

# Learning from Multiple Expert Annotators for Enhancing Anomaly Detection in Medical Image Analysis

Khiem H. Le<sup>1,†</sup>, Tuan V. Tran<sup>1,†</sup>, Hieu H. Pham<sup>1,2,3,\*</sup>  
Hieu T. Nguyen<sup>1</sup>, Tung T. Le<sup>1</sup>, Ha Q. Nguyen<sup>1,2</sup>

<sup>1</sup>*Smart Health Center, VinBigData JSC, Hanoi, Vietnam*

<sup>2</sup>*College of Engineering and Computer Science, VinUniversity, Hanoi, Vietnam*

<sup>3</sup>*VinUni-Illinois Smart Health Center, Hanoi, Vietnam*

---

## Abstract

Building an accurate computer-aided diagnosis system based on data-driven approaches requires a large amount of high-quality labeled data. In medical imaging analysis, multiple expert annotators often produce subjective estimates about “ground truth labels” during the annotation process, depending on their expertise and experience. As a result, the labeled data may contain a variety of human biases with a high rate of disagreement among annotators, which significantly affect the performance of supervised machine learning algorithms. To tackle this challenge, we propose a simple yet effective approach to combine annotations from multiple radiology experts for training a deep learning-based detector that aims to detect abnormalities on medical scans. The proposed method first estimates the ground truth annotations and confidence scores of training examples. The estimated annotations and their scores are then used to train a deep learning detector with a re-weighted loss function to localize abnormal findings. We conduct an extensive experimental evaluation of the proposed approach on both simulated and real-world medical imaging datasets. The experimental results show that our approach significantly outperforms baseline approaches that do not consider the disagreements among annotators, including methods in which all of the noisy annotations are treated equally as ground truth and the ensemble of different models trained on different label sets provided separately by annotators.

---

\*Corresponding author: [hieu.ph@vinuni.edu.vn](mailto:hieu.ph@vinuni.edu.vn) (Hieu H. Pham)

† These authors share the first authorship of this paper.

*Keywords:* Supervised learning, multiple annotators, object detection.

---

## 1. Introduction

Computer-aided diagnosis (CAD) systems for medical imaging analysis are getting more and more successful thanks to the availability of large-scale labeled datasets and the advances of supervised learning algorithms [1, 2]. To reach expert-level performance, those algorithms usually require high-quality label sets, commonly scarce because of the costly and intensive labeling procedures. A typical label collection process in medical imaging is “*repeated-labeling*”, where multiple clinical experts annotate each data instance to overcome human biases [3, 4, 5]. However, because of the differences from annotator biases and proficiency, annotations from the repeated-labeling process often suffer from high inter-reader variability [6, 7, 8], which could reduce leaning performance if we treat them as ground-truth.

Many prior works have been done to mitigate inter-reader variations in annotations, which can be categorized into two main groups: (i) one-stage approach and (ii) two-stage approach. The first group learns the model, annotators’ proficiency, and latent true labels jointly. Meanwhile, the second group first estimates the true label of each instance from its multiple label sets [9]. This process is known as “*truth inference*”. After that, a supervised learning model is trained on the estimated true labels. All of those approaches show impressive results on both classification and segmentation problems [10, 11].

This work aims at addressing a fundamental question “*How to train a deep learning-based detector effectively from a set of possibly noisy labeled data provided by multiple annotators?*” [12]. To this end, we introduce a novel approach that learns from multiple expert annotators to improve the performance of a deep neural network in detecting abnormalities from chest X-ray images. The proposed approach, as visualized in Figure 1, consists of two stages. The first one is truth inference using Weighted Boxes Fusion (WBF) algorithm [13] to estimate the true labels and their confidence scores. The second stage is to train

an object detector on estimated labels with a re-weighted loss function using implicit annotators' agreement, which is represented by the estimated confidence scores. For evaluation, we first simulate and test the proposed approach on a multiple-experts-detection dataset from MNIST [14] called MED-MNIST. We then validate our approach on a real-world chest X-ray dataset with radiologist's annotations. Experiments on those scenarios demonstrate that the proposed approach provides better detection performance in terms of mAP scores than the baseline of treating multiple annotations as ground truth and the ensemble of models supervised by individual expert annotations.

In summary, our main contributions in this work are two-folds:

- First, we introduce a simple yet effective method that allows a deep learning network to learn from multiple annotators to improve its performance in detecting abnormalities from medical images. The proposed approach aims at estimating the true annotations from multiple experts with confidence scores and uses these annotations to train a deep learning-based detector. This helps remove uncertainty in the learning process and provides higher label quality to train predictive models.
- Second, the proposed approach demonstrates its effectiveness on both simulated and real medical imaging datasets by surpassing current state-of-the-art methods on the context of learning with multiple annotators. In particular, our method is simple and can be applied for a wide range of applications in medical imaging and object detection in general. The codes used in the experiments are available on our Github page at <https://github.com/huyhieupham/learning-from-multiple-annotators>. We also have made the dataset used in this study available for public access on our project's webpage at <https://vindr.ai/datasets/cxr>.

