# Joint Learning of Localized Representations from Medical Images and Reports

Philip Müller[1], Georgios Kaissis[1,2,3], Congyu Zou[4], and Daniel Rueckert[1,3]

[1] Institute of Artificial Intelligence in Medicine, Technical University of Munich, 81675 Munich, Germany
philip.j.mueller@tum.de
[2] Institute of Radiology, Technical University of Munich, 81675 Munich, Germany
[3] Department of Computing, Imperial College London, London SW7 2BX, UK
[4] Department for Internal Medicine I, Klinikum Rechts der Isar, Technical University of Munich, 81675 Munich, Germany

**Abstract.** Contrastive learning has proven effective for pre-training image models on unlabeled data with promising results for tasks such as medical image classification. Using paired text (like radiological reports) during pre-training improves the results even further. Still, most existing methods target image classification downstream tasks and may not be optimal for localized tasks like semantic segmentation or object detection. We therefore propose **Lo**calized representation learning from **V**ision and **T**ext (LoVT), to our best knowledge, the first text-supervised pre-training method that targets localized medical imaging tasks. Our method combines instance-level image-report contrastive learning with local contrastive learning on image region and report sentence representations. We evaluate LoVT and commonly used pre-training methods on an evaluation framework of 18 localized tasks on chest X-rays from five public datasets. LoVT performs best on 10 of the 18 studied tasks making it the preferred method of choice for localized tasks.

**Keywords:** Representation Learning · Contrastive Learning · Text Supervision

## 1 Introduction and Motivation

In medical applications of computer vision, high-quality annotated data is scarce and expensive to acquire, as manually labeling samples typically requires trained physicians[72]. Therefore, the requirement for large labeled datasets can become quite problematic and may limit the applications of deep learning in this field. One approach to overcome this problem is to utilize radiological reports that are paired with medical images. Such reports are produced routinely in clinical practice and are typically written by medical experts (e.g. radiologists). They thus provide a valuable source of semantic information that is available with little additional costs. Rule-based Natural Language Processing (NLP) models like CheXpert[37] extract labels from these reports allowing the automatic creation of large datasets but they also have some significant limitations. Most

importantly, such approaches are typically limited to classification tasks. They generate overall labels for reports (and therefore the paired images) but relating these labels to specific image regions is nontrivial so they cannot be used for localized tasks like semantic segmentation or object detection. Also, rule-based NLP models have to be manually created and cannot generalize to different classification tasks or even different report writing styles[37]. Instead of using these reports to generate classification labels, the reports can be utilized directly in the pre-training method, as was first proposed in the ConVIRT method[96]. Here, the semantic information contained in the reports is used as weak supervision to pre-train image models that are then fine-tuned on labeled downstream tasks, where results can be improved or the number of labeled samples can be reduced. We argue that while this approach is quite promising it is not designed for localized downstream tasks. For example, ConVIRT[96] only works on per-sample image representations and does not explicitly provide more localized representations that might be beneficial for localized tasks like semantic segmentation and object detection. In this work, we therefore study how pre-training methods perform on localized tasks and develop a novel pre-training method designed for localized tasks.

Our contributions are as follows:

- We propose a local contrastive loss allowing to align local representations of sentences or image regions while encouraging spatial smoothness and sensitivity.
- We split each report into sentences and each image into regions (i.e. patches), compute representations for sentences and regions and align them using an attention mechanism and our proposed local contrastive loss.
- We compute global (i.e. per-image and per-report) representations using attention-pooling on the region and sentence representations, and then use a global contrastive loss to align them.
- We propose *Localized representation learning from Vision and Text (LoVT)*, a pre-training method that extends ConVIRT[96] using our proposed ideas and outperforms it on most localized downstream tasks.
- We evaluate our method trained using MIMIC-CXR[42,41,40,26] on a downstream evaluation framework[58] with 18 localized tasks on chest X-rays, including object detection and semantic segmentation on five public datasets. We compare it with several self- and text-supervised methods and with transfer from classification in more than 1400 evaluation runs. Our method LoVT proves as the most successful method outperforming all other methods on 10 out of 18 tasks.

## 2   Related Work

In recent years, contrastive learning[90,63,36,33,57,47,9,31,30,10,94,3,23,7,6], has become the state-of-the-art approach for self-supervised representation learning on images. It has been successfully applied as pre-training method in medical

imaging including downstream tasks such as image classification on chest X-rays[24,76,77].

Most contrastive learning approaches use, unlike our method, only instance-level contrast, i.e. represent each view of the image by a single vector. While the resulting representations are well-suited for global downstream tasks, they are not designed for localized downstream tasks. Therefore, there is a number of recent approaches that use region-level contrast[92,91,88,8,65,56], i.e. they act on representations of image regions. Unlike our method, these methods do not utilize paired text.

Recently however, there is much focus on self-supervised representation learning methods that pre-train image models for downstream tasks by taking advantage of the companion text[67,39,96,12,73,51]. VirTex[12] and ICMLM[73] use image captioning tasks (generative tasks). ConVIRT[96], CLIP[67] and ALIGN[39] on the other hand use multiview contrastive learning[2]. These approaches have been found to be more effective for discriminative downstream tasks[67]. ConVIRT, CLIP, and ALIGN all follow the same general framework where an image and a text encoder are trained jointly using the NT-Xent loss (which is also used in SimCLR) on image and text views. The text views are based on single sentences from companion text, in the case of ConVIRT it is a sentence sampled from the radiology report. The main difference between these methods is the datasets they are studied on, ConVIRT is trained on chest X-rays while the other methods use natural images. Additionally, CLIP uses attention pooling to compute image representations from feature maps while the other methods use the default pooling method from the image encoder (average pooling in the case of ResNet50[32]). Our method follows a similar framework but adds local contrastive losses for better performance on localized tasks. Also, it encodes the whole report instead of sampling a single sentence and uses attention pooling in the image and text encoders. LocTex[51] does localized pre-training on natural images with companion text and predicts alignment of text and image regions. Unlike our method, it uses supervision generated by mouse gazes instead of learning the alignment implicitly using a local contrastive loss. Most related to our work is the recently published local Mutual Information approach [48] that performs contrastive learning on report sentences and image regions but targets classification instead of localized tasks and does therefore neither encourage contrast between regions nor spatial smoothness.

## 3   Method

### 3.1   Assumptions and Intuition

As shown in Fig. 1, a radiology report is typically split into several sections, including a *Findings* section, describing related radiological images, and an *Assessment* section, interpreting the findings. As these sections describe medical aspects observed (*Findings*) in one or more related images and conclusions (*Assessment*) drawn from it, they provide supervision for identifying relevant patterns in the images and interpretations of these patterns. Both sections can be

| **EXAMINATION:** CHEST (PA and LAT) |
|---|
| **INDICATION:**      ___ year old woman with ?pleural effusion |
| **FINDINGS:**<br>Cardiac size cannot be evaluated.<br>Large left pleural effusion is new.<br>Small right effusion is new.<br>The upper lungs are clear.<br>Right lower lobe opacities are better seen in prior<br>CT.<br>There is no pneumothorax.<br>There are mild degenerative changes in the tho-<br>racic spine. |
| **IMPRESSION:**<br>Large left pleural effusion. |

**Fig. 1.** Example radiology report describing chest X-Rays. Taken from the MIMIC-CXR[42,41,26] dataset.

split into sentences and each of these sentences typically describes one or a few aspects of which we assume that most are related to one or a few very localized regions in a paired image. We randomly sample one of the images related to a given report and split it into $7 \times 7$ equally-sized regions. More precisely, we augment and resize the image to a size of $224 \times 224$, feed it into a convolutional neural network, and use the output feature map of size $7 \times 7$ as region representations. A language model encodes the tokens of the report as contextualized (i.e. considering their meaning in the whole report) vector representations from which we compute sentence representations. A many-to-many alignment model is then used to compute *cross-modal representations* from *uni-modal representations*, i.e. image region representations from sentence representations and vice-versa. We argue that by aligning cross-modal and uni-modal representations, the image region representations are encouraged to contain the high-level semantics present in the report.

### 3.2   Model Overview

Fig. 2 shows the general architecture of our proposed LoVT model. Each training sample $\boldsymbol{x}_i$ is a pair of an image $\boldsymbol{x}_i^{\mathcal{I}} \in \mathbb{R}^{224 \times 224}$ and the related report $\boldsymbol{x}_i^{\mathcal{R}}$ consisting of $M_i$ sentences. Both, $\boldsymbol{x}_i^{\mathcal{I}}$ and $\boldsymbol{x}_i^{\mathcal{R}}$, are encoded independently into two global representations, for image and report respectively, and multiple local representations per sample, corresponding to image regions and report sentences, respectively. An attention-based alignment model then computes cross-modal representations (i.e. sentence representations from image regions and vice-versa) which are aligned with the local uni-modal representations using local contrastive losses. Additionally, the global representations are aligned using a global contrastive loss. The encoders and the alignment model are trained jointly on batches of image-report pairs $\boldsymbol{x}_i$. The details of the model and the loss function will be described in the following sections.
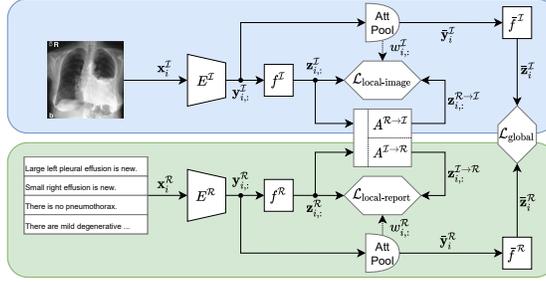
**Fig. 2.** Architecture of LoVT. Given an image $\boldsymbol{x}_i^{\mathcal{I}}$ and the related report $\boldsymbol{x}_i^{\mathcal{R}}$, the encoders $E^{\mathcal{I}}$ and $E^{\mathcal{R}}$ compute image region and report sentence representations, respectively, which are projected using $f^{\mathcal{I}}$ and $f^{\mathcal{R}}$. The alignment models $A^{\mathcal{R}\to\mathcal{I}}$ and $A^{\mathcal{I}\to\mathcal{R}}$ compute cross-modal report-to-image ($\boldsymbol{z}_{i,k}^{\mathcal{R}\to\mathcal{I}}$) and image-to-report ($\boldsymbol{z}_{i,m}^{\mathcal{I}\to\mathcal{R}}$) representations which are aligned with the uni-modal representations ($\boldsymbol{z}_{i,k}^{\mathcal{I}}$ and $\boldsymbol{z}_{i,m}^{\mathcal{R}}$) using the local losses $\mathcal{L}_{\text{local-image}}$ and $\mathcal{L}_{\text{local-report}}$, respectively. Global image ($\bar{\boldsymbol{y}}_i^{\mathcal{I}}$) and report ($\bar{\boldsymbol{y}}_i^{\mathcal{R}}$) representations are computed using attention pooling on the local representations, are then projected using $\bar{f}^{\mathcal{I}}$ and $\bar{f}^{\mathcal{R}}$ and aligned using the global loss $\mathcal{L}_{\text{global}}$.

### 3.3   Encoding

Each image $\boldsymbol{x}_i^{\mathcal{I}}$ is encoded into $K = H \times W$ (we use $K = 7 \times 7$) region representations $\boldsymbol{y}_{i,k}^{\mathcal{I}} \in \mathbb{R}^{d^{\mathcal{I}}}$ using the image encoder $E^{\mathcal{I}}$, where $k$ is the index of the image region, and $d^{\mathcal{I}}$ is the dimension of the image region representation space. Our approach is encoder agnostic, i.e. any model encoding image regions into vector representations can be used for $E^{\mathcal{I}}$. We use a ResNet50[32] and take the feature map before global average pooling as region representations. Similarly, each report $\boldsymbol{x}_i^{\mathcal{R}}$ is encoded into $M_i$ sentence representations $\boldsymbol{y}_{i,m}^{\mathcal{R}} \in \mathbb{R}^{d^{\mathcal{R}}}$ using the report encoder $E^{\mathcal{R}}$. Here $M_i$ is the number of sentences of report sample $i$, $m$ is the index of the sentence, and $d^{\mathcal{R}}$ is the dimension of the report sentence representation space. Note that while $K$ is constant, $M_i$ may be different for each sample. Any model encoding sentences into vector representations can be used for $E^{\mathcal{R}}$. We use BERT_base[14] to jointly encode the tokens of the concatenated sentences of each report and then perform max pooling over the token representations of each sentence to get sentence representations.

The global (i.e. per-sample) representations $\bar{\boldsymbol{y}}_i^{\mathcal{I}}$ and $\bar{\boldsymbol{y}}_i^{\mathcal{R}}$ are each computed by an attention pooling layer (not shared between modalities) on the region and sentence representations, respectively. It is implemented using multi-head query-key-value attention[82] where the query is computed from the globally averaged region or sentence representations. This pooling approach was first proposed for the image encoder of CLIP[67].

