

Knowledge Distillation with Adaptive Asymmetric Label Sharpening for Semi-supervised Fracture Detection in Chest X-rays

Yirui Wang¹, Kang Zheng¹, Chi-Tung Chang², Xiao-Yun Zhou¹, Zhilin Zheng³,
Lingyun Huang³, Jing Xiao³, Le Lu¹, Chien-Hung Liao², and Shun Miao¹

¹ PAII Inc., Bethesda, MD, USA

² Chang Gung Memorial Hospital, Linkou, Taiwan, ROC

³ Ping An Technology, Shenzhen, China

Abstract. Exploiting available medical records to train high performance computer-aided diagnosis (CAD) models via the semi-supervised learning (SSL) setting is emerging to tackle the prohibitively high labor costs involved in large-scale medical image annotations. Despite the extensive attentions received on SSL, previous methods failed to 1) account for the low disease prevalence in medical records and 2) utilize the image-level diagnosis indicated from the medical records. Both issues are unique to SSL for CAD models. In this work, we propose a new knowledge distillation method that effectively exploits large-scale image-level labels extracted from the medical records, augmented with limited expert annotated region-level labels, to train a rib and clavicle fracture CAD model for chest X-ray (CXR). Our method leverages the teacher-student model paradigm and features a novel adaptive asymmetric label sharpening (AALS) algorithm to address the label imbalance problem that specially exists in medical domain. Our approach is extensively evaluated on all CXR ($N = 65,845$) from the trauma registry of Chang Gung Memorial Hospital over a period of 9 years (2008-2016), on the most common rib and clavicle fractures. The experiment results demonstrate that our method achieves the state-of-the-art fracture detection performance, *i.e.*, an area under receiver operating characteristic curve (AUROC) of 0.9318 and a free-response receiver operating characteristic (FROC) score of 0.8914 on the rib fractures, significantly outperforming previous approaches by an AUROC gap of 1.63% and an FROC improvement by 3.74%. Consistent performance gains are also observed for clavicle fracture detection.

Keywords: Knowledge Distillation · Adaptive Asymmetric Label Sharpening · Semi-supervised Learning · Fracture Detection · Chest X-ray.

1 Introduction

Computer-aided diagnosis (CAD) of medical images has been extensively studied in the past decade. In recent years, substantial progress has been made

in developing deep learning-based CAD systems to diagnose a wide range of pathologies, *e.g.*, lung nodule diagnosis in chest computed tomography (CT) images [17], mass and calcification characterization in mammography [12], bone fracture detection/classification in radiography [16]. The state-of-the-art CAD solutions are typically developed based on large-scale expert annotations (*e.g.*, 128,175 labeled images for diabetic retinopathy detection [3], 14,021 labeled cases for skin condition diagnosis [10]). However, the labor cost of large-scale annotations in medical area is prohibitively high due to the required medical expertise, which hinders the development of deep learning-based CAD solutions for applications where such large-scale annotations are not yet available. In this work, we aim to develop a cost-effective semi-supervised learning (SSL) solution to *train a reliable, robust and accurate fracture detection model for chest X-ray (CXR) using limited expert annotations and abundant clinical diagnosis records.*

While expert annotations are expensive to obtain, medical records can often be efficiently/automatically collected retrospectively at large scale from a hospital’s information system. Motivated by the availability of retrospective medical records, a few large-scale X-ray datasets with natural language processing (NLP) generated image-level labels are collected and publicly released, *e.g.*, ChestXray-14 [15], CheXpert [5]. Previous works have subsequently investigated weakly-supervised learning to train CAD models using purely image-level labels [11,16]. However, since the image-level labels lack localization supervision, these methods often cannot deliver sufficient diagnosis and localization accuracy for clinical usage [7]. In addition, the public CXR datasets rely only on radiology reports, which are known to have substantial diagnostic errors and inaccuracies [2].