The rest of the paper is organized as follows. Related works on learning from multiple annotators and weighted training techniques are reviewed in Section 2. Section 3 presents the details of the proposed method with a focus on how to

estimate the ground truth annotations from multiple experts. Section 4 provides comprehensive experiments on a simulated object detection dataset and a real-world chest X-ray dataset. Section 5 discusses the experimental results, some key findings, and limitations of this work. Finally, Section 6 concludes the paper.

## 2. Related works

**Learning from multiple annotators.** There are two major lines of research on learning from multiple annotators: two-stage approaches [15, 9, 16] and one-stage approaches [10, 17, 18]. Two-stage approaches infer the true labels first, then train a model using the estimated ones. The most simple solution for label aggregation is majority voting, in which the choice of majority annotators regards as the truth [19]. However, when the skill levels of the annotators differ, the majority voting strategy may not work well. This is a common occurrence in the general “*learning from crowds*” problem when “spammers” are present. Later approaches typically incorporate other information into the truth inference procedure, such as the annotators’ proficiency [20], annotators’ confusion matrix [21, 22], or the difficulty of each sample [23]. While two-stage approaches have the advantage of simplicity in both implementing and debugging, they do not make use of the raw annotations in model learning. One-stage approaches address this issue by simultaneously estimating the hidden true labels and learning the desired model from noisy labels of multiple annotators. Earlier works in this group use Expectation Maximization (EM) algorithm [24] for jointly modeling the annotators’ ability and the latent ground-truth. More recent approaches employ end-to-end frameworks which enable the neural networks to learn directly from the noisy labels [12], and further developed by incorporating annotators’ confusion matrix [11, 10], or instance features [17].

**Weighted training examples.** In this paper, we propose a new re-weighted loss function in which we assign more weights to examples that we consider be more confident. Previous works on the use of weighted training examples can be briefly categorized into two groups: (i) emphasize hard examples and

(ii) emphasize easy examples. Methods in the group (i) include hard-example mining [25, 26], which is a bootstrapping technique over the difficult examples; boosting algorithms [27], where the misclassified examples in preceding weak classifiers are assigned with higher weights; and focal loss [28] that addresses class imbalance problems by adding a regulator to the cross-entropy loss for focusing on hard negative examples. Works in the group (ii) are instances of broader topics such as curriculum learning [29], which is biologically inspired by human gradual learning, with easier examples are preferred in early training stages; learning with noisy labels [30, 31], which prefers examples with smaller training losses as they are more likely to be clean.

Unlike any approaches above, we propose in this paper a new loss function that assigns more weights to more confident examples that determine by the consensus of multiple annotators. Our experimental results validate the correctness of this hypothesis.

### 3. Proposed Method

This section presents details of the proposed method. We first give a formulation on learning from multiple annotators (Section 3.1). We then introduce a simple way to estimate the true labels from multiple annotators (Section 3.2). Next, our network architecture and training methodology with a new re-weighted loss function are described (Section 3.3).

#### 3.1. Problem formulation

Given a set of  $N$  training images  $\{\mathbf{x}_i\}_{i=1}^N$  with corresponding bounding box annotations  $\{\tilde{y}_i^{(r)}\}_{i=1}^N$  from multiple annotators where  $\tilde{y}_i^{(r)}$  denotes the label for example  $\mathbf{x}_i$  given by annotator  $r \in \mathbb{S}(R)$ , which  $\mathbb{S}(R)$  is a set of  $R$  different expert annotators. In this study, we make use of those expert annotations  $\{\tilde{y}_i^{(r)}\}_{i=1}^N$  to estimate a single set of true labels with confidence scores  $\{\mathbf{y}_i; c_i\}_{i=1}^N$ . We then train a supervised object detector with the estimated labels using the proposed re-weighted loss function. In order to evaluate the effectiveness of the proposed

method, we use a gold-standard test set  $\mathcal{T} = \{(\mathbf{x}^{(j)}, \mathbf{y}^{(j)})\}_{j=1}^M$  containing  $M$  examples. In medical imaging scenarios, where the true labels are not available, we obtain the gold-standard test labels  $\mathbf{y}^{(j)}$  from the consensus of a group of experienced radiologists. Figure 1 below shows an overview of the proposed method.