Following previous works[96,9,30], we compute projected local representations $\boldsymbol{z}_{i,k}^{\mathcal{I}} \in \mathbb{R}^{d^{\mathcal{Z}}}$ and $\boldsymbol{z}_{i,m}^{\mathcal{R}} \in \mathbb{R}^{d^{\mathcal{Z}}}$, and projected global representations $\bar{\boldsymbol{z}}_i^{\mathcal{I}} \in \mathbb{R}^{\bar{d}^{\mathcal{Z}}}$ and $\bar{\boldsymbol{z}}_i^{\mathcal{R}} \in \mathbb{R}^{\bar{d}^{\mathcal{Z}}}$ from the representations $\boldsymbol{y}_{i,k}^{\mathcal{I}}$, $\boldsymbol{y}_{i,m}^{\mathcal{R}}$, $\bar{\boldsymbol{y}}_i^{\mathcal{I}}$, and $\bar{\boldsymbol{y}}_i^{\mathcal{R}}$, using the (non-shared) nonlinear transformations $f^{\mathcal{I}}$, $f^{\mathcal{R}}$, $\bar{f}^{\mathcal{I}}$, and $\bar{f}^{\mathcal{R}}$, respectively, where $d^{\mathcal{Z}}$

is the dimension of the shared local and $\bar{d}^{\mathcal{Z}}$ of the shared global representation space (we use 512 for both). Note that for local representations the projections are applied to each region $k$ or sentence $m$ independently.

### 3.4   Alignment Model

Following our assumptions (see Sec. 3.1), we compute an alignment of image regions and sentences and compute cross-modal representations using the alignment models $A^{\mathcal{I}\to\mathcal{R}}$ and $A^{\mathcal{R}\to\mathcal{I}}$, which are based on single-head query-key-value attention[82].

For each sentence $m$ the cross-modal representation $z_{i,m}^{\mathcal{I}\to\mathcal{R}}$ is computed by letting $z_{i,m}^{\mathcal{R}}$ attend to all image region representations $z_{i,k}^{\mathcal{I}}$ (of the related image). We therefore compute the probability $\alpha_{i,m,k}^{\mathcal{I}\to\mathcal{R}}$ that sentence $m$ is aligned with region $k$ based on the scaled dot product scores of their projected representations, i.e. $\alpha_{i,m,k}^{\mathcal{I}\to\mathcal{R}} = \text{softmax}_k\left(\frac{(\boldsymbol{Q}z_{i,m}^{\mathcal{R}})^T(\boldsymbol{Q}z_{i,k}^{\mathcal{I}})}{\sqrt{d^{\mathcal{Z}}}}\right)$, where the linear query-key projection $\boldsymbol{Q}$ is a learned matrix. Then the alignment model $A^{\mathcal{I}\to\mathcal{R}}$ uses $\alpha_{i,m,k}^{\mathcal{I}\to\mathcal{R}}$ to compute $z_{i,m}^{\mathcal{I}\to\mathcal{R}}$ as projected weighted sum of the image region representations $z_{i,k}^{\mathcal{I}}$:

$$z_{i,m}^{\mathcal{I}\to\mathcal{R}} = \boldsymbol{O}\left(\sum_{k=1}^{K}\alpha_{i,m,k}^{\mathcal{I}\to\mathcal{R}}\left(\boldsymbol{V}z_{i,k}^{\mathcal{I}}\right)\right), \tag{1}$$

where the value projection $\boldsymbol{V}$, and the output projection $\boldsymbol{O}$ are learned matrices.

In a similar fashion the cross-modal representations $z_{i,k}^{\mathcal{R}\to\mathcal{I}}$ are computed by $A^{\mathcal{R}\to\mathcal{I}}$:

$$z_{i,k}^{\mathcal{R}\to\mathcal{I}} = \boldsymbol{O}\left(\sum_{m=1}^{M_i}\alpha_{i,k,m}^{\mathcal{R}\to\mathcal{I}}\left(\boldsymbol{V}z_{i,m}^{\mathcal{R}}\right)\right), \tag{2}$$

with $\alpha_{i,k,m}^{\mathcal{R}\to\mathcal{I}} = \text{softmax}_m\left(\frac{(\boldsymbol{Q}z_{i,k}^{\mathcal{I}})^T(\boldsymbol{Q}z_{i,m}^{\mathcal{R}})}{\sqrt{d^{\mathcal{Z}}}}\right)$. Note that as $A^{\mathcal{R}\to\mathcal{I}}$ and $A^{\mathcal{I}\to\mathcal{R}}$ share the same matrices $\boldsymbol{Q}$, $\boldsymbol{V}$, and $\boldsymbol{O}$, the only difference between $\alpha_{i,k,m}^{\mathcal{R}\to\mathcal{I}}$ and $\alpha_{i,m,k}^{\mathcal{I}\to\mathcal{R}}$ is transposition and the index over which softmax is applied.

### 3.5   Loss Function

*Global Alignment* For global alignment we follow ConVIRT[96] and maximize the cosine similarity between paired image and report representations while minimizing the similarity between non-paired (i.e. from different samples) representations. The loss consists of a image-report part, where all non-paired report representations from the batch are used as negatives:

$$\ell_{\text{global}}^{\mathcal{I}\|\mathcal{R}} = -\log\frac{e^{\cos\left(\bar{z}_i^{\mathcal{I}},\bar{z}_i^{\mathcal{R}}\right)/\tau}}{\sum_j e^{\cos\left(\bar{z}_i^{\mathcal{I}},\bar{z}_j^{\mathcal{R}}\right)/\tau}}, \tag{3}$$

and a report-image part, defined analogously:

$$\ell_{\text{global}}^{\mathcal{R}\|\mathcal{I}} = -\log \frac{e^{\cos\left(\bar{z}_i^{\mathcal{R}}, \bar{z}_i^{\mathcal{I}}\right)/\tau}}{\sum_j e^{\cos\left(\bar{z}_i^{\mathcal{R}}, \bar{z}_j^{\mathcal{I}}\right)/\tau}} \ , \tag{4}$$

where $\tau$ is the similarity temperature (we use $0.1$) and all logarithms are natural. Both parts are combined using the hyperparameter $\lambda \in [0, 1]$ (we use $0.75$):

$$\mathcal{L}_{\text{global}} = \frac{1}{N} \sum_{i=1}^{N} \left[ \lambda \cdot \ell_{\text{global}}^{\mathcal{I}\|\mathcal{R}} + (1 - \lambda) \cdot \ell_{\text{global}}^{\mathcal{R}\|\mathcal{I}} \right] \ . \tag{5}$$

*Local Alignment* The global alignment loss does not only align the global representations but it also prevents the global representations from collapsing to a constant vector using negative samples to contrast the positive pairs. Similarly, we propose local alignment losses encouraging spatial (sentence) sensitivity through negatives from the same sample, i.e. preventing the local representations to be similar for all regions (sentences) of an image (report). We use two NT-Xent-based[9] local losses: $\mathcal{L}_{\text{local-image}}$, aligning region representations $z_{i,k}^{\mathcal{I}}$ with $z_{i,k}^{\mathcal{R}\to\mathcal{I}}$, and $\mathcal{L}_{\text{local-report}}$, aligning sentence representations $z_{i,m}^{\mathcal{R}}$ with $z_{i,m}^{\mathcal{I}\to\mathcal{R}}$.

Some regions or sentences may not be relevant for aligning a sample (e.g. background regions or sentences not related to the image). Therefore, we introduce region weights $w_{i,k}^{\mathcal{I}}$ and sentence weights $w_{i,m}^{\mathcal{R}}$, which are computed as the attention probabilities from the respective attention pooling layer (which was used to compute global representations), averaged over all attention heads. These weights are used in the local loss functions such that irrelevant representations do not have to be aligned. Note that we do not backpropagate through the region or sentence weights.

The loss $\mathcal{L}_{\text{local-image}}$ allows for having multiple positive pairs within each sample by giving each pair of regions $(k, l)$ a positiveness probability $p_{k,l}^{\mathcal{I}} \in [0, 1]$. We then treat each positive pair as its own (weighted) example and contrast it with all other pairs (again all logarithms are natural):

$$\ell_{\text{local-image}}^{\mathcal{I}\|\mathcal{R}\to\mathcal{I}} = -\sum_{l=1}^{K} p_{k,l}^{\mathcal{I}} \log \frac{e^{\cos\left(z_{i,k}^{\mathcal{I}}, z_{i,l}^{\mathcal{R}\to\mathcal{I}}\right)/\tau'}}{\sum_{k'} e^{\cos\left(z_{i,k}^{\mathcal{I}}, z_{i,k'}^{\mathcal{R}\to\mathcal{I}}\right)/\tau'}} \tag{6}$$

$$\ell_{\text{local-image}}^{\mathcal{R}\to\mathcal{I}\|\mathcal{I}} = -\sum_{l=1}^{K} p_{k,l}^{\mathcal{I}} \log \frac{e^{\cos\left(z_{i,k}^{\mathcal{R}\to\mathcal{I}}, z_{i,l}^{\mathcal{I}}\right)/\tau'}}{\sum_{k'} e^{\cos\left(z_{i,k}^{\mathcal{R}\to\mathcal{I}}, z_{i,k'}^{\mathcal{I}}\right)/\tau'}} \tag{7}$$

$$\mathcal{L}_{\text{local-image}} = \frac{1}{2N} \sum_{i=1}^{N} \sum_{k=1}^{K} w_{i,k}^{\mathcal{I}} \cdot \left[ \ell_{\text{local-image}}^{\mathcal{I}\|\mathcal{R}\to\mathcal{I}} + \ell_{\text{local-image}}^{\mathcal{R}\to\mathcal{I}\|\mathcal{I}} \right] \ . \tag{8}$$

Here $\tau'$ is the similarity temperature and is set to $0.3$. We assume that nearby image regions are often similar and that therefore nearby regions are more likely

to be positives while distant regions are more likely to be negatives. Thus, we define the positiveness probability $p_{k,l}^{\mathcal{I}}$ of two image regions as the complementary cumulative exponential distribution of $d_{\boldsymbol{x}}$ (their spatial $\ell_2$-distance in 2D space normalized by the length of the diagonal $\sqrt{H^2 + W^2}$) and set $p_{k,l}^{\mathcal{I}}$ to zero above cutoff threshold $T \in [0, \infty)$:

$$p_{k,l}^{\mathcal{I}} = \frac{\mathbb{1}_{[d_{\boldsymbol{x}}(k,l) \leq T]} \cdot e^{-d_{\boldsymbol{x}}(k,l)/\beta}}{\sum_{k'} \mathbb{1}_{[d_{\boldsymbol{x}}(k,k') \leq T]} \cdot e^{-d_{\boldsymbol{x}}(k,k')/\beta}} \ . \tag{9}$$

Here $\beta \in (0, \infty)$ is a sharpness hyperparameter. We set $\beta = 1$ and $T = 0.5$. Note that the normalization of $d_{\boldsymbol{x}}$ is equal to rescaling $T$ and $\beta$, i.e. it allows us to define both hyperparameters independently of the image size.

The definition of $p_{k,l}^{\mathcal{I}}$ is derived by modeling the occurrence of related features at specific distances in the image as a Poisson point process, such that the $\ell_2$-distance of related features follows the exponential distribution. We assume a Poisson process due to its property of being memoryless, i.e. knowing that a feature is already related to another feature at some distance does not change how distant additional related features can be found. Also, the probability density function of the exponential distribution is decreasing (with support on the interval $[0, \infty)$), which seems reasonable as it is typically more likely that related features are near than far. Its cumulative distribution function then describes the probability that two related features are within a given radius and its complementary function that of being outside a given radius. The threshold $T$ assures that very distant pairs do not count as positives. The loss $\mathcal{L}_{\text{local-image}}$ thus encourages spatial smoothness of image regions while maintaining spatial sensitivity through negative samples. Note that it is related to the pixel-contrast loss proposed in[92], where the main novelty of our work is the partly smooth definition of $p_{k,l}^{\mathcal{I}}$ based on the exponential distribution.

The local report loss $\mathcal{L}_{\text{local-report}}$ is defined similarly but we do not assume prior knowledge about the similarity of sentences and therefore only have a single positive pair per sentence (again all logarithms are natural):

$$\ell_{\text{local-report}}^{\mathcal{R} \| \mathcal{I} \to \mathcal{R}} = -\log \frac{e^{\cos\left(\boldsymbol{z}_{i,m}^{\mathcal{R}}, \boldsymbol{z}_{i,m}^{\mathcal{I} \to \mathcal{R}}\right)/\tau'}}{\sum_{m'} e^{\cos\left(\boldsymbol{z}_{i,m}^{\mathcal{R}}, \boldsymbol{z}_{i,m'}^{\mathcal{I} \to \mathcal{R}}\right)/\tau'}} \tag{10}$$

$$\ell_{\text{local-report}}^{\mathcal{I} \to \mathcal{R} \| \mathcal{R}} = -\log \frac{e^{\cos\left(\boldsymbol{z}_{i,m}^{\mathcal{I} \to \mathcal{R}}, \boldsymbol{z}_{i,m}^{\mathcal{R}}\right)/\tau'}}{\sum_{m'} e^{\cos\left(\boldsymbol{z}_{i,m}^{\mathcal{I} \to \mathcal{R}}, \boldsymbol{z}_{i,m'}^{\mathcal{R}}\right)/\tau'}} \tag{11}$$

$$\mathcal{L}_{\text{local-report}} = \frac{1}{2N} \sum_{i=1}^{N} \sum_{m=1}^{M_i} w_{i,m}^{\mathcal{R}} \cdot \left[ \ell_{\text{local-report}}^{\mathcal{R} \| \mathcal{I} \to \mathcal{R}} + \ell_{\text{local-report}}^{\mathcal{I} \to \mathcal{R} \| \mathcal{R}} \right] \tag{12}$$

*Total Loss* The total loss $\mathcal{L}$ is computed as the weighted sum of global and local losses:

$$\mathcal{L} = \gamma \cdot \mathcal{L}_{\text{global}} + \mu \cdot \mathcal{L}_{\text{local-image}} + \nu \cdot \mathcal{L}_{\text{local-report}} \ , \tag{13}$$

where $\gamma$, $\mu$, and $\nu$ are loss weights to balance the individual losses and are set to 1.0, 0.75, and 0.75, respectively. We determined these loss weights by running small grid searches (see Appendix E.1 for details).