A more promising and practical strategy for training CAD models is to use large-scale image-level labels extracted from the *clinical diagnosis reports* with a small number of expert annotated region-level labels. Different from radiology diagnoses made by the radiologist based on a single image modality, clinical diagnoses are made by the primary doctor considering all sources of information, *e.g.*, patient history, symptoms, multiple image modalities. Therefore, clinical diagnoses offer more reliable image-level labels for training CAD models. Previous SSL methods (*e.g.*, Π -model [6], Mean Teacher [13], Mix-Match [1]) have studied a similar problem setup, *i.e.*, training classification/segmentation models using a combination of labeled and unlabeled images. However, these general-purpose SSL methods assume that no label information is given for the *unlabeled* set. Therefore, they fail to take advantage of the clinical diagnosis reports that are available in our application. In addition, there is a unique data imbalance challenge in training CAD models using clinical diagnoses. In particular, due to the low prevalence of fractures in CXRs, the image-level diagnostic labels are imbalanced toward more negative (*e.g.*, 1:10 ratio). The region-level labeled positives and image-level diagnostic negatives are even more imbalanced (*e.g.*, 1:100 ratio). As a result, a specifically-designed SSL method is required to fully exploit the clinical diagnoses with imbalanced data distribution to effectively train CAD models.

To bridge this gap, we propose an effective SSL solution for fracture detection in CXR that better accounts for the imbalanced data distribution and exploits the image-level labels of the unannotated data. We adopt the teacher-student mechanism, where a teacher model is employed to produce pseudo ground-truths (GTs) on the image-level diagnostic positive images for supervising the training of the student model. Different from previous knowledge distillation methods where the pseudo GTs are directly used or processed with symmetric sharpening/softening, we propose an *adaptive asymmetric label sharpening (AALS)* to account for the teacher model’s low sensitivity caused by the imbalanced data distribution and to encourage positive detection responses on the image-level positive CXRs. The proposed method is evaluated on a real-world scenario dataset of all ($N = 65,843$) CXR images taken in the trauma center of Chang Gung Memorial Hospital from year 2008 to 2016. Experiments demonstrate that our method reports an area under receiver operating characteristic curve (AU-ROC) of 0.9318/0.9646 and an free-response receiver operating characteristic (FROC) score of 0.8914/0.9265 on the rib/clavicle fracture detection. Compared to state-of-the-art methods, our method significantly improves the AUROC by 1.63%/0.86% and the FROC by 3.74%/3.81% on rib/clavicle fracture detection, respectively.

2 Method

Problem Setup We develop a fracture detection model using a combination of image-level and region-level labeled CXRs. While the image-level labels can be obtained efficiently at a large scale by mining a hospital’s image archive and clinical records, the region-level labels are more costly to obtain as they need to be manually annotated by experts. We collected 65,845 CXRs from the trauma registry of a medical center. Diagnosis code and keyword matching of the clinical records are used to obtain image-level labels, resulting in 6,792 positive and 59,051 negative CXRs. Among positive CXRs, 808 CXRs with positive diagnosis are annotated by experts to provide region-level labels in the form of bounding-box. The sets of region-level labeled, image-level positive and image-level negative CXRs are denoted by \mathcal{R} , \mathcal{P} and \mathcal{N} , respectively. Our method aims to *effectively exploit both the region-level labels and the image-level labels under extremely imbalanced positive/negative ratio*.

2.1 Knowledge Distillation Learning

Similar to recent CAD methods [7], we train a neural network to produce a probability map that indicates the location of the detected fracture. Since the shape and scale of fractures can vary significantly, we employ feature pyramid network (FPN) [8] with a ResNet-50 backbone to tackle the scale variation challenge by fusing multi-scale features. The training consists of two steps: 1) supervised pre-training and 2) semi-supervised training. In the pre-training step, a fracture detection model is trained via supervised learning using $\mathcal{R} \cup \mathcal{N}$. In the semi-supervised training step, \mathcal{P} are further exploited to facilitate the training.

