

# Weakly Supervised Object Detection in Chest X-Rays with Differentiable ROI Proposal Networks and Soft ROI Pooling

Philip Müller, Felix Meissen, Georgios Kaissis and Daniel Rueckert, *Fellow, IEEE*

**Abstract**—Weakly supervised object detection (WSup-OD) increases the usefulness and interpretability of image classification algorithms without requiring additional supervision. The successes of multiple instance learning in this task for natural images, however, do not translate well to medical images due to the very different characteristics of their objects (i.e. pathologies). In this work, we propose Weakly Supervised ROI Proposal Networks (WSRPN), a new method for generating bounding box proposals on the fly using a specialized region of interest-attention (ROI-attention) module. WSRPN integrates well with classic backbone-head classification algorithms and is end-to-end trainable with only image-label supervision. We experimentally demonstrate that our new method outperforms existing methods in the challenging task of disease localization in chest X-ray images. Code: <https://github.com/philip-mueller/wsrpn>

**Index Terms**—Chest X-ray, Object detection, Pathology detection, Weak supervision

## I. INTRODUCTION

OBJECT localization is a vital task in computer vision. It is not only useful for many of the downstream tasks but is also a crucial factor for the interpretability of machine learning models. However, especially in medical images, localization labels such as bounding boxes are costly and difficult to obtain as they require vast amounts of working hours from trained professionals. Image labels, on the other hand, are easier to collect and can be mined from radiology reports associated with most existing medical images [1], [2]. This makes weakly supervised object detection (WSup-OD) a promising approach for the localization of diseases in medical images. It only requires image-level labels for training, allowing the use of such automatic collection approaches and thus making localization tractable for a wider range of medical applications. WSup-OD has a long history in natural images [3]–[5]. The SOTA methods here use multiple instance learning (MIL) [6],

Philip Müller (philip.j.mueller@tum.de), Felix Meissen (felix.meissen@tum.de), Georgios Kaissis (g.kaissis@tum.de), and Daniel Rueckert (daniel.rueckert@tum.de) are with the Chair for AI in Medicine and Healthcare (I31), School of Computation, Information and Technology, TU Munich, 85748 Garching

Georgios Kaissis is with the group for Reliable AI, Institute for Machine Learning in Biomedical Imaging, Helmholtz Munich, Germany

Daniel Rueckert is with the Biomedical Image Analysis Group, Imperial College London, London SW7 2AZ, UK

Philip Müller and Felix Meissen contributed equally.

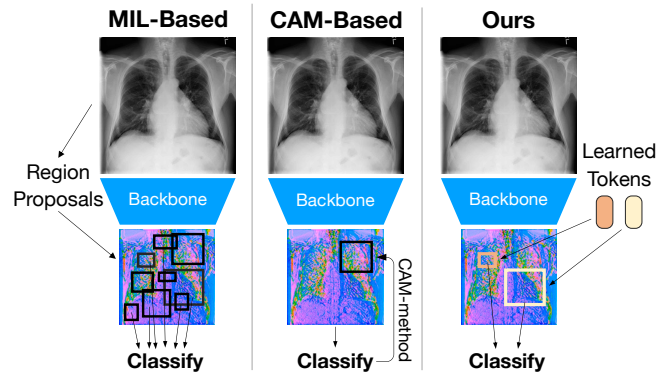


Fig. 1: Schematic illustration of MIL-based, CAM-based, and our novel WSRPN approach.

where bounding box proposals for each image are selected using algorithms such as Selective Search (SS) [7] or Edge Boxes (EB) [8]. These algorithms, however, generate box proposals based on heuristics for objects in natural images and are not suited for detecting diseases in chest X-ray images, as the latter ones have very different characteristics and are more subtle. Selective Search produces box proposals by over-segmenting the image based on pixel intensities. Since pathologies in chest X-rays are not characterized by unique local intensities, the Selective Search algorithm is likely to not focus on them. The Edge Boxes algorithm is based on the observation that in natural images, edges tend to correspond to object boundaries and, thus, searches for regions that wholly enclose edge contours. This method again delivers unsatisfactory results for chest X-rays, as diseases here oftentimes do not have clear edges, and even existing boundaries are often not visible in summation images because they are covered by dense, radiopaque masses along the viewing direction. That is why WSup-OD literature in medical images so far has mostly used CAM-based approaches [2], [9]–[12] that extract boxes from heatmaps. However, these approaches are known to exhibit sub-par performance [13].

To address this issue, we propose Weakly Supervised ROI Proposal Networks (WSRPN), a novel paradigm for WSup-OD in medical images. The bounding box proposals of our method are learned end-to-end and are predicted on the fly during the forward pass through an attention mechanism

similar to DETR [14]. To summarize, our contributions are the following:

- We propose WSRPN, a novel, learnable, end-to-end trainable, and fully differentiable box-proposal algorithm for weakly supervised object detection in medical images.
- We set a new state-of-the-art for weakly supervised object detection on the challenging and commonly used CXR8 [2] dataset.
- To the best of our knowledge, we provide the first multiple-instance learning method successfully trained on this dataset.

## II. RELATED WORK

### A. Weakly supervised object detection on natural images

So far, most works in WSup-OD have focused on natural images in datasets such as PASCAL VOC [15], [16], COCO [17], ILSVRC [18], and CUB-200-2011 [19]. The two dominant approaches in the field are Multiple Instance Learning (MIL) and generating bounding boxes from Class Activation Maps (CAM). Fig. 1 illustrates how these two approaches compare to our proposed method.

*a) MIL:* In MIL-based approaches, each image is considered a bag of instances (regions). Every bag with a positive class label contains at least one positive region. A MIL model is trained only with image labels by assigning every region the label of the whole bag. After training with a large corpus of diverse images, the model becomes invariant to uncorrelated variations and gives higher scores to the most discriminative regions in an image. To identify likely regions of objects in an image, region-proposal-algorithms, such as Selective Search [7] or Edge Boxes [8], are commonly used [4], [5], [20]–[22]. The seminal work here is by Bilen and Vedaldi [4], who extract a feature vector for each region from a backbone network using a Spatial Pyramid Pooling (SPP) layer [23] and subsequently classify each region with a detection (*is it an object?*) and a classification (*which class?*) branch. This method, however, tends to assign higher scores to the most discriminative regions in an image, which do not necessarily cover the whole extent of an object. Subsequent work has, thus, mainly focused on solving the most discriminative region problem by refining the predictions iteratively using multiple refinement streams [20], incorporating the scores of larger context around the region [5], clustering spatially adjacent regions of the same class [21], or maximizing the loss for the most discriminative region to force the model to focus on larger regions [24]. Very recently, Liao *et al.* [25] have proposed a novel method that uses Class Activation Maps as pseudo-ground-truth and cross-attention with learnable tokens to predict bounding boxes. Unlike our proposed WSRPN, however, their method is not fully differentiable, therefore limiting its use in more complex end-to-end models (c.f. Sec. V).

*b) CAM:* The idea of using Class Activation Mapping for weakly-supervised object detection was first proposed by Zhou *et al.* [26]. This method leverages the weights of the final classification layer to classify each patch in the unpooled feature map and create an activation heatmap for each class that can be thresholded and used for object detection. A

similar idea was proposed by Pinheiro *et al.* [27]. However, they first classified each patch in the feature map and then aggregated the resulting scores via LSE pooling, alleviating the need to create heatmaps via CAM. The authors of WELDON [28] use max-min pooling instead to incorporate negative evidence in the final classification and, thus, create better class contrast between the regions. Similar heatmaps are created via GradCAM [29], which uses the gradients w.r.t. the feature map instead of the classifier weights. Just like for MIL-based models, several approaches have been made to solve the most-discriminative-region problem for CAM-based models. In ACoL, for example, Zhang *et al.* [30] follow an idea similar to ICMWSD [24], masking out the most discriminative regions to make the model focus more on secondary features.