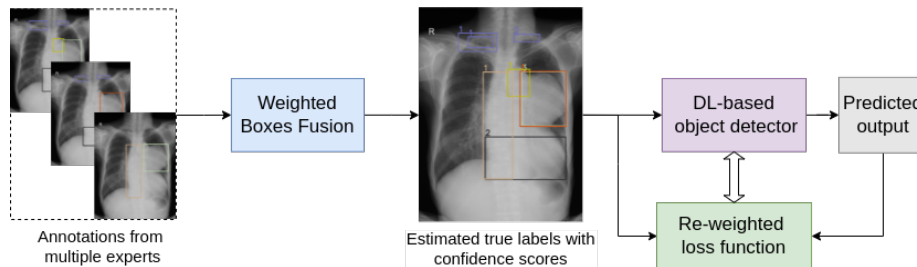


Figure 1: Illustration of the proposed approach that aims to build a deep learning system for abnormal detection on medical scans from multiple expert annotators. The training process contains two stages. The first stage focuses on truth inference, in which it estimates the true labels using the WBF algorithm [13] with the implicit annotator’s agreement as confidence scores. The second uses the estimated confidence scores to train a deep learning-based detector using a re-weighted object detection loss function. To provide abnormality analysis during the testing phase, only the fully trained image detector is required.

### 3.2. Estimating the true labels from multiple expert annotators

We firstly estimate the true labels using Weighted Boxes Fusion (WBF) algorithm [13]. This technique is used for combining predictions from multiple sources, i.e., using ensemble to achieve better prediction results or combining labels of different expert annotators. We describe the WBF algorithm in more detail in Algorithm 1. The final examples used to train deep learning detectors contain merged boxes with confidence scores. The visualization of fused boxes and the corresponding confidence scores are shown in Figure 2. Our fusion box algorithm emphasizes that the greater agreement between bounding boxes (e.g., two or three annotators have the same diagnosis for an abnormal finding on the image), the more likely the box annotation is correct.

---

**Algorithm 1: The WBF algorithm applied for multiple expert annotations**

---

**Input:** An image  $\mathbf{x}$  with a list of annotations  $\tilde{y}$  given by a set  $\mathbb{S}(R)$  of  $R$  experts.

The expert  $r \in \mathbb{S}(R)$  with proficiency  $p_r$  provides the annotations including  $r_x$  boxes,  $A_r = [\text{box}_1, \dots, \text{box}_{r_x}]$ . All of the experts' annotations being merged into a list  $A$ .

**Output:** A list of  $k$  fused boxes  $F = [\text{box}_1, \dots, \text{box}_k]$ .

- 1 Declare empty lists  $L$  and  $F$  for boxes clusters and fused boxes, respectively. Each position in the list  $L$  can have a cluster of boxes or a single box. Each position in  $F$  has only one box, which is the fused box from the corresponding cluster in  $L$ .
- 2 Iterate through all boxes in  $A$  in a cycle and attempt to find a matching box in the list  $F$ . Two boxes are defined matched if they have a high degree of overlap (e.g. IoU > 0.4). If there are more than one matching boxes in  $F$ , the one with the highest IoU will be chosen.
- 3 If the matching box is not found in step 1, add the current box to  $L$  and  $F$  as new entry for the new cluster before moving on to the next box in the list  $A$ .
- 4 If the match is found in step 1, add this box to the list  $L$  at the position  $pos$  which corresponds to the matching box in the list  $F$ .
- 5 Set the fused box's coordinates  $F[pos]$  to be the weighted average of  $T$  boxes accumulated in cluster  $L[pos]$  with the following formulas:

$$x_{1,2} := \frac{\sum_{i=1}^T p_i x_{i,2}}{\sum_{i=1}^T p_i}$$

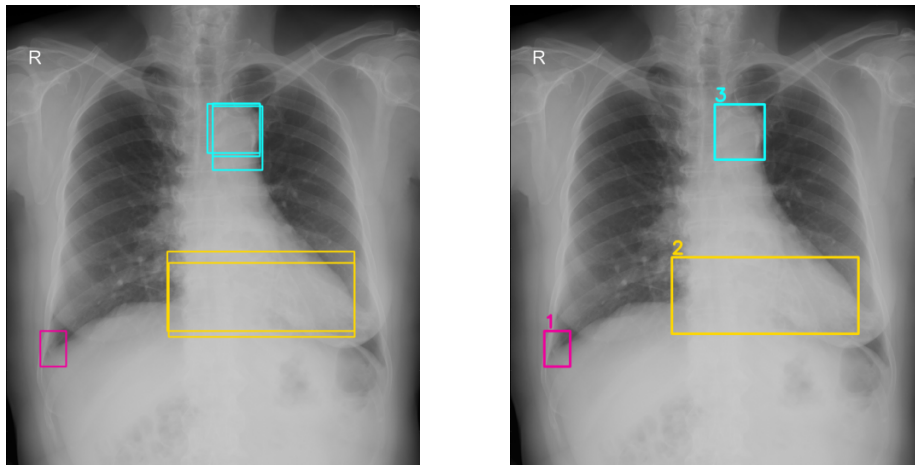
$$y_{1,2} := \frac{\sum_{i=1}^T p_i y_{i,2}}{\sum_{i=1}^T p_i}$$

- 6 Set the the fused boxes' confidence scores in  $F$  to the number of boxes in the corresponding cluster in  $L$  once all boxes in  $A$  have been processed.

$$c := c \min(T, N)$$

The fused boxes with confidence scores now represent the annotators' level of agreement.