## 4   Evaluation

### 4.1   Downstream Tasks and Experimental Setup

We evaluate our method on a downstream evaluation framework[58] with 18 localized tasks on chest X-rays, which we will shortly describe here. For more details, we refer to Appendix E.3.

*Evaluation Protocols* We only use the pre-trained ResNet50 (from the image encoder). For semantic segmentation tasks we evaluate in the following settings:: (i) **U-Net Finetune**: Here the ResNet50 is used as backbone of a U-Net[70] and is finetuned jointly with all other layers, (ii) **U-Net Frozen**: Here the ResNet50 is used as frozen backbone of a U-Net[70] and only the non-backbone layers are finetuned, and (iii) **Linear**: Here an element-wise linear layer is trained that is applied after the last feature map (before pooling) of the frozen ResNet50, before the results are upsampled to the segmentation resolution.

For object detection tasks we use the following protocols: (i) **YOLOv3 Finetune**: Here the ResNet50 is used as backbone of a YOLOv3[69] model and is finetuned jointly with the non-backbone layers, (ii) **YOLOv3 Frozen**: Here the ResNet50 is used as frozen backbone of a YOLOv3[69] model and only the non-backbone layers are finetuned, and (iii) **Linear**: Here the object detection ground truth is converted to segmentation masks and then the *Linear* segmentation protocol is used for evaluation.

*Downstream Datasets* We evaluate the pre-trained ResNet50 on several medical datasets, namely (i) **RSNA Pneumonia Detection**[86,74], with more than 260000 frontal-view chest X-rays with detection targets for pneumonia opacities. We use the *YOLOv3 Finetune*, *YOLOv3 Frozen*, and *Linear* protocols, each with 1%, 10%, and 100% of the training samples; (ii) **COVID Rural**[81,13], with more than 200 frontal-view chest X-rays with segmentation masks for COVID-19 lung opacity regions. We use the *UNet Finetune*, *UNet Frozen*, and *Linear* protocols; (iii) **SIIM-ACR Pneumothorax Segmentation**[75], with more than 12000 frontal-view chest X-rays with segmentation masks for pneumothorax. We use the *UNet Finetune*, *UNet Frozen* protocols, but due not use *Linear* due to the fine-grained nature of the segmentation masks; (iv) **Object CXR**[38] with 9000 frontal-view chest X-rays with detection targets for foreign objects. We use the *YOLOv3 Finetune*, *YOLOv3 Frozen*, and *Linear* protocols; (v) **NIH CXR**[86], with almost 1000 frontal-view chest X-rays with detection targets for eight pathologies (Atelectasis, Cardiomegaly, Effusion, Infiltrate, Mass, Nodule, Pneumonia, and Pneumothorax). Due to the limited data per class, we only use the *Linear* protocol. The different evaluation protocols complement each other: While the *U-Net Finetune* and *YOLOv3 Finetune* protocols evaluate how well

the pre-trained image models could be fine-tuned for practical applications, the *Linear* protocols directly evaluate the learned local representations (i.e. feature maps) while adding as few parameters as possible and therefore mostly omitting the variance introduced by random initialization during downstream evaluation. The *U-Net Frozen* and *YOLOv3 Frozen* protocols can be seen as middle ground between the two extremes, where representations are frozen but evaluated in a more practical setting (but with many randomly initialized layers). Overall this allows the analysis of many aspects of the pre-trained representations.

*Tuning and Evaluation Procedure* We tune all models on a single downstream task, *RSNA YOLOv3 Frozen 10%*. Other downstream tasks have not been evaluated during tuning to make sure that models are not biased towards the downstream tasks. After tuning, each model was evaluated on all downstream tasks. For each task the downstream learning rates were tuned individually per model (using single evaluation runs) before running five evaluations (all using the tuned learning rate). We report the average results of these five runs and their 95%-confidence interval.

*Pre-Training Dataset* We train our method on version 2 of MIMIC-CXR[40,41,42,26] as, to our best knowledge, it is the largest and most commonly used dataset of this kind. Since all downstream tasks contain only frontal views, we remove all lateral views, such that roughly 21000 training samples remain, each with a report and one or more frontal images.

*Baselines* We compare our method against several baseline methods:

- **Random Init.**: The ResNet50 is initialized using its default random initialization
- **ImageNet[71] Init.**: The ResNet50 is initialized with weights pre-trained on the ImageNet ILSVRC-2012 task[71];
- **CheXpert[37]**: The ResNet50 is pre-trained using supervised multi-label binary classification with CheXpert[37] labels on frontal chest X-rays of MIMIC-CXR
- **Global image pre-training methods**: The ResNet50 is pre-trained using the self-supervised pre-training methods SimCLR[9] or BYOL[30] on frontal chest X-rays of MIMIC-CXR. We decided to include SimCLR as is uses a similar loss function as LoVT and we include BYOL because of its widespread use.
- **Local image pre-training methods**: The ResNet50 is pre-trained using the self-supervised pre-training method PixelPro[92] on frontal chest X-rays of MIMIC-CXR. We include PixelPro to study the effect of local contrastive losses when using only images.
- **Global image-text pre-training methods**: The ResNet50 is pre-trained using the image-text methods ConVIRT[96] or CLIP[67] on frontal MIMIC-CXR. Note that for comparability we adapted CLIP to use the same image and text encoders as ConVIRT such that the main difference between CLIP

**Table 1.** Results on the RSNA pneumonia detection tasks with different training set sizes. All results are averaged over five evaluation runs and the 95%-confidence interval is shown. The best results per task are underlined, the second-best results are dash-underlined and the best results per pre-training category (general initialization, pre-training on 30% and 100%) are highlighted in bold. Note that the *YOLOv3 Frozen 10%* task (task 5) was used for tuning of all methods and may therefore not be representative as methods may overfit on this task.

| | RSNA YOLOv3 Finetune mAP (%) | | | RSNA YOLOv3 Frozen mAP (%) | | | RSNA Lin. Seg. Dice (%) | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1% | 10% | 100% | 1% | 10% | 100% | 1% | 10% | 100% |
| *General initialization methods* | | | | | | | | | |
| Random | 2.4±0.5 | 5.1±1.2 | 14.9±1.7 | 1.0±0.2 | 4.0±0.3 | 8.9±0.9 | 21.9±1.2 | 5.3±0.0 | 5.3±0.0 |
| ImageNet [71] | **5.0±0.7** | **12.4±0.8** | **19.0±0.2** | **3.6±1.4** | **8.0±0.1** | **15.7±0.3** | **27.5±0.6** | **38.3±0.0** | **43.3±0.0** |
| *Pre-Training on 30 % of frontal MIMIC-CXR* | | | | | | | | | |
| CheXpert [37] | **8.3±0.8** | 12.4±1.6 | **21.3±0.3** | 7.0±1.0 | 14.8±0.8 | 18.8±0.4 | 38.9±0.2 | 45.5±0.2 | 48.1±0.0 |
| BYOL [30] | 7.0±1.0 | 11.9±1.1 | 18.8±0.2 | 9.6±0.2 | 14.0±1.2 | **21.0±0.2** | 42.9±0.1 | 47.8±0.2 | 50.0±0.0 |
| SimCLR [9] | 6.7±0.5 | **12.9±0.5** | 20.4±1.8 | 7.9±1.0 | 11.9±0.1 | 19.9±0.2 | 43.1±0.0 | 46.0±0.0 | 48.2±0.0 |
| PixelPro [92] | 4.8±0.6 | 12.6±1.2 | 19.8±0.4 | 3.1±0.2 | 6.4±0.5 | 13.4±0.6 | 25.9±0.2 | 34.6±0.0 | 39.8±0.1 |
| ConVIRT [96] | 7.4±1.3 | 12.7±1.5 | 18.3±0.4 | **9.8±0.3** | 14.8±1.1 | 18.4±1.1 | 42.1±0.1 | 47.1±0.2 | 50.2±0.0 |
| CLIP [67]* | 7.2±0.8 | 12.8±1.2 | 19.7±0.5 | 9.3±0.4 | 16.1±1.1 | 19.6±1.4 | 44.3±0.1 | 48.8±0.1 | 50.7±0.0 |
| LoVT (Ours) | 7.7±1.0 | 11.7±0.5 | 17.2±1.3 | 8.6±1.5 | **17.9±0.4** | 18.0±0.1 | **46.0±0.0** | **49.4±0.0** | **51.5±0.0** |
| *Pre-Training on 100 % of frontal MIMIC-CXR* | | | | | | | | | |
| CheXpert [37] | 10.0±1.9 | 12.4±0.9 | **22.2±0.4** | 5.8±0.4 | 11.9±0.7 | 20.0±0.2 | 40.0±0.1 | 44.3±0.0 | 46.9±0.0 |
| BYOL [30] | 5.6±0.8 | 11.0±0.2 | 17.3±1.1 | 6.8±1.6 | 12.1±1.1 | 15.9±0.6 | 41.9±0.0 | 45.1±0.0 | 46.8±0.0 |
| SimCLR [9] | 7.1±0.7 | 12.2±0.8 | 18.8±1.0 | 5.4±0.2 | 13.1±0.2 | 17.3±1.6 | 43.0±0.0 | 45.1±0.0 | 47.0±0.0 |
| PixelPro [92] | 4.8±0.3 | 11.0±1.5 | 17.4±1.7 | 4.6±1.6 | 5.4±1.1 | 12.6±1.3 | 23.9±0.4 | 34.8±0.2 | 40.2±0.1 |
| ConVIRT [96] | **10.7±1.1** | **13.3±0.8** | 18.5±0.4 | 8.2±0.9 | 15.6±1.2 | 17.9±0.3 | 44.6±0.1 | 48.5±0.0 | 50.4±0.3 |
| CLIP [67]* | 7.0±1.5 | 10.7±1.1 | 19.9±0.8 | **11.9±0.7** | 15.0±1.1 | 18.7±0.0 | 45.2±0.0 | 49.3±0.1 | 51.1±0.0 |
| LoVT (Ours) | 8.5±0.8 | 13.2±0.6 | 18.1±3.2 | 9.6±1.2 | **16.4±1.3** | **20.5±1.0** | **46.3±0.0** | **50.1±0.0** | **51.8±0.0** |
| Task Nr. | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

\* Modified to use the same image and text encoders as ConVIRT and LoVT.

and ConVIRT is that CLIP uses attention pooling to compute the scan representation while ConVIRT uses average pooling. We include both methods as LoVT builds upon a similar general framework, where we include ConVIRT because it targets chest X-rays (like LoVT) and include CLIP because of its widespread use and as it uses (like LoVT) attention pooling in the image encoder. We decided not to include VirTex[12] and ICMLM[73] as they use generative tasks, which have been found to be less effective for discriminative downstream tasks[67].

## 4.2 Downstream Results

We present the downstream results of our model LoVT and the baselines, with pre-training on 100% and 30% of MIMIC-CXR. Tab. 1 shows the results on different subsets of the RSNA dataset and Tab. 2 shows the results on the remaining downstream datasets, i.e. on COVID Rural, SIM-ACR Pneumothorax, Object CXR, and NIH CXR.

*Comparison of Methods* We found that there is no single pre-training method performing best on all evaluated downstream tasks. On most tasks (15 out of 18) image-text self-supervised methods (i.e. LoVT, CLIP, or ConVIRT) outperform the other methods, such that they should be preferred if paired text is available.

Our model LoVT is the best method (over all pre-training settings) on 10 of 18 tasks, and significantly outperforms all other methods in 6 of these tasks,

**Table 2.** Results on downstream tasks on the COVID Rural, SIIM Pneumothorax, Object CXR, and NIH CXR datasets. All results are averaged over five evaluation runs and the 95%-confidence interval is shown. The best results per task are underlined, the second-best results are dash-underlined and the best results per pre-training category (general initialization, pre-training on 30% and 100%) are highlighted in bold.