### B. Weakly supervised object detection in medical images

WSup-OD is an underrepresented topic in the medical literature and is mainly focused on established CAM-based approaches from natural images. Along with the CXR8 dataset, Wang *et al.* [2] proposed a model for WSup-OD. It uses CAMs for detection and LSE pooling [27] instead of average- or max-pooling. The authors of CheXNet [9] also relied on the simple CAM approach for object localization in the chest X-ray images of the CXR8 dataset. To guide the initially unstable localization in early epochs, Hwang and Kim [10] start with training for classification and gradually shift the focus towards detection using a dedicated branch for each of the two tasks. Their detection branch outputs a heatmap as in [27] to localize tuberculosis in chest X-ray images. In [11], the authors extended the work of Pinheiro *et al.* [27] by using a multi-channel map for each class and employing max-min pooling as in [28] to better localize diseases in CXR8. Yu *et al.* [31] included anatomical information from radiology reports to guide localization. Lastly, Tang *et al.* [32] improve the results of [2] on CXR8 by employing a curriculum learning strategy based on Disease Severity Labels mined from radiology reports and using attention guidance to improve localization performance.

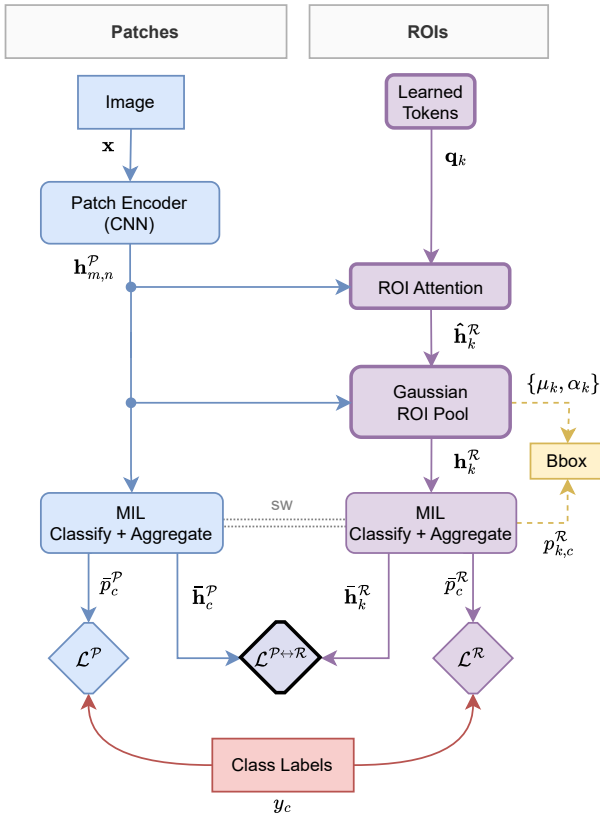
However, none of the above works in the medical domain provides quantitative results of standard metrics in object detection, such as mean Average Precision, limiting the comparability and quantification of their localization performance.

## III. METHOD

### A. Overview

In our weakly supervised object detection setting, we assume that we are given an image that is labeled with a set  $\mathcal{C}$  of non-exclusive classes, i.e. there is one binary classification label  $y_c \in \{0, 1\}$  per class  $c \in \mathcal{C}$  resulting in a multilabel binary classification task. Given only these per-image labels but without any bounding box supervision, we then learn an object detection model.

Fig. 2 provides an overview of our method WSRPN. It is based upon the MIL framework [6], [27], where *regions-of-interest (ROIs)*, i.e. bounding boxes, are predicted using a bounding box proposal algorithm. Following the findings of



**Fig. 2:** Overview of our model architecture. We show the patch branch (blue) and the ROI branch (purple), each with the encoding steps, MIL classification and aggregation, and the loss functions. Components typically used in a MIL model are colored in blue. Our key contributions are outlined with bold lines. “sw” stands for shared weights. Yellow denotes parts of the bounding box prediction.

Sec. I, we however cannot use one of the classical, heuristic bounding box proposal algorithms but instead learn the algorithm end-to-end as a fully differentiable component of our network. We, therefore, follow DETR [14] and use learned ROI query tokens attending to patch features (computed by a CNN backbone) and a box prediction network applied to the resulting ROI features. However, since we do not have supervision for the box proposals, the DETR loss function cannot be applied. To ensure that the predicted box parameters are meaningful (i.e. focus on relevant regions), we apply a Gaussian-based soft approximation of ROI pooling to aggregate ROI features from the patch features. Using a Gaussian distribution during soft ROI pooling introduces an inductive bias that assures that ROI features represent locally restricted regions around the predicted center coordinates of the ROI. The resulting ROI features are then classified and aggregated following the MIL framework, such that they can be trained using per-image class labels. Having only weak supervision, training the ROI proposals directly can lead to instabilities where the bad quality of box proposals during early training stages makes refining these proposals hard. We thus propose a two-branch approach where in the first branch the MIL

framework is applied to patches (we denote the patch branch by  $\mathcal{P}$ ), while the second branch (denoted by  $\mathcal{R}$ ) is designated to ROIs as described. We train both branches using a loss per branch and also introduce a consistency loss, assuring that the ROI proposals are aligned with discriminative patches.

In Sec. III-B and Sec. III-C, we describe the details of the patch and ROI branch, respectively, and in Sec. III-D, we describe how these branches can be trained using weak supervision from classification labels.

## B. Patch branch

a) *Patch encoder:* In the patch branch, we first encode each image into  $H \times W$  patches using the CNN backbone (we use DenseNet121 [33]). These patches are then projected to the model dimension  $d$ , and 2D cosine position encodings [34], [35] are added. We denote the resulting embeddings of patch  $(m, n)$  as  $\mathbf{h}_{m,n}^P \in \mathbb{R}^d$ , where  $m \in \{1, \dots, H\}$  is the  $y$ -index and  $n \in \{1, \dots, W\}$  is the  $x$ -index of the patch.

b) *Patch classification:* We now follow the MIL [6], [27] approach and classify each patch  $(m, n)$  into the classes in  $\mathcal{C}$ , but also predict an additional no-finding (i.e. background) class, denoted as  $\emptyset$ . We compute the class logits  $\tilde{p}_{m,n,c}^P$  of all classes in  $\mathcal{C}$  and the no-finding class  $\emptyset$  by applying a multi-layer perceptron (MLP) to the corresponding patch features  $\mathbf{h}_{m,n}^P$  and then compute the class probabilities  $p_{m,n,c}^P$  via

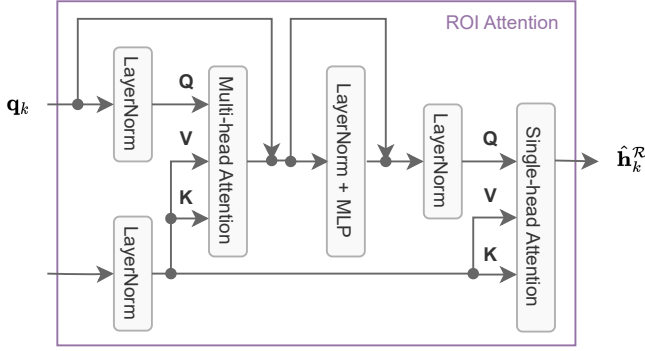
$$\begin{aligned} p_{m,n,\emptyset}^P &= \phi(\tilde{p}_{m,n,\emptyset}^P), \\ p_{m,n,c}^P &= (1 - p_{m,n,\emptyset}^P) \cdot \phi(\tilde{p}_{m,n,c}^P) \quad \forall c \in \mathcal{C}, \end{aligned} \quad (1)$$

where  $\phi$  is the sigmoid function. Patches with large no-finding probabilities  $p_{m,n,\emptyset}^P$  receive lower probabilities for other classes  $c \neq \emptyset$ . Note that the other classes do not influence each other (i.e. each class is considered as a binary classification task) and are thus non-exclusive. We found this approach more effective than having exclusive classes using Softmax.

c) *Aggregation of patch probabilities:* Further following the MIL framework, we now obtain a single per-image probability for each class  $c$  by aggregating the probabilities of all the patches using the *LogSumExp (LSE)* function [27] as a smooth approximation of max pooling as in [2], [27], where we set the scaling hyperparameter  $r$  to 5.0. We again assume multilabel binary classes, i.e. different classes  $c$  are treated independently of each other instead of being exclusive. The aggregated probabilities  $\bar{p}_c^P$  of classes  $c \in \mathcal{C}$  are thus computed as

$$\bar{p}_c^P = \text{LSE}_{m,n}(p_{m,n,c}^P) \quad \forall c \in \mathcal{C}. \quad (2)$$

The no-finding class  $\emptyset$  is considered a special case, and we aggregate it in two different ways: (i) following the OR logic, denoted by  $\vee \emptyset$ , where the class is considered positive if there is any positive patch (similar to the other classes), and (ii) following the AND logic, denoted by  $\wedge \emptyset$ , where it is considered positive only if all patches are positive, i.e. where there is no finding in the whole image. Case (ii) is implemented by inverting the probabilities of  $p_{m,n,\emptyset}^P$  before LSE pooling. The OR approach assures that there are always no-finding patches in an image, i.e. not all patches should



**Fig. 3:** ROI attention component from our ROI branch. Using cross-attention, ROI tokens  $\{q_k\}$  gather relevant information from the patch features  $\{h_{m,n}^P\}$  to compute the ROI features  $\{\hat{h}_k^R\}$ .

be assigned a class, while the AND approach assures that in samples without any other classes, there are only no-finding patches.