---



(a) The original annotations provided by multiple radiology experts. The same abnormal finding is represented by the sample color.

(b) Fused boxes with corresponding confidence scores after applied the WBF algorithm.

Figure 2: (a) Visualization of multiple expert annotations on a chest X-ray example from the VinDr-CXR dataset [5] and (b) the fused boxes with confidence scores obtained by the WBF algorithm.

### 3.3. Network architecture and training methodology

Object detection is a multi-task problem, in which the loss function consists of two parts: (1) the localization loss  $\mathcal{L}_{loc}$  for predicting bounding box offsets and (2) the classification loss  $\mathcal{L}_{cls}$  for predicting conditional class probabilities. In this work, we focus on one-stage anchor-based detectors. A general form of the loss function for those detectors can be written as

$$\mathcal{L}(p, p^*, t, t^*) = \mathcal{L}_{cls}(p, p^*) + \beta I(t) \mathcal{L}_{loc}(t, t^*)$$

$$I(t) = \begin{cases} 1 & \text{if } \text{IoU}\{a, a^*\} > \eta \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where  $t$  and  $t^*$  are the predicted and ground truth box coordinates,  $p$  and  $p^*$  are the class category probabilities, respectively;  $\text{IoU}\{a, a^*\}$  denotes the Intersection



over Union (IoU) between the anchor  $a$  and its ground truth  $a^*$ ;  $\eta$  is an IoU threshold for objectness, i.e. the confidence score of whether there is an object or not;  $\beta$  is a constant for balancing two loss terms  $\mathcal{L}_{cls}$  and  $\mathcal{L}_{loc}$  [32].

We use fused boxes confidence scores  $c_k^i$  obtained from Algorithm 1 to get a re-weighted loss function that emphasizes boxes with high annotators agreement. The new loss function, which we name it as Experts Agreement Re-weighted Loss (EARL) can now be written as

$$\mathcal{L}(p, p^*, t, t^*) = c\mathcal{L}_{cls}(p, p^*) + c\beta I(t)\mathcal{L}_{loc}(t, t^*), \quad (2)$$

where  $c$  is the fused box confidence score.

## 4. Experiments

We validate the proposed method in both synthetic and real-world scenarios: (1) the MED-MNIST, an object detection dataset, which was simulated from MNIST [14] with multiple expert annotations; (2) VinDr-CXR [5], a chest X-ray dataset with labels provided by multiple radiologists. In the following sections, we describe those two datasets and our experiment setup, as well as the experimental results.

### 4.1. Datasets

#### 4.1.1. MED-MNIST Dataset

Based on MNIST [14] – a database of handwritten digits, we synthesize a multiple-experts-detection dataset so-called MED-MNIST by two steps: (1) we construct the detection task by copying and pasting digits from MNIST into a black background with digit sizes are randomly chosen from a predefined range, the bounding box annotations would be the smallest rectangle that contains digits as visualized in Figure 4a; (2) we simulate  $R$  different expert opinions for each sample, assuming those  $R$  experts have the same proficiency  $p$ . The expert annotations are generated by varying two key factors that influence detection annotations: (i) class labels and (ii) object coordinates. To synthesize

the expert annotations on class labels, we use an unique transition matrix  $A_k (k \in \{1, \dots, R\})$  for each expert  $E_k$  to compute probability distributions that represent the expert mis-classification. The proposed transition matrix is shown in Figure 3. About the object coordinates, we simulate bounding box annotations that are highly overlapping with the true bounding box. Both factors (i) and (ii) are controlled by proficiency  $p$ . More specifically,  $A_k$  are diagonally dominant ( $a_{ii} > a_{ij}$  for all  $i \neq j$ ), and  $a_{ii} = \min(\max(0.5, \alpha), 1)$  with  $\alpha \sim \mathcal{N}(p, 0.05)$ . The simulated bounding boxes are subject to have  $IoU$  with the true bounding box being larger than  $p$ . Here we set the number of expert annotations per sample  $R$  to 3, and the proficiency  $p$  to 0.8. The simulated MED-MNIST dataset consists of 5,000 samples for training, 1,000 for hold-out validation and 1,000 for testing.