| | COVID Rural | | | SIIM-ACR Pneumoth. | | Object CXR | | | NIH CXR |
| | UNet Finetune | UNet Frozen | Linear | UNet Finetune | UNet Frozen | YOLOv3 Finetune | YOLOv3 Frozen | Linear | Linear |
| | Dice (%) | Dice (%) | Dice (%) | Dice (%) | Dice (%) | fROC (%) | fROC (%) | Dice (%) | Avg Dice (%) |
|---|---|---|---|---|---|---|---|---|---|
| *General initialization methods* | | | | | | | | | |
| Random | 34.0±1.1 | 32.2±1.8 | 6.0±0.0 | 23.2±1.0 | 23.9±1.6 | 49.5±1.2 | 28.4±1.4 | 6.9±0.0 | 0.5±0.4 |
| ImageNet [71] | **43.9±2.0** | **41.9±1.7** | **32.6±0.7** | **38.5±0.9** | **36.9±0.7** | **62.5±0.4** | **52.7±1.3** | **37.8±0.0** | **2.6±1.6** |
| *Pre-Training on 30 % of frontal MIMIC-CXR* | | | | | | | | | |
| CheXpert [37] | 43.5±4.9 | 44.1±3.2 | 32.1±2.0 | 38.9±0.9 | 40.7±0.7 | 62.2±0.6 | 46.3±1.9 | 16.5±7.7 | 8.7±0.6 |
| BYOL [30] | 46.2±1.6 | 47.5±1.6 | 36.9±1.7 | 43.1±0.6 | 42.9±0.3 | 59.6±1.0 | 55.7±1.0 | 32.3±0.1 | 6.0±0.1 |
| SimCLR [9] | 44.9±2.9 | 41.4±3.7 | 33.0±0.0 | 42.6±0.4 | 39.2±0.7 | 61.9±0.8 | 54.3±1.0 | 33.2±0.1 | 13.3±0.5 |
| PixelPro [92] | 47.0±3.4 | 38.5±3.9 | 26.6±0.4 | 39.3±0.8 | 39.1±0.3 | **63.1±0.7** | 46.3±0.2 | 29.9±0.2 | 1.8±0.0 |
| ConVIRT [96] | 48.8±2.2 | 44.2±3.1 | 45.0±3.0 | 42.5±1.0 | 42.5±0.2 | 62.5±0.1 | 54.0±0.7 | 37.7±0.1 | 11.4±0.8 |
| CLIP [67]* | 49.3±2.0 | 46.5±2.3 | 46.2±0.3 | 42.8±1.5 | 42.5±0.6 | 62.9±0.8 | 55.5±2.1 | **39.0±0.0** | 12.5±1.0 |
| LoVT (Ours) | **49.5±1.3** | **49.2±4.6** | **49.2±0.2** | **43.4±0.7** | **43.1±0.6** | 61.0±1.3 | **55.8±1.1** | 37.6±0.2 | **13.4±0.8** |
| *Pre-Training on 100 % of frontal MIMIC-CXR* | | | | | | | | | |
| CheXpert [37] | 46.2±1.7 | 45.9±3.9 | 37.7±0.4 | 34.2±0.8 | 37.7±0.3 | 57.5±1.1 | 39.8±2.4 | 19.4±0.1 | 15.2±0.0 |
| BYOL [30] | 50.7±2.7 | 42.0±3.0 | 32.9±0.0 | 42.6±0.7 | 40.7±0.7 | 60.6±1.1 | 53.1±0.8 | 21.8±0.1 | 5.7±0.0 |
| SimCLR [9] | 48.1±2.5 | 44.1±2.1 | 35.3±0.0 | 41.2±0.8 | 38.7±0.5 | 61.1±0.7 | 48.7±0.5 | 30.0±0.0 | 11.8±0.0 |
| PixelPro [92] | 42.4±4.4 | 37.7±1.0 | 18.9±6.4 | 39.4±1.2 | 38.7±0.6 | **65.0±0.5** | 46.2±1.2 | 29.7±0.1 | 1.8±0.0 |
| ConVIRT [96] | 47.9±0.7 | 46.0±1.1 | 42.7±2.0 | 39.3±0.3 | 43.1±0.3 | 60.6±1.2 | 52.5±1.0 | 36.0±0.0 | **18.6±0.1** |
| CLIP [67]* | 48.6±2.4 | 45.8±4.1 | 41.7±0.1 | 44.0±0.7 | **45.0±0.5** | 62.8±0.5 | 56.9±1.4 | 39.4±0.0 | 11.4±0.8 |
| LoVT (Ours) | **51.2±2.5** | **46.2±2.4** | 44.0±0.8 | **44.1±0.3** | 43.9±0.7 | 62.1±0.5 | **57.4±0.5** | **39.9±0.0** | 9.4±0.5 |
| Task Nr. | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 |

\* Modified to use the same image and text encoders as ConVIRT and LoVT.

while the second-best method CLIP significantly outperforms all other methods only on 2 tasks. LoVT outperforms image-only methods (i.e. BYOL, SimCLR, and PixelPro) on 14 tasks, where the localized image-only method PixelPro outperforms LoVT only on one task (task 15). On 11 tasks LoVT outperforms other text-supervised methods (i.e. ConVIRT and CLIP), on 14 tasks it outperforms CheXpert classification and on all but two tasks it outperforms ImageNet initialization. When using 100% of the pre-training data LoVT is the best pre-training method on 11 tasks (better by at least the confidence interval on 5 tasks) and when using 30% on 11 tasks (significantly the best on 4 tasks). LoVT performs best on all COVID Rural tasks, best on most *Linear* tasks, and quite well on the *Frozen* protocol, but does not perform well on the NIH CXR dataset and when finetuned on the RSNA dataset. As there is no single method performing best on all tasks and LoVT performs best in the majority of tasks, this makes LoVT the default method of choice for localized downstream tasks.

*Relevance of Pre-Training Dataset Size* In our experiments we do not observe a consistent benefit of using roughly 210000 pre-training samples (i.e. 100% of the data) over using roughly 63000 samples (i.e. 30%). While on some datasets like RSNA and Object CXR many methods often perform better when pretrained on 210000 samples (100%), on other datasets like COVID Rural, methods often perform better when pre-trained on 63000 samples (30%). When comparing LoVT pre-trained on 30% of the data with other methods pre-trained in both settings (i.e. 30% and 100%), we observe that LoVT outperforms image-only
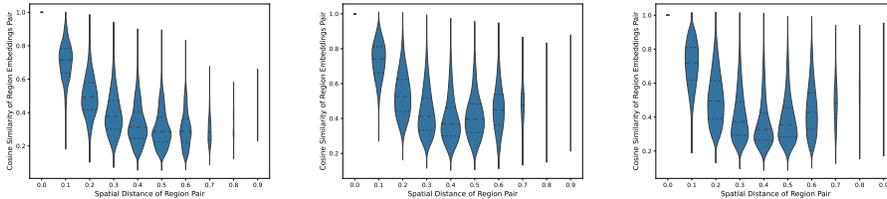
**Fig. 3.** Spatial smoothness and sensitivity of image region representations. **Left**: LoVT (Ours). **Middle**: No local losses. **Right**: No local losses and no attention pooling. Cosine similarities of image region pairs $\boldsymbol{y}_{i,k}^{\mathcal{I}}, \boldsymbol{y}_{i,k'}^{\mathcal{I}}$ (each from the same sample) plotted as violin plots (with their width representing the number of pairs and quartiles shown as dashed lines) over their spatial distance in the $7 \times 7$ image space (normalized and rounded to one decimal digit). We trained all models on 30% of the data and computed the representations on the test set.

methods (i.e. BYOL, SimCLR, and PixelPro) on 12 tasks, other text-supervised methods (i.e. ConVIRT and CLIP) on 7 tasks and CheXpert classification on 12 tasks, showing that LoVT effectively reduces the number of required pre-training samples.

*Relevance of Downstream Dataset Size* The results shown in Tab. 1 suggest that, as expected, larger downstream training sets lead to better results. However, we observe that for text-supervised methods (i.e. LoVT, CLIP, and ConVIRT), the downstream training set size is often less relevant compared to other methods. On the *RSAN YOLOv3 Frozen* tasks, LoVT (100%) outperforms ImageNet initialization by 31% when using 100% of the downstream samples, while it outperforms ImageNet initialization by even 167% when only using 1% of the samples.

*Spatial Smoothness and Sensitivity* We analyze the influence of the local losses and attention pooling on the spatial smoothness and sensitivity of image region representations and therefore plot in Fig. 3 the distributions of the cosine similarity of image region pairs over their spatial distances. For our LoVT model spatial smoothness and sensitivity can be observed as the quartiles and extreme points of the cosine similarity distributions decrease monotonously with increasing spatial distance, except for a few very distant region pairs with distances larger than 0.6. Note that these spatially very distant region pairs very likely represent opposite borders (or corners) of the image such that they both very likely contain background, explaining that they have more similar representations. Without local losses $\mathcal{L}_{\text{local-image}}$ and $\mathcal{L}_{\text{local-report}}$, the quartiles and extreme points decrease only for small spatial distances while increasing again for points further away, showing that spatial smoothness is only present for nearby regions and spatial sensitivity of more distant region is not optimal. When additionally replacing attention pooling with average (for image regions) and max (for

sentences) pooling, similar results can be observed except that the quartiles are decreasing faster and the maximum points do not decrease for nearby regions. We can therefore deduce that the local losses effectively encourage spatial smoothness and sensitivity while attention pooling alone has only little effect.

*Analysis of LoVT and Ablation Study* We refer to Appendix A for a detailed analysis of our method LoVT, including an ablation study (focusing on local weighting, global and local losses, and attention pooling), an analysis of the distribution and alignment of learned representations, and an analysis of the region weights $w_{i,k}^{\mathcal{I}}$.

## 5   Discussion

*Limitations of our Evaluation Procedure* We did not tune image encoder, downstream architectures, or preprocessing for downstream tasks, resized all inputs to only $224 \times 224$, and applied no data augmentation. Therefore, the presented downstream results are below results typically reported for these datasets. The evaluation procedure is kept simple to i) limit computational resources, ii) avoid bias induced by downstream tuning, and iii) allow for a fair comparison of pre-training methods, being the main purpose of this work. We assume that benefits observed in our evaluation procedure also indicate benefits for tuned real-world tasks, although they cannot be precisely quantified by our evaluation method.

*Limitations of LoVT* LoVT learns its alignment model implicitly based only on latent representations and instance-level pairing information. This makes the model sensitive to hyperparameters and hard to train. Also, it only uses local negatives from the same sample which restricts the number of negatives and may therefore limit its performance. Additionally, the alignment model is restricted to a simple attention mechanism and the regions are based on fixed patches that are not adaptive to the contents of the image. This may restrict the capabilities of the model and therefore of the pre-training method. For a detailed discussion of these limitations as well as of the potential negative societal impact we refer to Appendix C.

*Conclusion* We study pre-training for localized medical imaging on chest X-rays and propose a novel text-supervised method called LoVT, that combines instance-level contrastive learning with local contrastive learning. We evaluate our method on a novel evaluation framework consisting of 18 localized tasks on chest X-rays and compare it with typically used pre-training and initialization methods. While there is no single best method on all tasks, our method LoVT is the best method on 10 out of 18 studied tasks making it the method of choice for localized tasks.

We hope that our work provides valuable insights that encourage using pre-training for localized medical imaging and that our method inspires future work on localized text-supervised pre-training.

# References

1. Alsentzer, E., Murphy, J.R., Boag, W., Weng, W.H., Jin, D., Naumann, T., McDermott, M.B.A.: Publicly available clinical bert embeddings. In: ClinicalNLP (2019)
2. Bachman, P., Hjelm, R., Buchwalter, W.: Learning representations by maximizing mutual information across views. In: NeurIPS (2019)
3. Bardes, A., Ponce, J., LeCun, Y.: Vicreg: Variance-invariance-covariance regularization for self-supervised learning. In: ICLR (2022)
4. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. arXiv preprint arXiv: 1809.11096 (2019)
5. Caron, M., Bojanowski, P., Joulin, A., Douze, M.: Deep clustering for unsupervised learning of visual features. In: European Conference on Computer Vision (2018)
6. Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., Joulin, A.: Unsupervised learning of visual features by contrasting cluster assignments. In: NeurIPS (2020)
7. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: ICCV. pp. 9630–9640 (2021). https://doi.org/10.1109/ICCV48922.2021.00951
8. Chaitanya, K., Erdil, E., Karani, N., Konukoglu, E.: Contrastive learning of global and local features for medical image segmentation with limited annotations. In: NeurIPS (2020)
9. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: ICML (2020)
10. Chen, X., He, K.: Exploring simple siamese representation learning. In: CVPR. pp. 15745–15753 (2021). https://doi.org/10.1109/CVPR46437.2021.01549
11. Deng, C., Wu, Q., Wu, Q., Hu, F., Lyu, F., Tan, M.: Visual grounding via accumulated attention. In: CVPR (2018)
12. Desai, K., Johnson, J.: Virtex: Learning visual representations from textual annotations. In: CVPR. pp. 11157–11168 (2021). https://doi.org/10.1109/CVPR46437.2021.01101
13. Desai, S., Baghal, A., Wongsurawat, T., Al-Shukri, S., Gates, K., Farmer, P., Rutherford, M., Blake, G., Nolan, T., et al.: Data from chest imaging with clinical and genomic correlates representing a rural covid-19 positive population [data set]. The Cancer Imaging Archive (2020). https://doi.org/https://doi.org/10.7937/tcia.2020.py71-5978
14. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL. pp. 4171–4186 (2019). https://doi.org/10.18653/v1/N19-1423
15. Dinh, L., Krueger, D., Bengio, Y.: Nice: Non-linear independent components estimation. arXiv preprint arXiv: 1410.8516 (2015)
16. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real nvp. arXiv preprint arXiv: 1605.08803 (2017)
17. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: IEEE International Conference on Computer Vision. p. 1422–1430 (2015)
18. Donahue, J., Krähenbühl, P., Darrell, T.: Adversarial feature learning. arXiv preprint arXiv: 1605.09782 (2017)
19. Donahue, J., Simonyan, K.: Large scale adversarial representation learning. In: NIPS (2019)
20. Donahue, J., Hendricks, L.A., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: CVPR (2015)