### C. ROI branch

a) *ROI attention:* In the ROI branch, we use  $K$  learned ROI tokens  $q_k$  (where  $K$  is a hyperparameter, set to 10 in our experiments). Given a ROI token  $q_k$ , we now use our *ROI attention* component to gather relevant information from the patch features  $h_{m,n}^P$  to compute the ROI features  $\hat{h}_k^R$ . As shown in Fig. 3, the ROI attention component first performs multi-head cross attention [35] with ROI tokens used as queries and patch features used as keys and values. It then further processes the resulting token features using an MLP and a single-head cross-attention layer, where patch features are again used for keys and values.

b) *Box prediction and Gaussian ROI pooling:* Given the token features  $\hat{h}_k^R$  of token  $k$ , we now predict its box center coordinates  $\mu_k$  and size  $\sigma_k$ , each relative to the image size. We assume that relevant features within each ROI are roughly distributed following a normal distribution around the box center. Following this assumption, we now propose a smooth, and therefore differentiable, approximation of (hard) ROI pooling [36]. For each ROI  $k$ , we compute a soft receptive field (i.e. attention map)  $A_{k,m,n}$  over all patches  $(m,n)$ , centered over the ROI center  $\mu_k$  and with its scale (i.e. width and height) controlled by  $\sigma_k$ . We compute the receptive field  $A_{k,m,n}$ , which is proportional to the probability density function of a 2D multivariate Gaussian with independent  $x$  and  $y$  components (i.e. with zero covariance), as

$$A_{k,m,n} \propto \exp \left[ -\frac{1}{2} \left( \frac{\frac{m+0.5}{H} - \mu_{k,y}}{\sigma_{k,y}} \right)^2 \right] \times \exp \left[ -\frac{1}{2} \left( \frac{\frac{n+0.5}{W} - \mu_{k,x}}{\sigma_{k,x}} \right)^2 \right]. \quad (3)$$

Examples of such receptive fields  $A_{k,m,n}$  are shown in Fig. 6. Finally, we aggregate the patch features  $h_{m,n}^P$  for each ROI  $k$  using the receptive field  $A_{k,m,n}$  to get the final ROI features  $h_k^R$ .

c) *ROI classification:* We assign each ROI  $k$  a probability  $p_{k,c}^R$  for each  $c \in \mathcal{C} \cup \{\emptyset\}$  by using the classifier from the patch branch, including sharing the same weights, and applying it to the ROI features  $h_k^R$ .

d) *MIL aggregation of ROI probabilities:* As in the patch branch, we again follow the MIL framework to aggregate the ROI probabilities  $p_{k,c}^R$  over the whole image. However, instead of using LSE, we found the noisyOR [37]–[39] aggregation strategy more effective. noisyOR and its counterpart noisyAND are defined as follows:

$$\text{noisyOR}_k(\mathbf{p}) = 1 - \prod_k (1 - p_k), \quad (4)$$

$$\text{noisyAND}_k(\mathbf{p}) = \prod_k p_k. \quad (5)$$

Using these aggregation functions, we now compute the aggregated ROI probabilities  $\bar{p}_c^R$  for  $c \in \mathcal{C}$ :

$$\bar{p}_c^R = \text{noisyOR}_k(p_{k,c}^R) \quad \forall c \in \mathcal{C}. \quad (6)$$

We again consider the special cases for the no-finding class and aggregate with the OR ( $\bar{p}_{\vee\emptyset}^R$ ) and AND ( $\bar{p}_{\wedge\emptyset}^R$ ) logic.

### D. Weakly supervised loss function

Our weakly supervised loss function is defined as

$$\mathcal{L} = \mathcal{L}^P + \mathcal{L}^R + \mathcal{L}^{P \leftrightarrow R}, \quad (7)$$

where  $\mathcal{L}^P$  trains the patch branch,  $\mathcal{L}^R$  trains the ROI branch, and  $\mathcal{L}^{P \leftrightarrow R}$  assures that both branches are mutually consistent. The branch-specific loss functions ( $\mathcal{L}^P$  and  $\mathcal{L}^R$ ) each consist of two components: (i) a multilabel binary cross entropy loss  $\mathcal{L}_{\text{bce}}$  applied on aggregated patch or ROI probabilities for providing strong gradients, and (ii) a supervised contrastive loss  $\mathcal{L}_{\text{supcon}}$  [40] applied on per-class features from the patch or ROI branch for pushing the patches and ROIs to focus on discriminative regions. We therefore define the branch-specific loss functions as

$$\mathcal{L}^P = \mathcal{L}_{\text{bce}}^P + \mathcal{L}_{\text{supcon}}^P, \quad \mathcal{L}^R = \mathcal{L}_{\text{bce}}^R + \mathcal{L}_{\text{supcon}}^R. \quad (8)$$

a) *Multilabel binary cross entropy:* For the multilabel binary cross entropy losses  $\mathcal{L}_{\text{bce}}^P$  and  $\mathcal{L}_{\text{bce}}^R$ , we use the per-image binary labels  $y_c$  with  $c \in \mathcal{C}$  and  $y_c \in \{0, 1\}$ . We additionally define the no-finding label with AND logic  $y_{\wedge\emptyset}$  as true only if no other classes are true, i.e.  $y_{\wedge\emptyset} = 1 - \max_{c \in \mathcal{C}} y_c$ , as then all patches/ROIs should be classified as no-finding, and the no-finding label with OR logic  $y_{\vee\emptyset}$  as always true, i.e.  $y_{\vee\emptyset} = 1$ , as there should always be some patch/ROI that contains no finding. The losses  $\mathcal{L}_{\text{bce}}^P$  and  $\mathcal{L}_{\text{bce}}^R$  are weighted multilabel binary cross entropy losses over the classes  $\mathcal{C} \cup \{\wedge\emptyset, \vee\emptyset\}$  and are applied to the aggregated patch probabilities  $\bar{p}_c^P$  and ROI probabilities  $\bar{p}_c^R$ , respectively.

b) *Supervised contrastive loss:* The losses  $\mathcal{L}_{\text{supcon}}^P$  and  $\mathcal{L}_{\text{supcon}}^R$  are based on the supervised contrastive loss [40], which is an NTXent-based loss function [41] where positive pairs are defined based on label supervision. We consider each class  $c \in \mathcal{C}$  (but not the no-finding class  $\emptyset$ ) independently (as a



binary label) and define the set of positive samples  $j$  for each sample  $i$  and class  $c$  as

$$P(i, c) = \left\{ j \in \{1, \dots, N\} : y_c^{(i)} = y_c^{(j)} \right\}. \quad (9)$$

Following this setting, we require (sample-wise) per-class features for each  $c \in \mathcal{C}$ , which are computed once from the patch branch (for  $\mathcal{L}_{\text{supcon}}^{\mathcal{P}}$ ) and once from the ROI branch (for  $\mathcal{L}_{\text{supcon}}^{\mathcal{R}}$ ), and are denoted by  $\bar{h}_c^{\mathcal{P}} \in \mathbb{R}^d$  and  $\bar{h}_c^{\mathcal{R}} \in \mathbb{R}^d$ , respectively. We consider the class probabilities of each patch ( $p_{m,n,c}^{\mathcal{P}}$ ) or ROI ( $p_{k,c}^{\mathcal{R}}$ ) and compute the per-class features  $\bar{h}_c^{\mathcal{P}}$  and  $\bar{h}_c^{\mathcal{R}}$  as a weighted sum of all patch ( $h_{m,n}^{\mathcal{P}}$ ) and ROI ( $h_k^{\mathcal{R}}$ ) features, respectively, with weights computed as their (normalized) class probabilities. Finally, we project the results using an MLP. Similarly, we compute  $\bar{h}_c^{\mathcal{R}}$  from ROI features  $h_k^{\mathcal{R}}$  considering  $p_{k,c}^{\mathcal{R}}$ , where the MLP is shared between both branches.