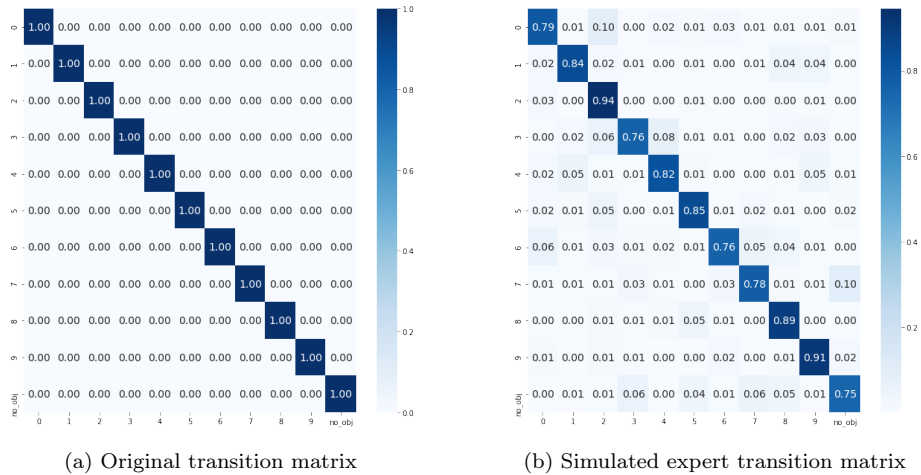


Figure 3: Visualization of the original and synthesized transition matrices. To simulate the false negative scenario, we use an additional class called `no_obj`.

#### 4.1.2. VinDr-CXR Dataset

VinDr-CXR [5], by far the largest public chest X-ray database with radiologist-generated annotations. It consists of 18,000 chest X-ray scans that come with both the localization of critical findings and the classification of common thoracic

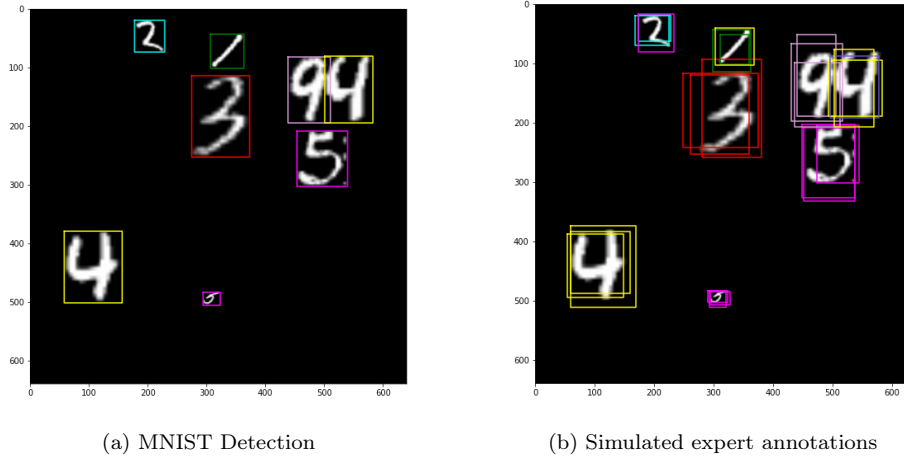


Figure 4: The MED-MNIST dataset with multiple expert annotations, obtained by perturbing boxes and classes from the MNIST dataset [14].

diseases. The dataset includes 15,000 scans for training and 3,000 scans for testing. In particular, the annotations were obtained by a group of 17 radiologists with at least eight years of experience. Each image in the training set was independently labeled by three radiologists, while the annotations in the test set were carefully treated and obtained by the consensus of 5 radiologists. Several examples from the VinDr-CXR dataset are shown in Figure 5.

#### 4.1.3. Rads-VinDr-CXR Dataset

One intriguing characteristic of the VinDr-CXR dataset [5] is that 94.28% of the abnormal scans in the training set (3,315 out of 3,516) were annotated by a group of three radiologists with their correspondence IDs being *R8*, *R9* and *R10*. As a result, we create Rads-VinDr-CXR, a sub-dataset that is only annotated by those three radiologists. The Rads-VinDr-CXR serves as a suitable multiple annotators dataset to validate the proposed approach.

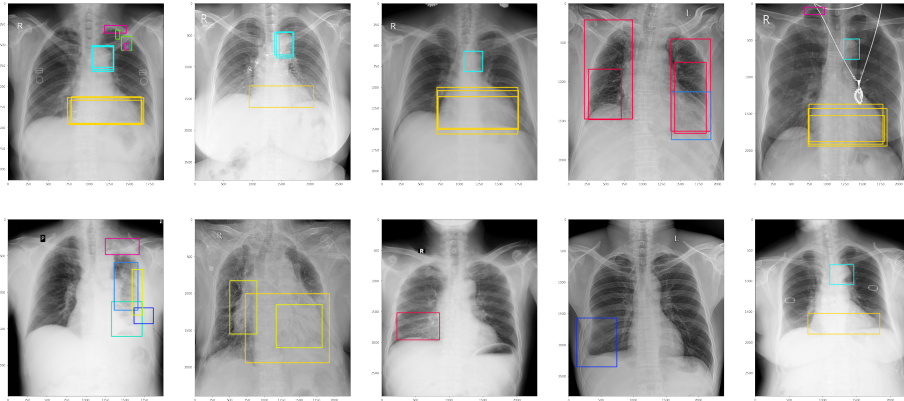


Figure 5: Visualization of abnormal findings (different bounding box colors represent different findings) from the VinDr-CXR dataset: (top) Each scan in the training set was annotated by three different radiologists; (bottom) Test set annotations were obtained from the consensus of five radiologists.

