
- 595 [26] Joseph, P.C. et al. COVID-19 image data collection: prospective predictions are the  
596 future. *Journal of Machine Learning for Biomedical Imaging*, early access (2020).
- 597 [27] Zhou, B.L., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Learning deep features  
598 for discriminative localization. In *Proc. IEEE Conference on Computer Vision and*  
599 *Pattern Recognition*, 2921-2929 (IEEE, 2016).
- 600 [28] Chetlur, S. et al. cuDNN: Efficient primitives for deep learning. Preprint at  
601 <https://arxiv.org/abs/1410.0759> (2014).
- 602 [29] He, K.M., Fan, H.Q., Wu, Y.X., Xie, S.N., & Girshick, R. Momentum contrast for  
603 unsupervised visual representation learning. In *Proc. IEEE Conference on Computer*  
604 *Vision and Pattern Recognition*, 9729–9738 (IEEE, 2020).
- 605 [30] Ba, J.L., Kiros, J.R. & Hinton, G.E. Layer normalization. In *International Conference*  
606 *on Learning Representations* (2016).

- 607 [31] Dahl, G.E., Sainath, T.N. & Hinton, G.E. Improving deep neural networks for LVCSR  
608 using rectified linear units and dropout. In *Proc. International Conference on Acous-*  
609 *tics, Speech and Signal Processing*, 8609-8613 (IEEE, 2013).
- 610 [32] Gutmann, M. & Hyvärinen, A. Noise-contrastive estimation: A new estimation prin-  
611 ciple for unnormalized statistical models. In *Proc. the Thirteenth International Con-*  
612 *ference on Artificial Intelligence and Statistics*, 297-304 (JMLR, 2010).
- 613 [33] Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of deep bidi-  
614 rectional transformers for language understanding. In *Proc. the North American Chap-*  
615 *ter of the Association for Computational Linguistics: Human Language Technologies*,  
616 4171-4186 (ACL, 2019).
- 617 [34] Sutskever, I., Martens, J., Dahl, G. & Hinton, G.E. On the importance of initializa-  
618 tion and momentum in deep learning. In *Proc. International Conference on Machine*  
619 *Learning*, 1139-1147 (PMLR, 2013).
- 620 [35] Goyal, P., Mahajan, D., Gupta, A. & Misra, I. Scaling and benchmarking self-  
621 supervised visual representation learning. In *Proc. International Conference on Com-*  
622 *puter Vision*, 6391-6400 (IEEE, 2019).
- 623 [36] Loshchilov, I. & Hutter, F. SGDR: Stochastic gradient descent with warm restarts. In  
624 *International Conference on Learning Representations* (2017).
- 625 [37] Paszke, A. et al. PyTorch: An imperative style, high-performance deep learning library.  
626 In *Proc. Advances in Neural Information Processing Systems*, 8024-8035 (2019).
- 627 [38] Micikevicius, P. et al. Mixed precision training. In *International Conference on Learn-*  
628 *ing Representations* (2018).
- 629 [39] Kingma, D.P. & Ba, J.L. Adam: A method for stochastic optimization. In *International*  
630 *Conference on Learning Representations* (2014).
- 631 [40] Zhou, H.Y., Chen, X.Y., Zhang, Y.H., Luo, R.B., Wang, L.S., & Yu, Y. Generalized  
632 Radiograph Representation Learning via Cross-supervision between Images and Free-  
633 text Radiology Reports. *Zenodo* <https://doi.org/10.5281/zenodo.5624117> (2021).