Given these aggregated patch features  $\bar{h}_c^{\mathcal{P}}$  and ROI features  $\bar{h}_c^{\mathcal{R}}$ , respectively, as representations for class  $c$  in sample  $i$ , the losses  $\mathcal{L}_{\text{supcon}}^{\mathcal{P}}$  and  $\mathcal{L}_{\text{supcon}}^{\mathcal{R}}$  follow the following form:

$$\mathcal{L}_{\text{supcon}} = \frac{1}{N|\mathcal{C}|} \sum_{i=1}^N \sum_{c \in \mathcal{C}} \frac{1}{|P(i, c)|} \sum_{j \in P(i, c)} \log \frac{e^{\cos(\bar{h}_c^{(i)}, \bar{h}_c^{(j)})/\tau}}{\sum_{j'=1}^N e^{\cos(\bar{h}_c^{(i)}, \bar{h}_c^{(j')})/\tau}}. \quad (10)$$

*c) Patch-ROI consistency regularizer:* To stabilize training and guide the generation of useful features in the ROI branch, we introduce a consistency regularization loss  $\mathcal{L}^{\mathcal{P} \leftrightarrow \mathcal{R}}$ . This loss ensures the agreement between the spatial distribution of class features of the ROI- and the patch branch. To calculate this agreement, we first need to compute spatial class-distribution probabilities (i.e. patch-wise class probabilities) for both the patch- and the ROI branch. While these class probability maps already exist for the patch branch (cf.  $p_{m,n,c}^{\mathcal{P}}$  from Eq. (1)), getting them for the ROI branch requires further steps:

For the ROI branch, we know the class probabilities  $p_{k,c}^{\mathcal{R}}$  and the spatial distribution (given by the soft receptive field  $A_{k,m,n}$ ) of each ROI. We use these to compute the spatial class map  $p_{m,n,c}^{\mathcal{R} \rightarrow \mathcal{P}}$  of each class  $c$  for each patch ( $m, n$ ) as follows:

$$p_{m,n,c}^{\mathcal{R} \rightarrow \mathcal{P}} = \text{noisyOR}_k(A_{k,m,n} \cdot p_{k,c}^{\mathcal{R}}) \quad \forall c \in \mathcal{C}. \quad (11)$$

For the no-finding class  $\emptyset$ , we consider the assigned patches of ROIs with high no-finding probabilities  $p_{k,\emptyset}^{\mathcal{R}}$  as well as patches where ROIs have low attention  $A_{k,m,n}$  and use noisyAND pooling over the ROIs:

$$p_{m,n,\emptyset}^{\mathcal{R} \rightarrow \mathcal{P}} = \text{noisyAND}_k(A_{k,m,n} \cdot p_{k,\emptyset}^{\mathcal{R}} + (1 - A_{k,m,n})) \quad (12)$$

This assures that patches that are only marginally considered during Gaussian ROI pooling but have high probabilities in real classes  $c \in \mathcal{C}$ , receive high probabilities for the no-finding class.

We now define the consistency loss  $\mathcal{L}^{\mathcal{P} \leftrightarrow \mathcal{R}}$  using the empirical KL-divergence  $D_{\text{KL}}$  from the newly computed spatial class

map  $p_{m,n,c}^{\mathcal{R} \rightarrow \mathcal{P}}$  (from the ROI branch) to the (original) spatial class map  $p_{m,n,c}^{\mathcal{P}}$  (from the patch branch):

$$\mathcal{L}^{\mathcal{P} \leftrightarrow \mathcal{R}} = \frac{1}{HW} \sum_{m,n} \sum_{c \in \mathcal{C} \cup \{\emptyset\}} D_{\text{KL}} \left[ p_{m,n,c}^{\mathcal{P}} \parallel p_{m,n,c}^{\mathcal{R} \rightarrow \mathcal{P}} \right] \quad (13)$$

## E. Inference

During inference, we initially predict one box for each ROI  $k$ . Center position  $\mu_k$  and box size  $\sigma_k$ , computed during box prediction, are used as box parameters. We compute the predicted class  $c_k^* \in \mathcal{C}$  of ROI  $k$  as  $c_k^* = \arg \max_{c \in \mathcal{C}} p_{k,c}^{\mathcal{R}}$  and use  $p_{k,c_k^*}^{\mathcal{R}}$  as its confidence score. Finally, we apply standard post-processing as it will be described in Sec. IV.

## IV. EXPERIMENTS

We show the effectiveness of our Weakly Supervised ROI Proposal Network on the task of disease localization in chest X-ray images.

*a) Dataset and evaluation metrics:* We follow previous works [2], [9], [11] and evaluate on the challenging ChestXray-8 (CXr8) dataset [2]. The dataset consists of 108 948 X-ray images from the National Institutes of Health Clinical Center in the US. The dataset contains labels for eight different disease types and “no-finding” ( $\emptyset$ ). Each image can have more than one positive label, turning the task into a multi-class classification problem. All labels were automatically mined from associated radiology reports with an algorithm that achieved an F1 score of 0.90 on an external dataset. The labels, thus, include a significant amount of noise, making the dataset challenging, even for classification. Additionally, the dataset contains 984 bounding boxes on 882 images from unique patients, hand-labeled by a board-certified radiologist. From the images with bounding boxes, we used 50% for validation and kept the other 50% as a held-out test set. The images of patients that were not included in the validation or test sets were used for training.

To compare the performance of our proposed model with the baselines, we report the Robust Detection Outcome (RoDeO) [42], a recently proposed metric for object detection in medical images, such as Chest X-rays, that reflects the clinical requirements for object detection methods better than other metrics and further gives insights about strengths and weaknesses of the models. Additionally, we report standard metrics, such as Average Precision (AP) and localization accuracy (loc-acc) at two different IoU thresholds (0.3 and 0.5). Note, however, that loc-acc is biased to favor models that predict fewer boxes.

*b) Implementation details:* As the backbone for our model and all baselines, we used a Densenet121 [33] as in [9], [11], pre-trained on ImageNet [43]. For all CAM-based methods that produce heatmaps, we adopted the bounding box generation method and parameters of Wang *et al.* [2], where the heatmaps are binarized, and box proposals are drawn around each connected component. We extended this method to also produce class probabilities and confidence scores per box (c.f. supplementary material). This is necessary to apply score-based postprocessing and compute the Average Precision metrics. Unless indicated otherwise, we post-process the

**TABLE I:** Mean and standard deviation of our method WSRPN against baseline methods on RoDeO, AP, and localization accuracy. The best method per metric is marked in bold. Our method outperforms all baselines on all object detection metrics, setting a new state-of-the-art for weakly supervised object detection on the challenging CXR8 [2] dataset.

Method	RoDeO [%]				AP [%]		loc-acc	
	cls	loc	shape	total	IoU@0.3	IoU@0.5	IoU@0.3	IoU@0.5
WSRPN (ours)	<b>31.7±2.4</b>	<b>44.1±1.6</b>	<b>29.4±1.0</b>	<b>34.0±1.3</b>	<b>9.44±0.90</b>	<b>6.34±0.86</b>	<b>0.78±0.00</b>	<b>0.77±0.00</b>
CheXNet [9]	19.8±1.1	19.9±0.7	10.9±0.4	15.6±0.5	8.26±0.81	3.32±0.56	0.55±0.00	0.52±0.00
↳ w/ noisyOR aggregation	23.9±1.1	20.1±0.7	12.2±0.4	17.3±0.5	8.45±0.92	1.13±0.33	0.59±0.00	0.55±0.00
↳ w/ $\mathcal{L}_{supcon}^P$	20.8±1.2	22.3±0.8	11.9±0.5	17.0±0.6	7.44±0.83	4.00±0.65	0.58±0.00	0.56±0.00
STL [10]	19.0±1.0	18.5±0.6	10.6±0.4	14.9±0.5	8.59±0.78	2.73±0.58	0.54±0.00	0.50±0.00
GradCAM [29]	17.6±1.2	17.5±0.6	9.8±0.4	13.9±0.5	7.07±0.87	0.18±0.14	0.54±0.00	0.51±0.00
CXR [2]	19.9±1.1	19.5±0.7	11.3±0.4	15.8±0.5	8.54±0.87	1.46±0.41	0.55±0.00	0.51±0.00
WELDON [28]	18.5±1.3	20.6±0.7	12.1±0.4	16.2±0.5	6.76±0.82	0.48±0.27	0.56±0.00	0.52±0.00
MultiMap Model [11]	21.0±1.2	20.0±0.7	11.6±0.4	16.3±0.5	7.53±0.80	1.53±0.41	0.57±0.00	0.53±0.00
LSE Model [27]	20.0±1.1	21.3±0.7	11.5±0.4	16.3±0.5	3.07±0.60	0.58±0.29	0.56±0.00	0.54±0.00
ACoL [30]	14.8±1.0	11.9±0.5	10.2±0.4	12.0±0.4	4.27±0.66	2.84±0.58	0.48±0.00	0.48±0.00