## 4.2. Experimental Setup

### 4.2.1. Evaluation metric

For all experiments, we report the detection performance using the standard mean average precision metric at a threshold of 0.4 (mAP@0.4) [33]. Specifically, a predicted object is a true positive if it has an IoU of at least 0.4 with a ground truth bounding box. The average precision (AP) is the mean of 101 precision values, corresponding to recall values ranging from 0 to 1 with a step size of 0.01. The final metric is the mean of AP over all lesion categories. We also employ mAP@[0.5:0.95:0.05] as an additional metric to assess the model’s performance on different IoU thresholds ranging from 0.5 to 0.95 with a step size of 0.05.

### 4.2.2. Implementation Details

The main detector used in our experiments is YOLOv5-S [34]. The network is built with PyTorch 1.7.1 (<https://pytorch.org/>) and trained on two NVIDIA RTX 2080 Ti GPUs. All training and test images are resized to the dimension of  $640 \times 640$  pixels. The detector is trained for 50 epochs with 1cycle learning rate decay [35] using the SGD optimizer [36]. The initial learning rate is set to  $1e-3$ .

To validate the robustness of the proposed approach across different deep learning detectors, we further train and evaluate EfficientDet [37] with sizes D3 and D4. Specifically, all images are resized to  $640 \times 640$  pixels and the model is trained for 30 epochs with constant learning rate  $3e-4$  using the AdamW optimizer [38].

#### 4.2.3. Comparison with the state-of-the-art

To the best of our knowledge, there is no existing multiple-annotators model for object detection tasks in the literature. Hence, we compare the performance of the proposed method against the baseline, which uses all experts’ annotations per example without taking into account the disagreement among annotators. On the Rads-VinDr-CXR dataset, we further compare our method with the Rads-ensemble, which is the ensemble of independent models trained on separate radiologists’ annotation sets. In this case, the WBF algorithm is used to combine the predictions of those models.

#### 4.3. Experimental Results

Table 1 and Table 3 report the experimental results of the YOLOv5-S detector on MED-MNIST and VinDr-CXR datasets, respectively. On both synthetic and real-world datasets, the proposed approach outperforms the chosen baselines, even with the ensemble of individual experts’ models. Specifically, on the test set of the MED-MNIST dataset, our method reports an overall  $mAP@0.4$  of 0.980 and an overall  $mAP@[0.5:0.95:0.05]$  of 0.849. These results are much higher the performance of the baseline with  $mAP@0.4 = 0.975$  and  $mAP@[0.5:0.95:0.05] = 0.815$ , boosting the  $mAP$  scores of the baseline by 0.51% and 4.2%, respectively. Experimental results on the VinDr-CXR and Rads-VinDr-CXR datasets also validate the effectiveness of the proposed method. We achieve an overall  $mAP@0.4$  of 0.200 on the VinDr-CXR dataset and an overall  $mAP@0.4$  of 0.158 on the Rads-VinDr-CXR dataset. We emphasize that these results outperform both the baseline model, individual model trained on label provided by individual annotator (i.e.  $R8$ ,  $R9$ ,  $R10$ ), as well as the ensemble

model.

The experimental results with EfficientDet detector are provided in Table 3. We found that better detection performances compared to the baseline have been reported. This evidence confirm the robustness of the proposed approach across deep learning detectors.

Table 1: Experimental results on the MED-MNIST dataset. The highest scores are highlighted in **red**.

Method	mAP@0.4	mAP@[0.5:0.95:0.05]
Baseline	0.975	0.815
WBF+EARL (ours)	<b>0.980</b>	<b>0.849</b>

Table 2: Experimental results on the VinDr-CXR and Rads-VinDr-CXR datasets with the YOLOv5-S detector. The highest scores are highlighted in **red**.

Dataset	Method	mAP@0.4
VinDr-CXR	Baseline	0.190
	WBF+EARL (ours)	<b>0.200</b>
Rads-VinDr-CXR	Baseline	0.148
	R8	0.121
	R9	0.132
	R10	0.124
	Rads-ensemble	0.154
	WBF+EARL (ours)	<b>0.158</b>

## 5. Discussions

### 5.1. Key findings and meaning

To the best of our knowledge, the proposed method is the first effort to train an image detector from labels provided by multiple annotators, which is crucial in constructing high-quality CAD systems for medical imaging analysis. In particular, we empirically showed a notable improvement in terms of mAP

Table 3: Experimental results on the VinDr-CXR dataset while EfficientDet is used as the detector. The scores are measured in mAP@[0.5:0.95:0.05], with highest values highlighted in **red**.

	Baseline	WBF+EARL
EfficientDet-D3	0.1142	<b>0.1353</b>
EfficientDet-D4	0.1223	<b>0.1431</b>

scores by estimating the true labels and then integrating the implicit annotators’ agreement into the loss function to emphasize the clean bounding boxes over the noisy ones. The idea is simple but effective, allowing the overall framework can be applied in training any image machine learning-based detectors.