21. Dosovitskiy, A., Springenberg, J.T., Riedmiller, M., Brox, T.: Discriminative unsupervised feature learning with convolutional neural networks. In: NIPS (2014)
22. Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., Courville, A.: Adversarially learned inference. arXiv preprint arXiv: 1606.00704 (2017)
23. Ermolov, A., Siarohin, A., Sangineto, E., Sebe, N.: Whitening for self-supervised representation learning. In: ICML. pp. 3015–3024 (2021)
24. Gazda, M., Plavka, J., Gazda, J., Drotár, P.: Self-supervised deep convolutional neural network for chest x-ray classification. IEEE Access pp. 151972–151982 (2021). https://doi.org/10.1109/ACCESS.2021.3125324
25. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. arXiv preprint arXiv: 1803.07728 (2018)
26. Goldberger, A., Amaral, L., Glass, L., Hausdorff, J., et al.: Physiobank, physiotoolkit, and physionet: Components of a new research resource for complex physiologic signals. Circulation [Online] **101**(23), 215—220 (2000)
27. Gomez, R., Gomez, L., Gibert, J., Karatzas, D.: Chapter 9 - self-supervised learning from web data for multimodal retrieval. In: Multimodal Scene Understanding (2019)
28. Goodfellow, I., Pouget-Abadie, J., et al.: Generative adversarial nets. In: NIPS (2014)
29. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In: CVPR (2017)
30. Grill, J.B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., Piot, B., kavukcuoglu, k., Munos, R., Valko, M.: Bootstrap your own latent - a new approach to self-supervised learning. In: NeurIPS (2020)
31. He, K., Fan, H., Wu, Y., et al.: Momentum contrast for unsupervised visual representation learning. In: CVPR. pp. 9726–9735 (2020). https://doi.org/10.1109/CVPR42600.2020.00975
32. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016). https://doi.org/10.1109/CVPR.2016.90
33. Hjelm, R.D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., Bengio, Y.: Learning deep representations by mutual information estimation and maximization. In: ICLR (2019)
34. Hu, R., Xu, H., Rohrbach, M., Feng, J., Saenko, K., Darrell, T.: Natural language object retrieval. In: CVPR (2016)
35. Huang, Z., Zeng, Z., Liu, B., Fu, D., Fu, J.: Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. arXiv preprint arXiv: 2004.00849 (2020)
36. Hénaff, O.J., Srinivas, A., et al.: Data-efficient image recognition with contrastive predictive coding. In: ICML. pp. 4182–4192 (2019)
37. Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., Seekins, J., Mong, D., Halabi, S., Sandberg, J., Jones, R., Larson, D., Langlotz, C., Patel, B., Lungren, M., Ng, A.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In: AAAI. pp. 590–597 (2019)
38. JF-Healthcare: object-cxr - automatic detection of foreign objects on chest x-rays. MIDL (2020), https://jfhealthcare.github.io/object-CXR/
39. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: ICML (2021)

40. Johnson, A., Lungren, M., Peng, Y., et al.: Mimic-cxr-jpg - chest radiographs with structured labels (version 2.0.0). PhysioNet (2019). https://doi.org/https://doi.org/10.13026/8360-t248

41. Johnson, A., Pollard, T., Berkowitz, S., et al.: Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. Sci Data **6**(317) (2019). https://doi.org/https://doi.org/10.1038/s41597-019-0322-0

42. Johnson, A., Pollard, T., Mark, R., Berkowitz, S., Horng, S.: Mimic-cxr database (version 2.0.0). PhysioNet (2019). https://doi.org/https://doi.org/10.13026/C2JT1Q

43. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: CVPR (2015)

44. Kingma, D.P., Dhariwal, P.: Glow: Generative flow with invertible 1x1 convolutions. In: NIPS. pp. 10215—-10224 (2018)

45. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2014)

46. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv: 1312.6114 (2014)

47. Li, J., Zhou, P., Xiong, C., Hoi, S.C.H.: Prototypical contrastive learning of unsupervised representations. In: ICLR (2021)

48. Liao, R., Moyer, D., Cha, M., et al.: Multimodal representation learning via maximization of local mutual information. In: MICCAI. pp. 273–283 (2021). https://doi.org/10.1007/978-3-030-87196-3_26

49. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV (2014)

50. Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., Tang, J.: Self-supervised learning: Generative or contrastive. IEEE Trans Knowl Data Eng (2021). https://doi.org/10.1109/TKDE.2021.3090866

51. Liu, Z., Stent, S., Li, J., Gideon, J., Han, S.: Loctex: Learning data-efficient visual representations from localized textual supervision. In: ICCV. pp. 2147–2156 (2021). https://doi.org/10.1109/ICCV48922.2021.00217

52. Loshchilov, I., Hutter, F.: Sgdr: stochastic gradient descent with warm restarts. In: ICLR (2017)

53. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2019)

54. Lu, J., Batra, D., Parikh, D., Lee, S.: Vilbert: pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In: NeurIPS (2019)

55. van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. Journal of Machine Learning Research **9**, 2579–2605 (2008)

56. Mahendran, A., Thewlis, J., Vedaldi, A.: Cross pixel optical-flow similarity for self-supervised learning. In: ACCV 2018 (2019)

57. Misra, I., van der Maaten, L.: Self-supervised learning of pretext-invariant representations. In: CVPR. pp. 6706–6716 (2020). https://doi.org/10.1109/CVPR42600.2020.00674

58. Müller, P., Kaissis, G., Zou, C., Rueckert, D.: Radiological reports improve pretraining for localized imaging tasks on chest x-rays. In: [to be published at] MICCAI (2022)

59. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European Conference on Computer Vision. p. 69–84 (2016)

60. van den Oord, A., Vinyals, O., et al.: Neural discrete representation learning. In: NIPS. pp. 6306—-6315 (2017)

61. den Oord, A.V., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al.: Conditional image generation with pixelcnn decoders. In: NIPS. pp. 4790—4798 (2016)
62. Oord, A.V., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. In: ICML. pp. 1747—-1756 (2016)
63. van den Oord, A., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv: 1807.03748 (2019)
64. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: NeurIPS (2019)
65. Pinheiro, P.O., Almahairi, A., Benmalek, R.Y., Golemo, F., Courville, A.: Unsupervised learning of dense visual representations. In: NeurIPS (2020)
66. Qi, P., Zhang, Y., Zhang, Y., Bolton, J., Manning, C.D.: Stanza: A Python natural language processing toolkit for many human languages. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations (2020), https://nlp.stanford.edu/pubs/qi2020stanza.pdf
67. 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. In: ICML. pp. 8748–8763 (2021)
68. Razavi, A., van den Oord, A., Vinyals, O.: Generating diverse high-fidelity images with vq-vae-2. In: NIPS. pp. 14837—-14847 (2019)
69. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. arXiv preprint arXiv: 1804.02767 (2018)
70. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241 (2015). https://doi.org/10.1007/978-3-319-24574-4_28
71. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. IJCV **115**(3), 211–252 (2015). https://doi.org/10.1007/s11263-015-0816-y
72. Saraf, V., Chavan, P., Jadhav, A.: Deep learning challenges in medical imaging. In: ICACTA. pp. 293–301 (2020). https://doi.org/10.1007/978-981-15-3242-9_28
73. Sariyildiz, M.B., Perez, J., Larlus, D.: Learning visual representations with caption annotations. In: ECCV. pp. 153–170 (2020). https://doi.org/10.1007/978-3-030-58598-3_10
74. Shih, G., Wu, C.C., Halabi, S.S., Kohli, M.D., Prevedello, L.M., Cook, T.S., Sharma, A., Amorosa, J.K., Arteaga, V., Galperin-Aizenberg, M., et al.: Augmenting the national institutes of health chest radiograph dataset with expert annotations of possiblepneumonia. Radiology: Artificial Intelligence **1** (2019). https://doi.org/https://doi.org/10.1148/ryai.2019180041
75. Society for Imaging Informatics in Medicine: Siim-acr pneumothorax segmentation. https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation (2019)
76. Sowrirajan, H., Yang, J., Ng, A.Y., Rajpurkar, P.: Moco pretraining improves representation and transferability of chest x-ray models. In: MIDL (2021)
77. Sriram, A., Muckley, M., Sinha, K., Shamout, F., Pineau, J., Geras, K.J., Azour, L., Aphinyanaphongs, Y., Yakubova, N., Moore, W.: Covid-19 prognosis via self-supervised representation learning and multi-image prediction. arXiv preprint arXiv: 2101.04909 (2021)

78. Su, W., et al.: Vl-bert: pre-training of generic visual-linguistic representations. In: ICLR (2020)
79. Sun, C., Myers, A., Vondrick, C., Murphy, K., Schmid, C.: Videobert: a joint model for video and language representation learning. In: ICCV (2019)
80. Tan, H., Bansal, M.: Lxmert: learning cross-modality encoder representations from transformers. In: EMNLP (2019)
81. Tang, H., Sun, N., Li, Y.: Segmentation model of the opacity regions in the chest x-rays of the covid-19 patients in the us rural areas and the application to the disease severity. medRxiv (2020). https://doi.org/https://doi.org/10.1101/2020.10.19.20215483
82. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: NIPS (2017)
83. Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.A.: Extracting and composing robust features with denoising autoencoders. In: ICML (2008)
84. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: CVPR (2015)
85. Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. In: CVPR (2016)
86. Wang, X., Peng, Y., Lu, L., et al.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: CVPR. pp. 3462–3471 (2017). https://doi.org/10.1109/CVPR.2017.369
87. Wang, X., Xu, Z., Tam, L., Yang, D., Xu, D.: Self-supervised image-text pre-training with mixed data in chest x-rays. arXiv preprint arXiv: 2103.16022 (2021)
88. Wang, X., Zhang, R., Shen, C., Kong, T., Li, L.: Dense contrastive learning for self-supervised visual pre-training. In: CVPR. pp. 3023–3032 (2021). https://doi.org/10.1109/CVPR46437.2021.00304
89. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: EMNLP (2020)
90. Wu, Z., Xiong, Y., Yu, S., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. In: CVPR. pp. 3733–3742 (2018). https://doi.org/10.1109/CVPR.2018.00393
91. Xie, E., Ding, J., Wang, W., Zhan, X., Xu, H., Sun, P., Li, Z., Luo, P.: Detco: Unsupervised contrastive learning for object detection. In: ICCV. pp. 8372–8381 (2021). https://doi.org/10.1109/ICCV48922.2021.00828
92. Xie, Z., Lin, Y., Zhang, Z., Cao, Y., Lin, S., Hu, H.: Propagate yourself: Exploring pixel-level consistency for unsupervised visual representation learning. In: CVPR. pp. 16679–16688 (2021). https://doi.org/10.1109/CVPR46437.2021.01641
93. Yakubovskiy, P.: Segmentation models pytorch. `https://github.com/qubvel/segmentation_models.pytorch` (2020)
94. Zbontar, J., Jing, L., Misra, I., LeCun, Y., Deny, S.: Barlow twins: Self-supervised learning via redundancy reduction. In: ICML (2021)
95. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: European conference on computer vision. p. 649–666 (2016)
96. Zhang, Y., Jiang, H., Miura, Y., Manning, C.D., Langlotz, C.P.: Contrastive learning of medical visual representations from paired images and text. arXiv preprint arXiv: 2010.00747 (2020)

97. Zhou, L., Palangi, H., Zhang, L., Hu, H., Corso, J., Gao, J.: Unified vision-language pre-training for image captioning and vqa. In: AAAI (2020)

**Table 3.** Results of the ablation study evaluated on the *RSNA YOLOv3 Frozen 10%* task. Results are averaged over five evaluation runs and the 95%-confidence interval is shown. The best results are highlighted in bold.

| Method | Global loss | Local losses | Local weighting | Pool | LR scheduler | RSNA YOLOv3 Frozen 10 % |
|---|---|---|---|---|---|---|
| LoVT (Ours) | ✓ | ✓ | ✓ | attention | cyclic-cosine | **17.9±0.4** |
| – | ✓ | ✓ | ✗ | attention | cyclic-cosine | 16.9±1.3 |
| – | ✓ | ✓ | ✗[1] | avg/max | cyclic-cosine | 15.7±0.4 |
| – | ✗ | ✓ | ✗[1] | – | cyclic-cosine | 12.3±0.7 |
| – | ✓ | non-smooth $\mathcal{L}_{\text{local-scan}}$ | ✓ | attention | cyclic-cosine | 15.6±1.4 |
| – | ✓ | $\mathcal{L}_{\text{local-report}}$ only | ✓ | attention | cyclic-cosine | 16.2±0.7 |
| – | ✓ | $\mathcal{L}_{\text{local-scan}}$ only | ✓ | attention | cyclic-cosine | 13.7±1.2 |
| – | ✓ | ✗ | – | attention | cyclic-cosine | 17.4±0.9 |
| – | ✓ | ✗ | – | avg/max | cyclic-cosine | 14.2±1.0 |
| CLIP[2] | single sentence | ✗ | – | attention | cyclic-cosine | 16.1±1.1 |
| – | single sentence | ✗ | – | avg/max | cyclic-cosine | 15.8±1.3 |
| ConVIRT | single sentence | ✗ | – | avg/max | reduce-on-plateau | 14.8±1.1 |

[1] Not realizable.
[2] Modified to use the same image and text encoders as ConVIRT and LoVT.

# A    Analysis and Ablation Study

In this section we analyze the relevance of different components of LoVT (as proposed in the main paper) and study its learned representations. In order to save computational resources, all further analysis and the ablation study are conducted on 30% of the pre-training data.

*Ablation Study* We conduct an ablation study to analyze the relevance of components of LoVT and the effects of the changes we made compared to ConVIRT. Focus of our ablation study are (i) the local weighting, (ii) the global and local losses, (iii) and attention pooling. We compare different model variants and their results on the *RSNA YOLOv3 Frozen 10%* task in Tab. 3. Note that we focus our ablation study on this single task as this is the task used for tuning all models and baselines while the other tasks are only used for the final evaluation (see Appendix E.3).

Starting from the unmodified LoVT model we first remove the local (region and sentence) weighting in the local losses, i.e. we use constant weights $w_{i,k}^{\mathcal{I}}$ and $w_{i,m}^{\mathcal{R}}$, and observe inferior results, showing the relevance of these weights. We then also remove attention pooling and replace it by average (avg) pooling for images and max pooling for reports. The performance further decreases highlighting the importance of attention pooling. Note that the local weights cannot be computed without attention pooling, making a model with local weighting but without attention pooling non-realizable. We further remove the global loss $\mathcal{L}_{\text{global}}$, i.e. set $\gamma = 0$ and only use the local losses without local weighting, and observe a large drop in downstream performance. We assume that this is caused by missing contrast between samples. Without the global loss, local weights can again not be computed, making a model with local weighting but without global loss non-realizable.