predictions of all models by keeping only the most confident predicted box per class (top1-per-class), a valid assumption in the CXR8 dataset that has at maximum one box for any class per image. We implemented all models in PyTorch [44] and optimized them using AdamW [45] with a learning rate of  $1.5 \cdot 10^{-4}$ , weight decay of  $10^{-6}$ , and gradient clipping at norm 1.0. All models were trained for a maximum number of 50000 iterations with early stopping (patience set to 10000) and a batch size of 128. Finally, the checkpoint with the highest mAP on the validation set was chosen. The images were resized to  $224 \times 224$  pixels and normalized with the mean and standard deviation of the training dataset. During training, we augmented the data by applying random color jitter and random Gaussian blurring, each with a probability of 50%, using the Albumentations library [46]. We applied two different random augmentations to each image to guarantee always at least one positive sample for the  $\mathcal{L}_{supcon}$ . During validation or testing, no data augmentation was applied. All of our experiments were performed on a single Nvidia RTX A6000 GPU. Our model trained for roughly 8 hours, requiring about 11 GB of GPU memory.

#### A. Comparison with the baselines

Table I shows the results of our method and the baselines. WSRPN significantly outperforms all weakly supervised baselines on all metrics by a large margin (Welch’s t-test,  $p < 0.001$ ). Compared to the best baseline CheXNet w/ noisyOR aggregation, WSRPN achieves a relative improvement of 96.5% in RoDeO score, setting a new state-of-the-art. Especially, the box quality of our method is better than the baselines. For the submetrics  $\text{RoDeO}_{loc}$  and  $\text{RoDeO}_{shape}$ , the relative improvements to CheXNet w/ noisyOR aggregation are 119.4% and 141.0%, respectively. Also, in terms of AP and loc-acc, our method outperforms the baselines by a large margin, especially when more accurate localization is required. At the IoU-threshold of 0.5, we notice a relative improvement of 58.5% in AP. On the test set, WSRPN predicts, on average, 1.049 boxes per sample, which much more closely resembles the 1.098 true boxes per sample than CheXNet with 3.411 (even after applying top1-per-class filtering). This quality is

also expressed by the much better loc-acc across all thresholds compared to the baselines.

All baselines in Table I are CAM-based since MIL-based training (WSDDN [4]) did not converge on this challenging dataset and was thus excluded from the table. A likely reason for the failure of WSDDN are the box proposal algorithms available for this method (SS and EB). The box proposals of these algorithms have a significantly lower overlap with the objects in chest X-ray images than with those of natural images (c.f. Table II). This strongly limits the detection performance of models building upon these proposals.

**TABLE II:** We computed the average IoU of the target boxes in CXR8 [2] and PASCAL VOC 2007 [15] and the boxes produced by the Selective Search (SS) [7] and Edge Boxes (EB) [8] algorithms. We only considered the predicted box with the highest IoU for every target box, making these numbers an upper bound for methods using SS or EB.

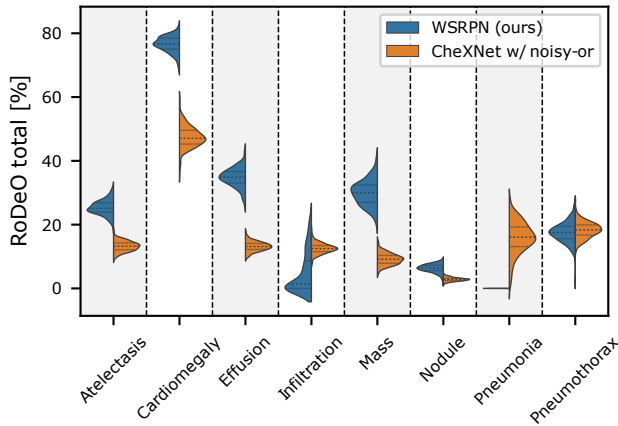
Algorithm	CXR8	PASCAL VOC 2007
Selective Search (fast)	0.31	0.76
Selective Search (quality)	0.37	0.82
Edge Boxes	0.50	0.75

#### B. Performance on different pathologies

In Fig. 4 and Table III, we study the results individually for each of the eight pathologies (atelectasis, cardiomegaly, effusion, infiltration, mass, nodule, pneumonia, pneumothorax) of the bootstrapped ( $N = 250$ ) test dataset. We compare our model WSRPN with the best baseline (CheXNet with noisyOR-aggregation). We observe (cf. Fig. 4) that on five pathologies (atelectasis, cardiomegaly, effusion, mass, and nodule), our method WSRPN performs significantly better than the baseline, often by large margins. On pneumothorax, it is competitive with the baseline, while on two pathologies (infiltration and pneumonia), it performs notably worse. Table III provides further explanations by distinguishing between the quality of classification, localization, and (box) shape similarity. We observe that the localization and shape quality of our model WSRPN outperforms the baseline for

**TABLE III:** Results per pathology of our model WSRPN and the best baselines. Overall, our method WSRPN performs significantly better on five pathologies (atelectasis, cardiomegaly, effusion, mass, and nodule). On pneumothorax, it is competitive with the baselines, while on two pathologies (infiltration and pneumonia), it performs notably worse. However, WSRPN outperforms the baselines on all pathologies when considering localization and shape similarity.

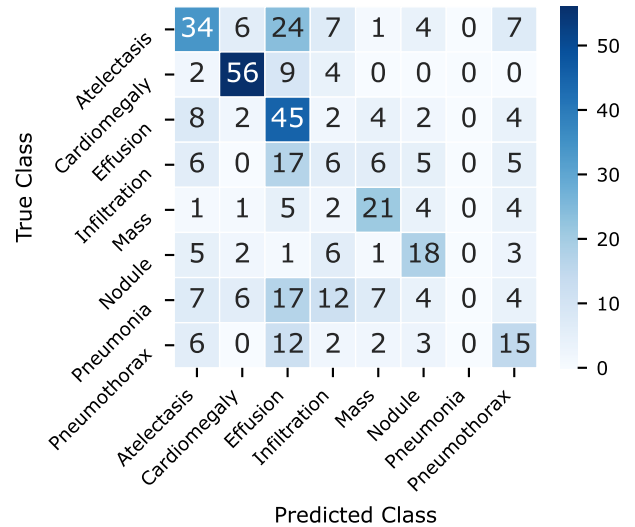
Pathology	RoDeO cls [%]		RoDeO loc [%]		RoDeO shape [%]		RoDeO total [%]	
	WSRPN	CheXNet noisyOR	WSRPN	CheXNet noisyOR	WSRPN	CheXNet noisyOR	WSRPN	CheXNet noisyOR
Atelectasis	<b>28.4</b> ±5.3	24.2±3.0	<b>34.7</b> ±3.5	17.1±2.1	<b>18.7</b> ±1.9	7.8±1.0	<b>25.3</b> ±2.3	13.2±1.6
Cardiomegaly	<b>73.9</b> ±5.7	46.9±5.5	<b>95.0</b> ±1.4	65.0±4.1	<b>66.3</b> ±1.5	38.0±2.4	<b>76.5</b> ±2.5	47.4±3.3
Effusion	<b>53.6</b> ±6.6	26.6±2.7	<b>34.0</b> ±3.6	11.2±1.5	<b>26.4</b> ±2.3	9.9±1.1	<b>34.8</b> ±3.0	13.1±1.5
Infiltration	2.2±3.0	<b>16.7</b> ±2.6	<b>50.4</b> ±5.0	15.2±1.9	<b>30.6</b> ±2.8	8.8±1.1	4.9±6.3	<b>12.5</b> ±1.6
Mass	<b>43.6</b> ±8.7	17.7±3.4	<b>36.3</b> ±5.5	12.1±2.4	<b>20.5</b> ±3.0	5.4±1.2	<b>29.9</b> ±3.9	9.2±1.8
Nodule	<b>39.5</b> ±8.6	19.9±3.9	<b>6.2</b> ±2.6	3.6±1.3	<b>3.7</b> ±0.4	1.4±0.2	<b>6.3</b> ±1.3	2.8±0.6
Pneumonia	0.0±0.0	<b>8.5</b> ±3.6	<b>57.7</b> ±4.3	45.0±4.3	<b>34.6</b> ±2.4	24.7±2.5	0.0±0.0	<b>16.0</b> ±5.0
Pneumothorax	22.6±7.2	<b>37.5</b> ±5.0	<b>14.1</b> ±3.4	13.8±2.4	<b>20.2</b> ±2.6	15.6±2.1	17.6±3.1	<b>18.3</b> ±2.5