### 5.2. Limitations

Despite the higher predictive performance over the relevant baselines, we acknowledge that the proposed method has some limitations. First, the overall architecture is not end-to-end. It may not fully exploit the benefits of combining truth inference and training the desired image detector. Second, applying the WBF algorithm to annotation sets with a high level of noise may produce low-quality training data. This case is quite impractical in the medical imaging field when the annotators are experienced clinical experts, but it frequently occurs in the general *learning from crowds* problems.

## 6. Conclusion

This paper concentrates on the use of annotations from multiple experts to build a robust deep learning system for abnormality detection on medical images. We propose using Weighted Boxes Fusion (WBF) algorithm to obtain the aggregated annotations with the implicit annotators’ agreement as confidence scores. The estimated annotations are then used to train a deep learning detector with a re-weighted loss function that incorporates the confidence scores to localize abnormal findings. We empirically demonstrate that the proposed approach

outperforms current state-of-the-art baseline approaches in both synthetic and real-world scenarios. To the best of our knowledge, we introduce for the first time an effective method that trains an object detector from multiple annotators. We believe our method is simple and can be applied widely in medical imaging.

## 7. Acknowledgements

This work was supported by Smart Health Center at VinBigData JSC. The authors gratefully acknowledge all anonymous reviewers for their valuable comments and suggestions.

## References

- [1] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, C. I. Sánchez, A survey on deep learning in medical image analysis, *Medical Image Analysis* 42 (2017) 60–88.
- [2] H.-P. Chan, L. M. Hadjiiski, R. K. Samala, Computer-aided diagnosis in the era of deep learning, *Medical Physics* 47 (5) (2020) e218–e227.
- [3] M. I. Razzak, S. Naz, A. Zaib, Deep learning for medical image processing: Overview, challenges and the future, *Classification in BioApps* (2018) 323–350.
- [4] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Illcus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya, et al., CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33, 2019, pp. 590–597.
- [5] H. Q. Nguyen, K. Lam, L. T. Le, H. H. Pham, D. Q. Tran, D. B. Nguyen, D. D. Le, C. M. Pham, H. T. Tong, D. H. Dinh, et al., VinDr-CXR: An open dataset of chest X-rays with radiologist’s annotations, *arXiv preprint arXiv:2012.15029*.



- [6] T. Watadani, F. Sakai, T. Johkoh, S. Noma, M. Akira, K. Fujimoto, A. A. Bankier, K. S. Lee, N. L. Müller, J. W. Song, J. S. Park, D. A. Lynch, D. M. Hansell, M. Remy-Jardin, T. Franquet, Y. Sugiyama, Interobserver variability in the CT assessment of honeycombing in the lungs, *Radiology* 266 (3) (2013) 936–944.
- [7] A. B. Rosenkrantz, R. P. Lim, M. Haghighi, M. B. Somberg, J. S. Babb, S. S. Taneja, Comparison of interreader reproducibility of the prostate imaging reporting and data system and likert scales for evaluation of multiparametric prostate MRI, *AJR Am J Roentgenol* 201 (4) (2013) W612–618.
- [8] E. Lazarus, M. B. Mainiero, B. Schepps, S. L. Koelliker, L. S. Livingston, BI-RADS lexicon for US and mammography: interobserver variability and positive predictive value, *Radiology* 239 (2) (2006) 385–391.
- [9] V. S. Sheng, J. Zhang, Machine learning with crowdsourcing: A brief summary of the past research and future directions, *Proceedings of the AAAI Conference on Artificial Intelligence* 33 (01) (2019) 9837–9843.
- [10] R. Tanno, A. Saeedi, S. Sankaranarayanan, D. C. Alexander, N. Silberman, Learning from noisy labels by regularized estimation of annotator confusion, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11244–11253.
- [11] L. Zhang, R. Tanno, M.-C. Xu, C. Jin, J. Jacob, O. Ciccarrelli, F. Barkhof, D. Alexander, Disentangling human error from ground truth in segmentation of medical images, in: H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, H. Lin (Eds.), *Advances in Neural Information Processing Systems*, Vol. 33, Curran Associates, Inc., 2020, pp. 15750–15762.
- [12] F. Rodrigues, F. Pereira, Deep learning from crowds, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32, 2018.
- [13] R. Solovyev, W. Wang, T. Gabruseva, Weighted boxes fusion: Ensembling