We also study the relevance of the local losses $\mathcal{L}_{\text{local-image}}$ and $\mathcal{L}_{\text{local-report}}$. Starting again from unmodified LoVT, we first adapt the local image loss $\mathcal{L}_{\text{local-image}}$

by redefining the positiveness score in a non-smooth way with $p_{k,l}^{\mathcal{I}} \propto \mathbb{1}_{[d_{\boldsymbol{x}}(k,l) \leq T]}$. The performance drops showing the relevance of the smoothness. When removing any of the local losses completely, i.e. either setting $\mu = 0$ or $\nu = 0$ and keeping only the global and one of the local losses ($\mathcal{L}_{\text{local-image}}$ or $\mathcal{L}_{\text{local-report}}$), the performance also drops compared to unmodified LoVT, showing that both local losses are required for optimal results. Note that removing $\mathcal{L}_{\text{local-report}}$ leads to a larger drop in downstream performance than removing only $\mathcal{L}_{\text{local-image}}$, indicating that $\mathcal{L}_{\text{local-report}}$ is more relevant for alignment. When both local losses are fully removed by setting $\mu = 0$ and $\nu = 0$, such that only the global loss remains, the performance slightly drops compared to unmodified LoVT showing that the local losses are relevant components of the model. However, the drop in performance is smaller than when removing only one of the local losses, which indicates that the symmetry of the local losses is essential for them to work. If we further replace attention pooling by avg/max pooling, a large drop in performance is observed, which again highlights the importance of attention pooling. Note that without avg/max pooling the local losses provide more improvements than when using attention pooling.

We also study the differences to ConVIRT[96] and (modified) CLIP[67]. Starting from ConVIRT, replacing their learning rate schedule (reducing on plateaus) by our cyclic cosine schedule (see Appendix E.1) improves the results. Further replacing their avg/max pooling (to compute global representations) by attention pooling improves the results even further. This settings corresponds to (modified) CLIP. In ConVIRT (and CLIP), only a single sentence per report is sampled when computing report representations. Replacing this sampling method by ours, where all sentences of a report are used to compute its representation, the results are improved if attention pooling is used. If no attention pooling is used, the performance degrades when using all sentences instead of a single randomly sampled one.

In our ablation study we highlighted the importance of all components of LoVT. We also showed that some of our proposed changes can also be used to improve the ConVIRT or CLIP models.

*Representation Distribution and Alignment Analysis* We analyze how representations are distributed and how well they are aligned. In Fig. 4 we show the standard deviation (std) of image $\bar{\boldsymbol{y}}_i^{\mathcal{I}}$ and image region $\boldsymbol{y}_{i,k}^{\mathcal{I}}$ representations of LoVT and variants of it without local losses or global loss. It can be observed that the (total) std of image region representations $\boldsymbol{y}_{i,k}^{\mathcal{I}}$ is similar in all three studied settings, indicating that the local and global losses have little influence on the overall variance of local representations. We additionally analyze the mean per-sample std and std of per-sample centroids of $\boldsymbol{y}_{i,k}^{\mathcal{I}}$ to study how representations are distributed within a sample and between different samples, respectively. The per-sample std of $\boldsymbol{y}_{i,k}^{\mathcal{I}}$ is smallest when only using the global loss (no local losses) and largest when only using the local losses (no global loss). Vice versa, the centroids std is largest when only using the global loss (no local losses) and smallest when only using the local losses (no global loss). Therefore, the local losses encourage the representations to differ within each
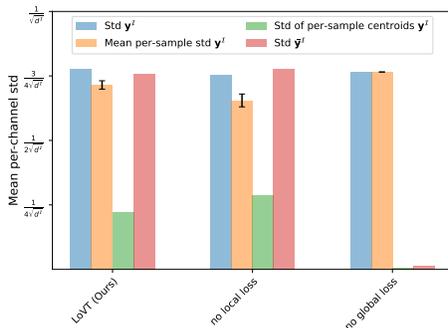
**Fig. 4.** Standard deviation (std) of image ($\bar{\boldsymbol{y}}_i^{\mathcal{I}}$) and image region ($\boldsymbol{y}_{i,k}^{\mathcal{I}}$) representations. **Left**: LoVT (Ours) **Middle**: No local losses. **Right**: No global loss. For $\boldsymbol{y}_{i,k}^{\mathcal{I}}$ we additionally show the mean std per-sample, i.e. how different representations are within a sample, and the std of the per-sample centroids. The models were trained on 30% of frontal MIMIC-CXR and then evaluated on the whole test set.

sample, i.e. ensure spatial sensitivity, while the global loss encourages them to differ between samples, i.e. prevents the collapse of per-image representations to a constant vector. The std of global image representations $\bar{\boldsymbol{y}}_i^{\mathcal{I}}$ behaves similarly to the centroids std of $\boldsymbol{y}_{i,k}^{\mathcal{I}}$, except that the local losses have only little influence on it. Note that the centroids std and std of global image representations almost completely vanish without the global loss, while there is still notable per-sample std present without the local losses. This highlights the importance of the global loss for preventing the collapse of representations.

In Fig. 5 we plot the alignment quality of local representations, i.e. the $\ell_2$-distance of report-to-image ($\boldsymbol{z}_{i,k}^{\mathcal{R}\to\mathcal{I}}$) with their related image region representations ($\boldsymbol{z}_{i,k}^{\mathcal{I}}$) and of image-to-report ($\boldsymbol{z}_{i,m}^{\mathcal{I}\to\mathcal{R}}$) with report sentence representations ($\boldsymbol{z}_{i,m}^{\mathcal{R}}$). We compare the representations learned by LoVT with (default) and without global loss. It can be observed that in both cases the image-to-report representations are better aligned than the report-to-image representations. This can be expected, as most of the information contained in the report is based on the image, making it easy to compute sentence representations from image region representations, while images contain additional details not described by the reports, making it harder to compute report-to-image representations. Both, report-to-image and image-to-report representations, are slightly better aligned when the global loss is used additionally to the local losses during training (as in the unmodified LoVT model). One can therefore deduce that the global loss supports the learning of local representations.

In Fig. 6 we plot a t-SNE[55] projection of local representations learned by LoVT. Sentence representations are similarly distributed to image-to-report representations confirming, as already observed in the alignment analysis, that the model is able to align both distributions. Only one cluster of sentence representations is separated from image-to-report representations. We assume that these
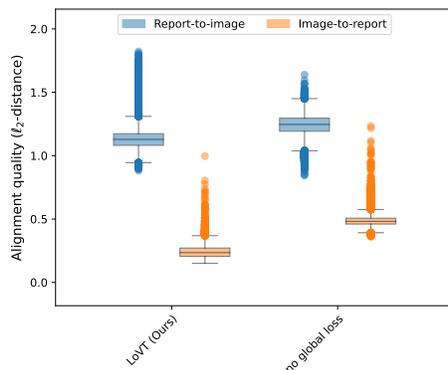
**Fig. 5.** Alignment quality of local representations. **Left**: LoVT (Ours) **Right**: No global loss. Measured by the $\ell_2$-distance of uni-modal with their related cross-modal representations. **Blue**: Report-to-image ($z_{i,k}^{\mathcal{R}\to\mathcal{I}}$) with image region ($z_{i,k}^{\mathcal{I}}$) representations. **Orange**: Image-to-report ($z_{i,m}^{\mathcal{I}\to\mathcal{R}}$) with report sentence ($z_{i,m}^{\mathcal{R}}$) representations. All representations are $\ell_2$-normalized before the distance is computed. The models were trained on 30% of frontal MIMIC-CXR and then evaluated on the whole test set.

are sentences that do not describe features present in the image, e.g. describing features from lateral views or differences to previous studies of the patient. Image region representations and report-to-image representations are distributed differently, which again confirms that these could not be fully aligned. In the t-SNE[55] projection the image region representations are split into many clusters. We assume that this is a result of the negatives in the local image loss encouraging contrast between (spatially distant) region representations of each sample, such that our model behaves similarly to a clustering algorithm.

To further study the effect of the local losses, we plot t-SNE[55] projections of image region representations from samples of the RSNA pneunomia detection dataset in Fig. 7. We compare the representations learned using our unmodified LoVT model, our model without local losses, and our model without local losses and attention pooling. For the unmodified LoVT model, image region clusters can again be observed while such clusters cannot be observed without the local losses. This confirms our assumption that these clusters are a result of the local losses. It can also be observed that in the unmodified LoVT model the representations of pneumonia positive regions are distributed in a very confined area of space and are therefore probably easily separable from non-pneumonia regions. Without local losses the positive region representations are more spread over the space making them harder to separate. If attention pooling is not used as well, the positive region representations are distributed around multiple areas in space which may also hurt separability. Therefore, using the local losses and attention pooling improves separability of downstream representations which is confirmed by the results shown in Tab. 3.
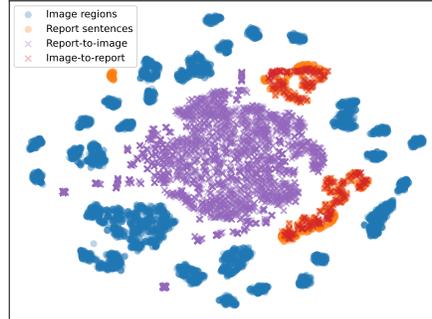
**Fig. 6.** t-SNE[55] plot of projected local uni-modal representations (points) and the aligned cross-modal representations (crosses). **Blue**: Image regions $(z_{i,k}^{\mathcal{I}})$. **Orange**: Report sentences $(z_{i,m}^{\mathcal{R}})$. **Purple**: Report-to-image $(z_{i,k}^{\mathcal{R}\to\mathcal{I}})$. **Red**: Image-to-report $(z_{i,m}^{\mathcal{I}\to\mathcal{R}})$. We trained our model on 30% of frontal MIMIC-CXR and computed the representations on the first 100 samples of the test set.
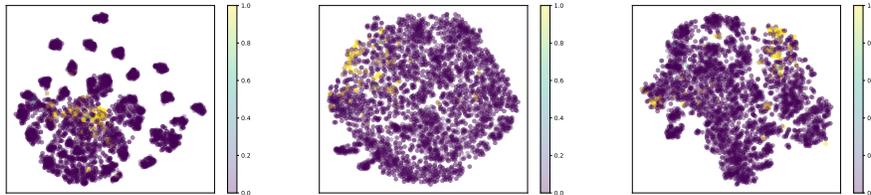


**Fig. 7.** t-SNE[55] plots of image region representations from samples of the RSNA pneumonia detection dataset. The color of each point indicates the overlap of the related region with a pneumonia opacity region. **Left**: LoVT (Ours). **Middle**: No local losses. **Right**: No local losses and no attention pooling. We trained all models on 30% of frontal MIMIC-CXR and computed the representations on the first 100 samples of the RSNA test set.

*Local Weights* In order to understand how the weighting of image regions works, we study how the region weights $w_{i,k}^{\mathcal{I}}$ are distributed. In Fig. 8 we therefore plot the mean region weights on the mean image of the pre-training test set. The weights are distributed horizontally symmetrically around the center of the images and most focus is on the lungs and around the spine. This indicates that the weighting works as expected, as most pathologies in a frontal chest X-ray are typically observed at lungs or heart.

## B   Discussion of the Limitations of LoVT

*Weak Supervision and Sensitivity to Hyperparameters* As no supervision for the alignment of image regions and report sentences is available, we implicitly
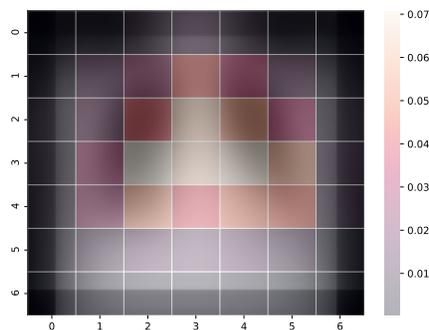
**Fig. 8.** Local image region weights $w_{i,k}^{\mathcal{I}}$ of different regions averaged over all samples and plotted on top of the mean image. The model was trained on 30% of frontal MIMIC-CXR. Weights and mean image were computed on the whole test set.

learn an alignment model in the latent representation space. We jointly learn this alignment model and the latent representations of image regions and report sentences, having only the global alignment information of image-report pairs as supervision. Therefore, we suspect that the model tends to converge to local optima, explaining its sensitivity to hyperparameters, especially to the learning rate. While using the cyclic-cosine learning rate schedule helps, our method is still hard to train and tune. We leave studying more explicit supervision, e.g. by including generative losses, to future work.

*Limited Negatives for Local Alignment* We only use local negatives from the same sample. By design, the number of local negatives is therefore very limited and many of these negatives may be very easy. This may limit the model performance on downstream tasks[9,31]. In preliminary experiments we also included negatives from other samples but could not observe a benefit. We leave the study of losses with more negatives (e.g. based on the MoCo[31] approach) or without explicit negatives (e.g. based on the BYOL[30] approach) to future work.

*Limited Alignment Model* We decided to use single-head scaled dot product attention with linear projections as our alignment model. While this keeps the alignment model simple and computationally cheap, it also limits its capabilities. We leave studying more complex alignment models, like multi-head attention or (one or more) transformer[82] layers, to future work.