**Fig. 4:** Comparison of the results per pathology between our method WSRPN and the best baseline on the bootstrapped ( $N = 250$ ) test set. On five pathologies (atelectasis, cardiomegaly, effusion, mass, and nodule), our WSRPN method performs significantly better, on pneumothorax, it is competitive with the baselines, while on two pathologies (infiltration and pneumonia), it performs worse.

all eight pathologies by large margins. For localization, we observe relative improvements of 103% for atelectasis, 46% for cardiomegaly, 204% for effusion, 232% for infiltration, 200% for mass, 72% for nodule, 28% for pneumonia, 2% for pneumothorax. For shape similarity, these improvements are similarly significant.

We further show the confusion matrix for our proposed WSRPN in Fig. 5. Since confusion matrices are not trivial to generate for object detection problems, we computed them from the 1-to-1 correspondences between predicted and ground-truth boxes after the matching step in RoDeO [42]. The figure confirms that the model often confuses infiltration with pneumonia and that it seems to fail to predict cases of pneumonia. Both classes, however, are not well defined in CXR images. Infiltration is an imprecise descriptive term used for accumulations of an abnormal substance in the lung, while pneumonia is a clinical diagnosis that can not solely be made from an X-ray image. Pneumonia is further often detected by



**Fig. 5:** Confusion matrix for our proposed WSRPN. The matrix was generated from the 1-to-1 correspondences between predicted and ground-truth boxes after the matching step in RoDeO [42].

symptoms such as infiltrations and related to pleural effusions, in part explaining the confusions in the figure.

From these observations, we conclude the following: i) Our model WSRPN performs exceptionally well at localizing pathologies, while its classification capabilities reveal limitations on some classes. This can especially be observed for pneumonia and infiltration, where no or only a few bounding boxes are correctly classified, and on pneumothorax, where the baseline performs particularly well at classification. ii) Good localization and classification capabilities do not necessarily correlate between pathologies. For example, pneumonia is localized well but classified incorrectly, while nodules are classified quite well but are not located well.

### C. Qualitative results and failure cases

Fig. 6 shows example predictions of our model. The first two columns show correctly detected pathologies. The quality of these samples reflects the performance of our model for



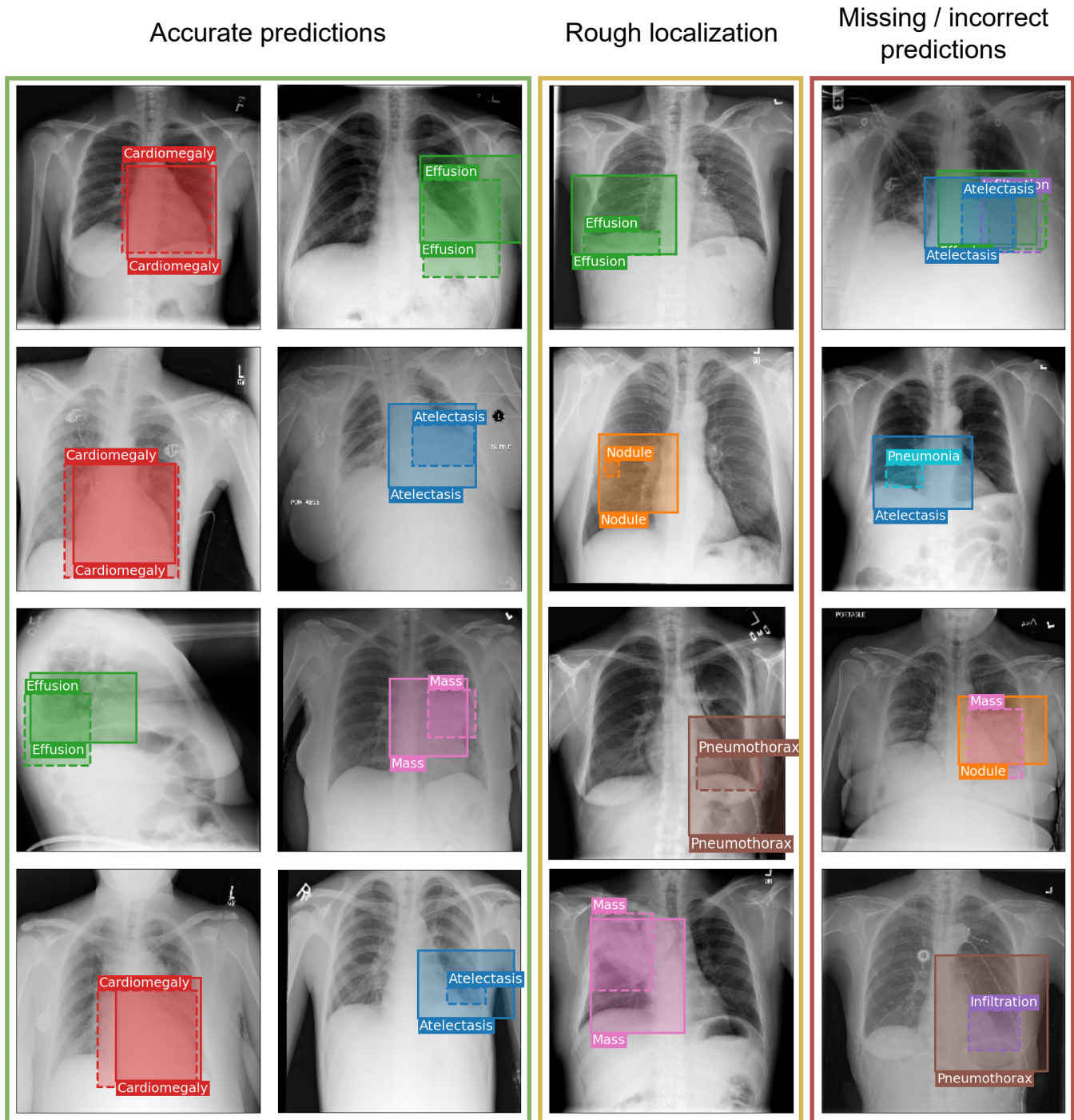


Fig. 6: Qualitative results of some exemplary images. Left: successfully detected pathologies. Middle: Roughly localized correct predictions. Right: failure cases. Solid boxes are predictions. Dashed boxes are human-annotated targets.

each class (c.f. Table III): Cardiomegaly is detected nearly perfectly, but also effusion, atelectasis, and mass are often successfully detected by the proposed WSRPN. Besides the successful cases, we mainly identified two types of failure cases, namely (i) imprecise prediction of the exact extent of the pathology and (ii) miss-classification or partial detection.

Examples of failure type (i) are shown in Fig. 6, column three (yellow column). Here, pathologies are detected and roughly localized. However, the predicted bounding boxes do

not match the target boxes because their aspect ratios differ or the predicted box is too small or too large. Such cases may be hard to tackle and may require semi-supervision, especially as the exact extent of pathologies can be hard to define and is often subjective. Generally, our model tends to produce larger boxes, which is especially problematic for classes with small boxes, such as nodules (c.f. Table III). However, the low performance of the baseline models indicates that this is a common problem of WSUp-OD models.



Failure cases of type (ii) are shown in Fig. 6, column four (red column) and include cases where bounding boxes are predicted approximately correctly but with incorrect classes, especially if classes have similar clinical meaning (e.g. mass and nodule, row 3, col 4) or are correlated (e.g. pneumonia increasing the likelihood of atelectasis, row 2, col 4). In other such cases, classes are not detected at all, especially in samples with multiple overlapping boxes (row 1, col 4). We assume a significant part of this category of failure cases can be tackled by improving the model’s classification performance on the dataset.