- boxes from different object detection models, *Image and Vision Computing* 107 (2021) 104117.
- [14] Y. LeCun, C. Cortes, C. Burges, MNIST handwritten digit database, ATT Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist> 2.
- [15] Y. Zheng, G. Li, Y. Li, C. Shan, R. Cheng, Truth inference in crowdsourcing: Is the problem solved?, *Proc. VLDB Endow.* 10 (5) (2017) 541–552.
- [16] Y. Jin, M. Carman, Y. Zhu, Y. Xiang, A technical survey on statistical modelling and design methods for crowdsourcing quality control, *Artificial Intelligence* (2020) 103351.
- [17] J. Li, H. Sun, J. Li, Z. Chen, R. Tao, Y. Ge, Learning from multiple annotators by incorporating instance features (2021). [arXiv:2106.15146](https://arxiv.org/abs/2106.15146).
- [18] L. Zhang, R. Tanno, M.-C. Xu, C. Jin, J. Jacob, O. Ciccarrelli, F. Barkhof, D. Alexander, Disentangling human error from ground truth in segmentation of medical images, in: H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, H. Lin (Eds.), *Advances in Neural Information Processing Systems*, Vol. 33, Curran Associates, Inc., 2020, pp. 15750–15762.
- [19] V. S. Sheng, F. Provost, P. G. Ipeirotis, Get another label? Improving data quality and data mining using multiple, noisy labelers, in: *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '08*, Association for Computing Machinery, New York, NY, USA, 2008, p. 614–622.
- [20] D. Karger, S. Oh, D. Shah, Iterative learning for reliable crowdsourcing systems, *Advances in Neural Information Processing Systems* 24.
- [21] A. Dawid, A. Skene, Maximum likelihood estimation of observer error-rates using the em algorithm, *Journal of The Royal Statistical Society Series C-applied Statistics* 28 (1979) 20–28.

- [22] P. Smyth, U. Fayyad, M. Burl, P. Perona, P. Baldi, Inferring ground truth from subjective labelling of venus images, in: Proceedings of the 7th International Conference on Neural Information Processing Systems, NIPS'94, MIT Press, Cambridge, MA, USA, 1994, p. 1085–1092.
- [23] J. Whitehill, T.-f. Wu, J. Bergsma, J. Movellan, P. Ruvolo, Whose vote should count more: Optimal integration of labels from labelers of unknown expertise, *Advances in Neural Information Processing Systems* 22 (2009) 2035–2043.
- [24] V. C. Raykar, S. Yu, L. H. Zhao, A. Jerebko, C. Florin, G. H. Valadez, L. Bogoni, L. Moy, Supervised learning from multiple experts: Whom to trust when everyone lies a bit, in: Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09, Association for Computing Machinery, New York, NY, USA, 2009, p. 889–896.
- [25] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Vol. 1, 2005, pp. 886–893 vol. 1.
- [26] A. Shrivastava, A. Gupta, R. Girshick, Training region-based object detectors with online hard example mining, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 761–769.
- [27] R. E. Schapire, Explaining AdaBoost, in: *Empirical Inference*, Springer, 2013, pp. 37–52.
- [28] T.-Y. Lin, P. Goyal, R. Girshick, K. He, P. Dollár, Focal loss for dense object detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 2980–2988.
- [29] Y. Bengio, J. Louradour, R. Collobert, J. Weston, Curriculum learning, in: Proceedings of the 26th Annual International Conference on Machine

- Learning, ICML '09, Association for Computing Machinery, New York, NY, USA, 2009, p. 41–48.
- [30] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning (still) requires rethinking generalization, *Communications of the ACM* 64 (3) (2021) 107–115.
- [31] D. Arpit, S. Jastrzundebdski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, S. Lacoste-Julien, A closer look at memorization in deep networks, in: *Proceedings of the 34th International Conference on Machine Learning, ICML'17, JMLR.org*, 2017, p. 233–242.
- [32] Z. Zou, Z. Shi, Y. Guo, J. Ye, Object detection in 20 years: A survey, arXiv preprint arXiv:1905.05055.
- [33] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The Pascal Visual Object Classes (VOC) challenge, *International Journal of Computer Vision* 88 (2) (2010) 303–338.
- [34] G. Jocher, A. Stoken, J. Borovec, NanoCode012, A. Chaurasia, TaoXie, L. Changyu, A. V, Laughing, tkianai, yxNONG, A. Hogan, lorenzomammana, AlexWang1900, J. Hajek, L. Diaconu, Marc, Y. Kwon, oleg, wanghaoyang0106, Y. Defretin, A. Lohia, ml5ah, B. Milanko, B. Fineran, D. Khromov, D. Yiwei, Doug, Durgesh, F. Ingham, ultralytics/yolov5: v5.0 - YOLOv5-P6 1280 models, AWS, Supervise.ly and YouTube integrations (Apr. 2021).
- [35] L. N. Smith, Cyclical learning rates for training neural networks, in: *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2017, pp. 464–472.
- [36] S. Ruder, An overview of gradient descent optimization algorithms, arXiv preprint arXiv:1609.04747.

- [37] M. Tan, R. Pang, Q. V. Le, Efficientdet: Scalable and efficient object detection, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 10781–10790.
- [38] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, arXiv preprint arXiv:1711.05101.