*Non-Adaptive Regions* In LoVT the image region representations are computed for fixed regions, i.e. patches. Their boundaries are arbitrary and relevant features may therefore be spread across regions or multiple neighboring features may be in the same region, making it hard to learn region representations. We leave studying other techniques for finding content-based regions and computing their representations to future work.

## C    Discussion of Negative Societal Impact

In this section we discuss the possible negative societal impact of our work. We identified three primary aspects: i) usage of our method in medical applications, ii) data privacy issues, and iii) energy consumption.

*Usage in Medical Applications* As our method is targeted towards medical applications, potential issues in our method may lead to harm through misdiagnosis. Most of the potential issues, including interpretability and reliability issues, are not specific to our method but apply to most deep learning methods in medicine and we therefore do not discuss them here. Still, we identified another potential issue: data bias learned during pre-training. While bias from data may be learned by most machine learning methods, in our case the bias might be learned from both, pre-training and fine-tuning data. During fine-tuning the pre-training dataset might not be available making it hard to identify such a data bias. As possible mitigation strategies the pre-training dataset (if available) should also be analyzed for data bias or the fine-tuned model should be tested for learned bias before using it in medical applications. Note that this issue applies to most transfer learning approaches including other pre-training methods.

*Data Privacy Issues* Information learned from the pre-training dataset is contained in the pre-trained model weights, which may include information about individuals in the dataset. When distributing such models to make them available for others to fine-tune, this information is distributed as well. If the pre-training dataset is non-public but the pre-trained model is made publicly available, this may lead to data leakage and therefore a privacy breach. This is especially problematic if individuals can be re-identified. Therefore, pre-trained models should be distributed only under the same conditions as the pre-training dataset or other precaution measures, like privacy-preserving machine learning, should be taken to prevent data leakage. We thus decided not to release our pre-trained model weights publicly. Note that this issue applies to most transfer learning approaches including other pre-training methods.

*Energy Consumption and Environmental Impact* Training of deep learning models consumes substantial amounts of energy and may therefore have environmental impacts. In our experiments, we observe that pre-training (including LoVT and the baselines) typically takes 0.5-2 days while downstream tasks typically only train for minutes to a few hours per run. While we did not study the exact energy consumption, we use the training times as an estimate and conclude that the energy consumption, and therefore the environmental impact, during pre-training is substantially higher than during finetuning. We observe that in our setting most studied pre-training methods, including LoVT, have similar runtimes (1-2 days, depending on the exact hyperparameters, for training on the full dataset) and only CheXpert pre-training is significantly faster (typically taking 0.5-1 days for pre-training). Besides training, also hyperparameter tuning

needs to be considered, which may be required when applying LoVT to another pre-training dataset.

While, as we observed, the high energy consumption of pre-training is an effect that is general to many methods, it should still be considered when deciding whether and how to use LoVT. An approach to reduce the energy consumption is to limit the hyperparameter tuning of LoVT on the pre-training dataset (e.g. only tune the learning rate) and use the defaults from our paper for most hyperparameters, although tuning other hyperparameters may improve downstream results. Instead, hyperparameter tuning could be more focused on the finetuning of the model. Additionally, pre-trained models should be made publicly available where this does not lead to privacy issues.

## D      Detailed Discussion of Related Work

### D.1     Self-supervised Representation Learning

State-of-the-art methods for pre-training image models using self-supervised representation learning can be categorized into *generative* and *discriminative* approaches. Generative models learn a distribution over the training images and a latent representation space. Typically, these approaches are autoencoding models[83,46,60,68], which learn to reconstruct the input image (or parts of it), adversarial models[28,18,22,4,19], where data and representation are modeled jointly, autoregressive models[62,61], where image regions are conditioned on previous image regions, or flow-based models[15,16,44], which estimate high-dimensional densities from data. Generative models can recover the original data distribution without the need for assumptions on downstream tasks and are therefore well-suited for a wide range of applications, most notably for generative tasks[50]. However, they have some inherent problems, most notably, they model the distribution in the data space (e.g. in pixel-space) and therefore focus too much on low-level details (like pixels) instead of encouraging high-level abstractions that are typically required for discriminative downstream tasks like classification[50].

Discriminative approaches are better suited for such tasks as they define discriminative objectives based on pretext tasks created from the unlabeled data. Early discriminative approaches relied on heuristics when defining the pretext tasks[21,17,95,59,25], limiting the generality of the resulting representations. In recent years, contrastive approaches[90,63,36,33,57,47,9,31,30,10,94,3,23], have become the state-of-the-art discriminative approaches for self-supervised representation learning. Contrastive methods act in the representation space and try to align representations of similar images (e.g. different views from the same image) while spreading representations of different images. Clustering approaches like DeepCluster[5] also belong to the contrastive approaches. Discriminative approaches are in general very lightweight and contrastive methods are currently state-of-the art for discriminative downstream tasks. However, they are not suited for generative tasks and many aspects, like the need for negative sampling, are not well-understood yet although being tackled by approaches like

BYOL[30], SimSiam[10], BarlowTwins[94], and VICReg[3]. Contrastive methods have also been successfully applied to medical imaging including image classification on chest X-rays[24,76,77].

Most contrastive learning approaches use instance-level contrast, i.e. represent each view of the image by a single vector. While the resulting representations are well-suited for global downstream tasks like image classification they lack properties like spatial sensitivity or smoothness required for more localized downstream tasks (like segmentation or detection)[92]. Therefore, there is a number of recent approaches that use region-level contrast[92,91,88,8,65,56], i.e. they act on representations of image regions. These approaches are more suited for localized tasks and therefore typically outperform instance-level methods on such tasks.

## D.2  Multimodal Representation Learning

While self-supervised representation learning on a single modality (e.g. images) already achieves great results, in some settings more modalities are available. Utilizing such additional modalities can improve the downstream results as additional information is available that can be utilized during representation learning. One form of such additional modalities is text, that often accompanies images in the form of captions or linked reports. Early works on combining image and text modalities did not focus on pre-training for downstream tasks but instead on learning aligned representations for cross-modal retrieval[85,27], on encoder-decoder tasks like image captioning[84,43,20], and on joint prediction tasks like visual question answering[29] and visual grounding[11,34]. In recent years, several works utilized transformer[82] models to compute joint representations of image content and text[79,54,78,80,97,35]. They pre-trained their models using self-supervised tasks like multi-modal alignment prediction on paired image-text datasets and then finetuned them on multi-modal downstream tasks like image retrieval or visual question answering. While these methods can effectively pre-train joint image-text models, these models cannot be used for image-only downstream tasks. Recently, there is much focus on self-supervised representations learning methods that pre-train image models for downstream tasks by taking advantage of the companion text[67,39,96,12,73]. VirTex[12] and ICMLM[73] use image captioning tasks (generative tasks), ConVIRT[96], CLIP[67] and ALIGN[39] use multiview contrastive learning[2] (contrastive tasks). In[87] generative and contrastive losses are combined to train on mixed chest X-ray data, i.e. where only for some images paired reports is available. LocTex[51] does localized pre-training on natural images with companion text using a dot product based model to predict alignment of text and image regions. Unlike our method it uses supervision generated by mouse gazes instead of learning the alignment implicitly using a local contrastive loss. Most related to our work is the recently published local-mi[48] that does contrastive learning on report sentences and image regions but aligns each sentence with its most related region instead of using an alignment model like our method. Also, it targets classifi-

cation instead of localized tasks and does therefore neither encourage contrast between regions nor spatial smoothness.

# E   Experiment Details

In all our experiments we use PyTorch[64] Version 1.10 (BSD-style license[5]) and train on a single NVIDIA Quadro RTX 8000.

## E.1   Pre-Training

*Pre-training Data and Pre-Processing* Our method can be used with any dataset containing pairs of medical images and reports supposing that the reports contain multiple sentences and the sentences in the reports provide a semantically useful description of the contents in the image. We use version 2 of MIMIC-CXR[42,41,40,26] as, to our best knowledge, it is the largest and most commonly used dataset of this kind conating more than $200,000$ imaging studies, each with one or more frontal or lateral chest X-rays and one semi-structured free-text radiology report, written by a practicing radiologist during routine clinical care, describing radiological findings of the images.

We download the already pre-processed images from its JPG-version[6] and remove all except the frontal views, i.e. we only keep the *antero-posterior (AP)* and *postero-anterior (PA)* views. We download the reports from MIMIC-CXR[7] and extract the text from the *Findings* and *Impression* sections. Reports containing none of these sections are removed. For each report we concatenate the extracted text from both sections and remove reports where the extracted text contains less than three tokens (based on tokenizing it using Stanza[66]). We split the extracted text into sentences using Stanza[66] again. Finally, we remove all samples that contain no images or no report (after the previous steps) and then apply the training/validation/test splits provided by MIMIC-CXR-JPG[40] such that we have 210228/1712/2867 training/validation/test samples (i.e. studies with one report and one or more images each), respectively.

*Encoders and Model Details* In the image encoder we use the ResNet50 implementation from Torchvision[8] and initialize it with ImageNet[71] weights[9]. In the report encoder we use the BERT_base PyTorch implementation from Hugging-face Transformers[89][10] and initialize it with weights from ClinicalBERT[1][11] that was trained on clinical notes.

---

[5] `https://github.com/pytorch/pytorch/blob/master/LICENSE`
[6] `https://physionet.org/content/mimic-cxr-jpg/2.0.0/` (PhysioNet Credentialed Health Data License)
[7] `https://physionet.org/content/mimic-cxr/2.0.0/` (PhysioNet Credentialed Health Data License)
[8] `https://github.com/pytorch/vision` (BSD 3-Clause License)
[9] `https://pytorch.org/hub/pytorch_vision_resnet/` (BSD 3-Clause License)
[10] `https://github.com/huggingface/transformers` (Apache-2.0 License)
[11] `https://huggingface.co/emilyalsentzer/Bio_ClinicalBERT` (MIT License )

We model the nonlinear transformations $f^{\mathcal{I}}$, $f^{\mathcal{R}}$, $\bar{f}^{\mathcal{I}}$, and $\bar{f}^{\mathcal{R}}$ as shallow neural networks without shared parameters, consisting of a (element-wise) linear layer with output size 2048, batch norm, and ReLU followed by another linear layer with output size 512. This follows previous works[9,30,96] except for the batch norm which we found beneficial.

*Data Augmentation* For image augmentations we first randomly sample one of the frontal chest X-rays of the sample (i.e. study) and then follow the augmentation scheme of ConVIRT[96], i.e. random cropping (resized to $224 \times 224$), horizontal flipping, affine transformations, contrast and brightness jittering and Gaussian blur. We also tried removing geometric augmentations but found this setting to perform worse. For text augmentations we concatenate all sentences of the Findings and Assessment sections of the report in the sample but randomly change the order by swapping pairs of sentences with a probability of 0.6. We also tried randomly removing or duplicating sentences but did not find it to be beneficial.

*Training Details and Cyclic Cosine Learning Rate Schedule* For pre-training, we experimented with different learning rate schedules and found that a cyclic cosine learning rate schedule[52] where the restarts also follow the cosine function and with a cycle length of two epochs (i.e. one decreasing and one increasing epoch) is beneficial. As both modalities have different properties (e.g. type of contained information) and the encoders have different architectures, they may converge at different speeds making it hard for both to adapt to each other. Therefore, we assume that the decreasing phase of the schedule allows the encoders to catch up and adapt to each other while the increasing phase allows them to learn faster and escape local optima. We use the AdamW[53] optimizer with a batch size of 32, with 16 gradient accumulation steps, an initial learning rate of $1 \times 10^{-4}$, and weight decay $1 \times 10^{-6}$ and train until the validation loss does not decrease for 10 consecutive epochs.

*Hyperparameter Tuning* We tune the hyperparameters of LoVT using only the *RSNA YOLOv3 Frozen 10%* task (see Appendix E.3). For the hyperparameters of the global loss, i.e. $\tau$ and $\lambda$, and for the global representation dimension $\bar{d}^{\mathcal{Z}}$ we use the default values from ConVIRT[96]. In preliminary experiments, we tried different values for the local representation dimension $d^{\mathcal{Z}}$ but did not find that small changes to it have significant influence on the results, and therefore set $d^{\mathcal{Z}} = \bar{d}^{\mathcal{Z}}$ (i.e. 512). We determined the hyperparameters of the local losses, i.e. $\tau'$, $\beta$, and $T$, in preliminary experiments including grid searches and manual tuning. The loss weights $\gamma$, $\mu$, and $\nu$ were determined by running small grid searches in the following way: We first set $\gamma = \mu = \nu = 1$ and run a grid search to balance the local loss weights $\mu$ and $\nu$ while keeping $\gamma$ fixed, i.e. trying ($\mu = 0.5, \nu = 1.5$), ($\mu = 1.0, \nu = 1.0$), and ($\mu = 1.5, \nu = 0.5$). After we found that ($\mu = 1.0, \nu = 1.0$) performs best, we ran a grid search to balance local and global losses while keeping $\mu = \nu$, i.e. trying ($\gamma = 0.75, \mu = 1.0, \nu = 1.0$), ($\gamma = 1.0, \mu = 1.0, \nu = 1.0$),

$(\gamma = 1.0, \mu = 0.75, \nu = 0.75)$, and $(\gamma = 1.0, \mu = 0.25, \nu = 0.25)$. We found that $(\gamma = 1.0, \mu = 0.75, \nu = 0.75)$ performs best.

All hyperparameters except the learning rate are tuned using 30% of the pre-training dataset and we tune the learning rate individually on 30% and 100% of the pre-training data. Note that we also slightly tune the learning rate when tuning other hyperparameters and in our ablation study.