#### D. Ablation studies

**TABLE IV:** Ablation study on the loss function. We experimented with different combinations of the individual loss components. Our default configuration (using all components) is highlighted in grey.

Method	$\mathcal{L}_{bce}^P$	$\mathcal{L}_{supcon}^P$	$\mathcal{L}_{bce}^R$	$\mathcal{L}_{supcon}^R$	$\mathcal{L}^{P \leftrightarrow R}$	RoDeO	AP@0.3
WSRPN	✓	✓	✓	✓	✓	<b>34.0±1.3</b>	<b>9.44±0.90</b>
no $\mathcal{L}^P$			✓	✓	✓	12.0±0.5	0.34±0.19
no $\mathcal{L}^R$	✓	✓			✓	18.7±0.7	9.00±0.94
no $\mathcal{L}^{P \leftrightarrow R}$	✓	✓	✓	✓		18.0±0.8	2.74±0.56
no $\mathcal{L}_{bce}$		✓		✓	✓	6.3±0.4	0.16±0.12
no $\mathcal{L}_{supcon}$	✓		✓		✓	23.2±0.9	6.31±0.83
only $\mathcal{L}_{bce}$	✓		✓			14.2±0.6	0.59±0.23
only $\mathcal{L}_{supcon}$		✓		✓		6.2±1.1	2.85±0.66

$K$	5	8	10	12	16
		= $ \mathcal{C} $			= $2 \mathcal{C} $
RoDeO	25.0±1.0	21.7±0.8	<b>34.0±1.3</b>	23.8±0.9	23.6±0.9
AP@0.3	6.88±0.80	5.40±0.78	<b>9.44±0.90</b>	7.01±0.80	5.53±0.69

**TABLE V:** Ablation study on the number of ROI tokens  $K$ .

$\beta$	2.0 (normal)	3.0	4.0	5.0
RoDeO	<b>34.0±1.3</b>	32.2±1.3	31.9±1.3	30.5±1.1
AP@0.3	<b>9.44±0.90</b>	8.55±0.87	8.07±0.91	8.61±0.91

**TABLE VI:** Ablation study on the shape parameter  $\beta$  of the generalized Gaussian distribution.

We conduct extensive ablation studies to quantify the relevance of different loss functions (Table IV), the influence of the number of ROI tokens  $K$  (Table V), the treatment of the no-finding class (Table VII), the assumed distribution of the soft receptive field (Table VI), and the patch size (Table VIII).

a) *Loss functions:* In Table IV, we observe that without the patch branch loss components ( $\mathcal{L}^P = \mathcal{L}_{bce}^P + \mathcal{L}_{supcon}^P$ ), the performance drops substantially in both RoDeO and AP, highlighting the importance of the patch branch for stabilizing the training. If, instead, the ROI branch loss components ( $\mathcal{L}^R = \mathcal{L}_{bce}^R + \mathcal{L}_{supcon}^R$ ) are removed, the performance drops as well, but the model is still competitive with the best baselines. Here, the consistency loss  $\mathcal{L}^{P \leftrightarrow R}$  trains the ROI

**TABLE VII:** Ablation study on the usage of the no finding class  $\emptyset$  during MIL aggregation (and in the BCE losses  $\mathcal{L}_{bce}^P$  and  $\mathcal{L}_{bce}^R$ ). We experimented with ignoring it and using only the classes in  $\mathcal{C}$ , or additionally using either the AND-aggregation ( $\wedge \emptyset$ ) or the OR-aggregation ( $\vee \emptyset$ ) of the no-finding class but found that using both of them (marked in grey) is most effective for both BCE losses.

	$\mathcal{L}_{bce}^P$		$\mathcal{L}_{bce}^R$	
	RoDeO	AP@0.3	RoDeO	AP@0.3
$\mathcal{C}$	18.0±0.7	3.95±0.51	32.3±1.3	9.03±0.99
$\mathcal{C} \cup \{\wedge \emptyset\}$	23.2±0.9	3.05±0.69	33.3±1.2	8.75±0.92
$\mathcal{C} \cup \{\vee \emptyset\}$	22.8±0.8	5.05±0.56	32.3±1.3	8.42±0.89
$\mathcal{C} \cup \{\wedge \emptyset, \vee \emptyset\}$	<b>34.0±1.3</b>	<b>9.44±0.90</b>	<b>34.0±1.3</b>	<b>9.44±0.90</b>

**TABLE VIII:** Ablation study on different patch sizes. We experimented with two options to reduce the patch size from  $32 \times 32$  to  $16 \times 16$ : (a) using lower-level features from *denseblock3* instead of *denseblock4*, and (b) skipping the last pooling features. For both cases, we observe no significant differences in performance.

	RoDeO	AP@0.3
default ( $32 \times 32$ )	<b>34.0±1.3</b>	<b>9.44±0.90</b>
denseblock3 ( $16 \times 16$ )	32.7±1.3	9.12±0.82
skip last pooling layer ( $16 \times 16$ )	31.9±1.3	8.95±0.86

branch based on the predicted patch classes. Training without the consistency loss  $\mathcal{L}^{P \leftrightarrow R}$  leads to poor performance, again confirming the relevance of the consistency loss for stabilizing the box predictions based on the patch branch.

Additionally, we study the relevance of the BCE ( $\mathcal{L}_{bce}^P, \mathcal{L}_{bce}^R$ ) and supervised contrastive ( $\mathcal{L}_{supcon}^P, \mathcal{L}_{supcon}^R$ ) loss functions. Removing the BCE losses always leads to a performance collapse, independently of using the consistency loss  $\mathcal{L}^{P \leftrightarrow R}$ . Removing the supervised contrastive losses leads to a performance drop, but the performance does not collapse entirely if the consistency loss is used.

b) *Number of ROI tokens:* The number of ROI tokens  $K$  determines the maximum number of proposed boxes and is a crucial parameter of our proposed method. Table V shows the detection performance for varying values for this parameter. The optimum is found at  $K = 10$  ( $= |\mathcal{C}| + 2$ ) tokens. A notably anomaly is at  $K = 8$  ( $= |\mathcal{C}|$ ) tokens. We hypothesize that in this case, the single tokens can get too class-specific, which hurts performance.

c) *No-finding handling:* In Table VII, we study the use of the no-finding class  $\emptyset$  during MIL-aggregation (and in the BCE losses  $\mathcal{L}_{bce}^P$  and  $\mathcal{L}_{bce}^R$ ). While in our standard setting, we use both OR ( $\vee \emptyset$ ) and AND ( $\wedge \emptyset$ ) interpretations of this class, we also experiment with using only one of them and ignoring the no-finding class completely (considering only the classes in  $\mathcal{C}$ ). While this hyperparameter has minimal influence on the ROI branch, and the other settings are still competitive with the baselines, the performance degrades when changing this hyperparameter for the patch branch.

d) *Receptive field distribution*: We assume that the relevant features for a pathology roughly follow a Gaussian distribution. We check the validity of this assumption by gradually switching to a more “box-like” distribution by increasing the parameter  $\beta$  of the generalized Gaussian distribution. Table VI shows that performance is maximal at  $\beta = 2$  which corresponds to a standard Gaussian distribution.

e) *Patch size*: Table VIII shows results for changing the patch size of the encoder from  $32 \times 32$  to  $16 \times 16$ . We employed two distinct strategies to investigate the influence of the patch size: First, we used the features of *denseblock3* instead of *denseblock4* in the DenseNet121 encoder. This feature map contains twice as many patches with half the size. Since the features of *denseblock3* may differ significantly from those of *denseblock4*, we experimented with another technique of skipping the final average pooling layer. This results in the same patch size, but the features are closer to those of our proposed models. Neither alternative improved over our default model.

f) *Single- vs multi-box*: Lastly, we measure how our model performs on images with a single (single-box) and multiple target classes (multi-box). In the former, WSRPN achieves a high RoDeO score of  $36.4 \pm 1.6$ , while for the latter, more difficult cases, the score drops to  $23.7 \pm 1.5$ .

## V. DISCUSSION AND CONCLUSION

### A. Clinical applicability

Our model WSRPN shows promising results for pathology localization on chest X-rays. It can provide precise or rough localization for most of the studied pathologies, even if bounding boxes are sometimes too huge. In clinical practice, even such rough localizations can provide massive value as they can assist clinicians in quickly spotting pathologies, especially in time-critical situations like emergency units. However, we also found some limitations that restrict its current clinical applicability. Most importantly, it often misclassifies some of the pathologies. Note that the risks of misclassification differ between pathologies. For example, misclassifying a mass as a nodule does not have severe consequences as one of them is detected since both are indicators of cancer and require further examination. Misclassifying (or missing) pneumothorax, on the other hand, is more critical as immediate clinical intervention may be required. Therefore, future work may focus on improving the classification capabilities of our WSRPN model.

### B. A novel approach towards weakly supervised pathology detection

We propose the first WSup-OD method that can directly optimize (i.e., is differentiable w.r.t.) the box parameters (position and size). Existing WSup-OD methods rely on unsupervised, non-differentiable region proposals (MIL-based methods) or predict bounding boxes using thresholding (CAM-based methods). On the other hand, our Gaussian ROI pooling enables the box parameters to be optimized directly by different kinds of supervision signals, even simultaneously, which is impossible with current other approaches. This enables a wide range

of applications beyond WSup-OD, including, but not limited to, the integration into multimodal large language models, contrastive learning with text, or semi-supervised learning with bounding boxes for a subset of samples. We are convinced that – besides setting a new state-of-the-art on this challenging task – we open up a new research direction without the need for thresholding or external box proposals, which enables this underexplored field (of weakly supervised pathology detection) to progress beyond the existing approaches on which research has mostly stagnated in recent years.

### C. Conclusion

We have proposed WSRPN – a new paradigm for WSup-OD using learned box proposals – after identifying weaknesses in the established box proposal algorithms when applied to X-ray images. While further clinical validation is required, we set a new state-of-the-art in disease detection on the challenging CXR8 [2] dataset and significantly improve upon existing methods. MIL-based methods for natural images have improved dramatically over several years, and we expect a similar evolution for RPN-MIL methods. We deem incorporating other forms of weak supervision like text, anatomy information, or semi-supervision into our framework as promising future research.

## REFERENCES

- [1] J. Irvin *et al.*, “Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison,” in *AAAI*, 2019, pp. 590–597.
- [2] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, “Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases,” in *CVPR*, 2017, pp. 2097–2106.
- [3] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, “Is object localization for free?—weakly-supervised learning with convolutional neural networks,” in *CVPR*, 2015, pp. 685–694.
- [4] H. Bilen and A. Vedaldi, “Weakly supervised deep detection networks,” in *CVPR*, 2016, pp. 2846–2854.
- [5] V. Kantorov, M. Oquab, M. Cho, and I. Laptev, “Contextlocnet: Context-aware deep network models for weakly supervised localization,” in *ECCV*, 2016, pp. 350–365.
- [6] J. Ramon and L. De Raedt, “Multi instance neural networks,” in *ICML 2000 workshop on attribute-value and relational learning*, 2000, pp. 53–60.
- [7] J. R. R. Uijlings, K. E. A. Van De Sande, T. Gevers, and A. W. M. Smeulders, “Selective search for object recognition,” *IJCV*, vol. 104, pp. 154–171, 2013.
- [8] C. L. Zitnick and P. Dollár, “Edge boxes: Locating object proposals from edges,” in *ECCV*, 2014, pp. 391–405.
- [9] P. Rajpurkar *et al.*, “Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning,” *arXiv preprint arXiv:1711.05225*, 2017.
- [10] S. Hwang and H.-E. Kim, “Self-transfer learning for weakly supervised lesion localization,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 239–246.
- [11] C. Yan, J. Yao, R. Li, Z. Xu, and J. Huang, “Weakly supervised deep learning for thoracic disease classification and localization on chest x-rays,” in *ACM BCB*, 2018, pp. 103–110.
- [12] S. Hu *et al.*, “Weakly supervised deep learning for covid-19 infection detection and classification from ct images,” *IEEE Access*, vol. 8, pp. 118 869–118 883, 2020.
- [13] F. Shao *et al.*, “Deep learning for weakly-supervised object detection and localization: A survey,” *Neurocomputing*, 2022.
- [14] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *ECCV*, 2020, p. 213–229.

- [15] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [16] —, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results," <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [17] T.-Y. Lin *et al.*, "Microsoft coco: Common objects in context," in *ECCV*, 2014, pp. 740–755.
- [18] O. Russakovsky *et al.*, "ImageNet Large Scale Visual Recognition Challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.
- [19] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," California Institute of Technology, Tech. Rep. CNS-TR-2011-001, 2011. [Online]. Available: <https://resolver.caltech.edu/CaltechAUTHORS:20111026-120541847>
- [20] P. Tang, X. Wang, X. Bai, and W. Liu, "Multiple instance detection network with online instance classifier refinement," in *CVPR*, 2017, pp. 2843–2851.
- [21] P. Tang *et al.*, "Pcl: Proposal cluster learning for weakly supervised object detection," *IEEE TPAMI*, vol. 42, no. 1, pp. 176–191, 2018.
- [22] Z. Chen, Z. Fu, R. Jiang, Y. Chen, and X.-S. Hua, "Slv: Spatial likelihood voting for weakly supervised object detection," in *CVPR*, 2020, pp. 12 995–13 004.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE TPAMI*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [24] Z. Ren *et al.*, "Instance-aware, context-focused, and memory-efficient weakly supervised object detection," in *CVPR*, 2020, pp. 10 598–10 607.
- [25] M. Liao *et al.*, "End-to-end weakly supervised object detection with sparse proposal evolution," in *ECCV*, 2022.
- [26] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *CVPR*, 2016, pp. 2921–2929.
- [27] P. O. Pinheiro and R. Collobert, "From image-level to pixel-level labeling with convolutional networks," in *CVPR*, 2015, pp. 1713–1721.
- [28] T. Durand, N. Thome, and M. Cord, "Weldon: Weakly supervised learning of deep convolutional neural networks," in *CVPR*, 2016, pp. 4743–4752.
- [29] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *CVPR*, 2017, pp. 618–626.
- [30] X. Zhang, Y. Wei, J. Feng, Y. Yang, and T. S. Huang, "Adversarial complementary learning for weakly supervised object localization," in *CVPR*, 2018, pp. 1325–1334.
- [31] K. Yu, S. Ghosh, Z. Liu, C. Deible, and K. Batmanghelich, "Anatomy-guided weakly-supervised abnormality localization in chest x-rays," in *MICCAI*, 2022, pp. 658–668.
- [32] Y. Tang *et al.*, "Attention-guided curriculum learning for weakly supervised classification and localization of thoracic diseases on chest radiographs," in *MIDL*, 2018.
- [33] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *CVPR*, 2017, pp. 2261–2269.
- [34] Z. Wang and J.-C. Liu, "Translating math formula images to latex sequences using deep neural networks with sequence-level training," *IJDAR*, vol. 24, pp. 1–13, 06 2021.
- [35] A. Vaswani *et al.*, "Attention is all you need," in *NeurIPS*, 2017.
- [36] R. Girshick, "Fast r-cnn," in *ICCV*, 2015, pp. 1440–1448.
- [37] O. Maron and T. Lozano-Pérez, "A framework for multiple-instance learning," in *NIPS*, 1997, pp. 570–576.
- [38] P. Viola, J. C. Platt, and C. Zhang, "Multiple instance boosting for object detection," in *NIPS*, 2005, pp. 1417–1424.
- [39] B. Babenko, P. Dollár, Z. Tu, and S. Belongie, "Simultaneous learning and alignment: Multi-instance and multi-pose learning," in *Workshop on Faces in Real-Life Images: Detection, Alignment, and Recognition*, 2008.
- [40] P. Khosla *et al.*, "Supervised contrastive learning," in *NeurIPS*, vol. 33, 2020, pp. 18 661–18 673.
- [41] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, 2020.
- [42] F. Meissen, P. Müller, G. Kaissis, and D. Rueckert, "Robust detection outcome: A metric for pathology detection in medical images," *arXiv preprint arXiv:2303.01920*, 2023.
- [43] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [44] A. Paszke *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *NeurIPS*, vol. 32, 2019.
- [45] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *ICLR*, 2019.
- [46] A. Buslaev, V. I. Iglovikov, E. Khvedchenya, A. Parinov, M. Druzhinin, and A. A. Kalinin, "Albumentations: Fast and flexible image augmentations," *Information*, vol. 11, no. 2, 2020.