### E.2   Baselines

*Random and ImageNet Init.* For random initialization we do not pre-train the ResNet50 backbone but instead initialize it randomly following its default initialization. For the ImageNet initialization we use the weights[9] provided by Torchvision.

*CheXpert* We train the ResNet50 backbone using multi-label binary classification on MIMI-CXR. We use five CheXpert[37] labels (Cardiomegaly, Edema, Consolidation, Atelectasis, and Pleural Effusion), which are included in the MIMIC-CXR-JPG[40] dataset, and convert them to binary labels following the *U-Ones* mapping[37] (i.e. mapping all uncertain labels to positive labels). During CheXpert pre-training we use the full ResNet50 model including the average pooling and the fully connected (FC) layer but throw away the latter two for downstream tasks. All layers except the FC layer are initialized from ImageNet weights[9] and we randomly initialize the FC layer such that it has an output dimension of five (matching the number of classes). We use the sigmoid activation on the outputs, multi-label binary cross-entropy loss and the Adam[45] optimizer and train with batch size 64 and weight decay $1 \times 10^{-6}$ until the validation *Area Under Receiver Operating Characteristic (AUROC)* does not increase for 10 consecutive epochs after which we select the checkpoint with the best validation AUROC. We tuned the initial learning rate and set it to $3 \times 10^{-4}$ ($1 \times 10^{-4}$ when trained on 30% of the data). If the validation AUROC does not increase for three consecutive epochs we multiply the current learning rate by 0.5.

*SimCLR[9]* We use the PyTorch implementation available at `https://github.com/spijkervet/SimCLR` (MIT License) with the default image augmentations from the paper (with images resized to $224 \times 224$) except for color jittering where we do not adjust saturation and hue due to the monochrome nature of chest X-rays. Following [96] we set the output dimension to 128 and hidden size to 4096, use batch size 128, and weight decay $1 \times 10^{-4}$  For training, we use the Adam[45] optimizer and the cosine decay learning rate schedule[52], without restarts, over 100 epochs, with a single warm-up epoch. We tuned the initial learning rate and set it to $3 \times 10^{-4}$

*BYOL[30]* We use the PyTorch implementation available at `https://github.com/lucidrains/byol-pytorch` (MIT License) with the default image augmentations from the paper (with images resized to $224 \times 224$) except for color jittering where we do not adjust saturation and hue due to the monochrome nature of

chest X-rays. We set the output dimension to 128 and hidden size to 4096, use decay rate 0.99, batch size 64, and weight decay $1 \times 10^{-4}$  For training, we use the Adam[45] optimizer and the cosine decay learning rate schedule[52], without restarts, over 100 epochs, with a single warm-up epoch. We tuned the initial learning rate and set it to $1 \times 10^{-4}$ ($3 \times 10^{-5}$ when trained on 30% of the data).

*PixelPro[92]* We use the PyTorch implementation available at `https://github.com/lucidrains/pixel-level-contrastive-learning` (MIT License) with the default image augmentations from the paper (with images resized to $224 \times 224$) except for color jittering where we do not adjust saturation and hue due to the monochrome nature of chest X-rays.

We set the output dimension to 512 and hidden size to 2048, use batch size 64, and weight decay $1 \times 10^{-5}$  For training, we use the Adam[45] optimizer and the cosine decay learning rate schedule[52], without restarts, over 100 epochs, with a single warm-up epoch. We tuned the initial learning rate and set it to $1 \times 10^{-3}$

*ConVIRT[96]* We use our own implementation of ConVIRT (as the general framework of ConVIRT is similar to LoVT) and train until the validation loss does not decrease for 15 consecutive epochs after which we use the checkpoint with the lowest validation loss. We tuned the learning rate and set it to $1 \times 10^{-4}$ ($1 \times 10^{-5}$ when trained on 30% of the data). If the validation loss does not decrease for 12 consecutive epochs we multiply the current learning rate by 0.5. We use the default values from the paper for all other hyperparameters.

*CLIP[67]* We use our own implementation of CLIP (as the general framework of CLIP is similar to LoVT). For better comparability with LoVT and the other baselines, we use ResNet50 and BERT_base as encoders. Following the framework of CLIP we only encode single sentences and therefore randomly sample a sentence from the report (as in ConVIRT). We use the AdamW[53] optimizer with a batch size of 32 (the same as used in ConVIRT and LoVT), with 16 gradient accumulation steps, and the cyclic cosine learning rate scheduler (as in LoVT) with an initial learning rate of $1 \times 10^{-4}$ , and weight decay $1 \times 10^{-6}$ and train until the validation loss does not decrease for 10 consecutive epochs.

*Batch Sizes of the Baselines* Most contrastive learning methods are very sensitive to the used batch size, therefore the batch size is an important hyperparameter when comparing such methods. However, increasing the batch size also increases the GPU memory consumption and different methods have different memory requirements, such that using the same batch size for all methods does not allow for a fair comparison as in practice available GPU memory is typically limited. We therefore decided to use three different batch sizes: The smallest batch size ($32^{12}$) is used for all text-supervised methods (i.e. ConVIRT, CLIP

---

[12] We use this batch size as it was used in ConVIRT and as memory requirements are then kept below 24GB allowing training on widely used GPUs.

and our LoVT) as they require much memory due to their language model and they are also less sensitive to the batch size[96]. For image-only methods with a momentum encoder (i.e. BYOL and PixelPro) we use a larger batch size (64) and for SimCLR we further increase the batch size (128) as it does not have a momentum encoder and is very sensitive to the used batch size.

### E.3  Downstream Evaluation

*Datasets*

- **RSNA Pneumonia Detection**[86,74] (Licensed following the competition rules[13]): We download the dataset from its Kaggle page[14] but use only their training set which we randomly split into our training, validation, and test set resulting in 16010/5337/5337 training/validation/test samples, respectively. For each sample we compute a segmentation mask (used in the *Linear* evaluation) from all the ground truth detection boxes of that sample.
- **COVID Rural**[81,13] (TCIA Data Usage Policy and CC BY 4.0 License): We download the dataset from its Github repository[15] and randomly split it into training, validation, and test set of sizes 133/44/44, respectively.
- **SIIM-ACR Pneumothorax Segmentation**[75] (Licensed following the competition rules[16]): We download the dataset from its Kaggle re-upload[17], which is officially recommended on the original challenge website[18], but use only their training set which we randomly split into our training, validation, and test set resulting in 7229/2409/2409 training/validation/test samples, respectively.
- **Object CXR**[38] (CC BY-NC 4.0 License): We download the dataset from a re-upload[19] as it is no longer available at its original source[20]. We randomly split their training set into our training and validation set and use their development set as our test set such that we have 6400/1600/1000 training/validation/test samples, respectively. For each sample we compute a segmentation mask (used in the *Linear* evaluation) from all the ground truth detection boxes of that sample.
- **NIH CXR**[86] (Licensed for public use with attribution[21]): We download the ChestX-ray8 dataset provided by the NIH Clinical Center from its official website[22] but use only the samples where bounding boxes are provided as

---

[13] https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/rules
[14] https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/data
[15] https://github.com/haimingt/opacity_segmentation_covid_chest_X_ray/tree/master/covid_rural_annot
[16] https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation/rules
[17] https://www.kaggle.com/seesee/siim-train-test/
[18] https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation/overview/siim-cloud-healthcare-api-tutorial
[19] https://academictorrents.com/details/fdc91f11d7010f7259a05403fc9d00079a09f5d5
[20] https://jfhealthcare.github.io/object-CXR/
[21] https://nihcc.app.box.com/v/ChestXray-NIHCC/file/249502714403
[22] https://nihcc.app.box.com/v/ChestXray-NIHCC/

ground truth. We randomly split these samples into our training, validation, and test set such that we have 588/196/196 training/validation/test samples, respectively.

### E.4 Evaluation Protocols

In this section we describe the details of the evaluation protocols used in the evaluation framework[58], including downstream model architectures and training details. Note that we do not use image augmentations in any of the evaluation protocols but resize and pad the input images to size $224 \times 224$.

*U-Net Finetune* We do not use the original UNet[70] architecture but instead build a UNet-like model[23] based on the pre-trained ResNet50 model. Therefore, we use the ResNet50 (except its avg pooling and FC layer) as the contracting path (left side) of our UNet-like model. The last feature map of ResNet50 has a size of $7 \times 7$ and dimension 2048. Here we add two more convolutional blocks (each with a $3 \times 3$ convolution followed by batchnorm and ReLU) with output dimension 2048 to the contracting path. For the expansive path (right side) we closely follow the architecture of the original UNet but use five instead of four upsampling blocks (each with $2 \times 2$ transposed convolution, concatenation, and two $3 \times 3$ convolutions each followed by batchnorm and ReLU) which have output dimensions 1024, 512, 256, 128, and 64, respectively. For concatenation the ResNet50 blocks conv4, conv3, conv2, conv1, and the input image are used. We add a single $1 \times 1$ convolution that predicts the positive class scores and use the binary Dice loss from the Segmentation Models Pytorch library[93].

For training, we use the Adam[45] optimizer with weight decay $1 \times 10^{-6}$ The learning rate is tuned individually for each model and task based on the best validation Dice[24]. The learning rate is multiplied by 0.5 if the validation Dice does not decrease for three consecutive epochs. We use a warmup period in which we do not train the ResNet50 backbone but only the other, randomly initialized, layers with a learning rate of $1 \times 10^{-3}$ after which we train the whole model (including the ResNet50). On the COVID Rural dataset we use batch size eight, a warmup period of 20 iterations and do early stopping (based on validation Dice) after 20 epochs. On the SIIM-ACR Pneumothorax dataset we use batch size 64, a warmup period of 100 iterations and do early stopping after 10 epochs. Finally we report the test Dice of the epoch with the best validation Dice.

*U-Net Frozen* We use the same architecture and loss function as in the *U-Net Finetune* protocol but freeze the pre-trained ResNet50 weights and never train

---

[23] Our implementation is based on `https://github.com/kevinlu1211/pytorch-unet-resnet-50-encoder/blob/master/u_net_resnet_50_encoder.py` (MIT License)

[24] We use the micro-averaged Dice score based on this implementation: `https://torchmetrics.readthedocs.io/en/latest/references/modules.html#f1` (Apache-2.0 License)

them. Instead we only train the other layers using the same hyperparameters as in the *U-Net Finetune* protocol (except for the warmup period which is not relevant in this setting).

*Linear* We use the frozen pre-trained ResNet50 (except for its avg pooling and FC layer) to compute $7 \times 7$ feature maps. A randomly initialized element-wise linear layer (i.e. a $1 \times 1$ convolution) is applied to these feature maps and the results are upsampled to the segmentation resolution using bilinear interpolation to predict the class scores. We then use the binary Dice loss from the Segmentation Models Pytorch library[93]. For the NIH CXR dataset we train each class independently using the binary Dice loss.

For detection tasks we first compute segmentation masks from the detection ground truth using the union of all target bounding boxes per sample and then interpret the task as a segmentation task. Note that for the Object CXR dataset we create bounding box masks only for box and ellipse detection targets but use polygon masks for polygon detection targets.

For training, where we only train the linear layer, we use the Adam[45] optimizer with weight decay $1 \times 10^{-6}$ The learning rate is tuned individually for each model and task based on the best validation Dice[24]. The learning rate is multiplied by 0.5 if the validation Dice does not decrease for three consecutive epochs. On the *COVID Rural Linear* and the *RSNA Lin. Seg. 1%* tasks we use batch size eight and do early stopping (based on validation Dice) after 20 epochs. On all other *Linear* tasks we use batch size 64 and do early stopping after 10 epochs. Finally we report the test Dice of the epoch with the best validation Dice. Note that for the *NIH CXR Linear* task we use the *macro averaged Dice (Avg Dice)* as metric.

*YOLOv3 Finetune* We closely follow the architecture[25] of the original YOLOv3[69] but use the pre-trained ResNet50 as its backbone (replacing the Darknet-53 backbone) while randomly initializing all other layers. The backbone features for the three prediction scales are extracted from the outputs of the conv3 (highest resolution), conv4, and conv5 (lowest resolution) blocks of ResNet50, respectively. We use the default anchors presented in their paper but scale them according to our image input size of $224 \times 224$.

For training, we use the losses and loss weights from the YOLOv3 paper and train with the Adam[45] optimizer with weight decay $1 \times 10^{-6}$ The learning rate is tuned individually for each model and task based on the best validation *mean Average Precision (mAP)*. We compute[26] the mAP score following the COCO[49] mAP and with the following Intersection over Union (IoU) thresholds: $0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75$. The learning rate is multiplied by 0.5 if the validation mAP does not decrease for three consecutive epochs. We use a warmup period of 100 iterations in which we do not train the ResNet50 backbone but only the other layers with a learning rate of $1 \times 10^{-3}$ after which we train the

---

[25] Our implementation is based on `https://github.com/BobLiu20/YOLOv3_PyTorch`

[26] `https://github.com/bes-dev/mean_average_precision` (MIT License)

whole model (including the ResNet50). On the *RSNA YOLOv3 Finetune 1%* task we use batch size eight and do early stopping (based on validation mAP) after 20 epochs. On all other *YOLOv3 Finetune* tasks we use batch size 64 and do early stopping after 10 epochs. Finally we report the test mAP of the epoch with the best validation mAP.

*YOLOv3 Frozen* We use the same architecture and loss functions as in the *YOLOv3 Finetune* protocol but freeze the pre-trained ResNet50 weights and never train them. Instead we only train the other layers using the same hyper-parameters as in the *YOLOv3 Finetune* protocol (except for the warmup period which is not relevant in this setting).