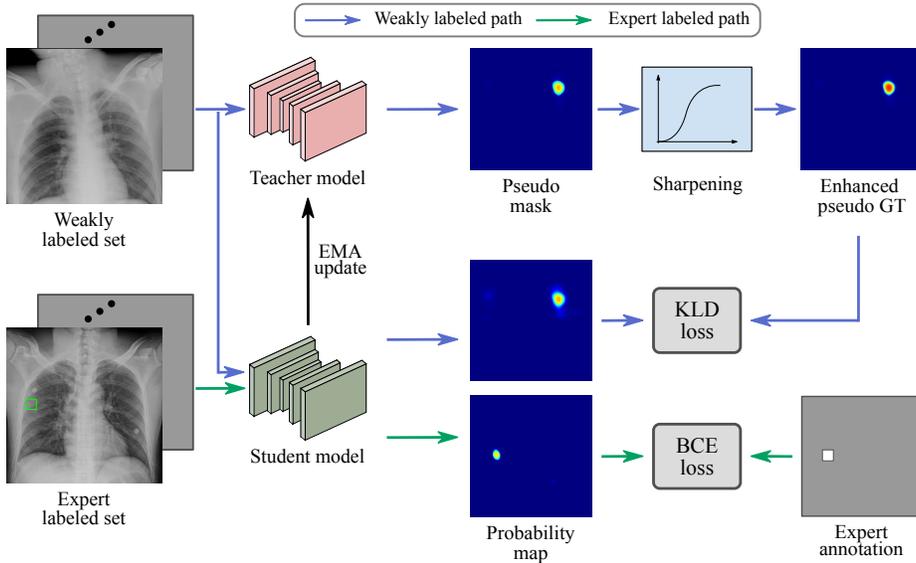


Fig. 1. An overview of the proposed knowledge distillation with AALS. The student model is trained via back-propagation. The teacher model is updated by the EMA.

Supervised pre-training We train the network using only CXRs in \mathcal{R} and \mathcal{N} , where pixel-level supervision signals can be generated. Specifically, for CXRs in \mathcal{R} , GT masks are generated by assigning *one* to the pixels within the bounding-boxes and *zero* elsewhere. For CXRs in \mathcal{N} , GT masks with all *zeros* are generated. During training, we use the pixel-wise binary cross-entropy (BCE) loss between the predicted probability map and the generated GT mask, written as:

$$\mathcal{L}_{sup} = \sum_{x \in (\mathcal{R} \cup \mathcal{N})} \text{BCE}(f_{\theta}(x), y), \quad (1)$$

where x and y denote the CXR and its pixel-level supervision mask. $f_{\theta}(x)$ denotes the probability map output of the network parameterized by θ . Due to the extreme imbalance between \mathcal{R} and \mathcal{N} (e.g., 808 vs. 59,861), the pre-trained model tends to have a low detection sensitivity, i.e., producing low probabilities on fracture sites.

Semi-supervised training To effectively leverage \mathcal{P} in training, we adopt a teacher-student paradigm where the student model learns from the pseudo GT produced by the teacher model on \mathcal{P} . The teacher and student models share the same network architecture, i.e., ResNet-50 with FPN, and are both initialized using the pre-trained weights from the supervised learning step. Inspired by the Mean Teacher method [13], we train the student model via back propagation and iteratively update the teacher model using the exponential moving average (EMA)

of the student model weights during training. Denoting the weights of the teacher and student models at training step t as θ'_t and θ_t , respectively, the weights of the teacher model are updated following:

$$\theta'_t = \alpha\theta'_{t-1} + (1 - \alpha)\theta_t, \quad (2)$$

where α is a smoothing coefficient to control the pace of knowledge update. α is set to 0.999 in all our experiments following [13].

CXRs in the region-level labeled set (\mathcal{R}), image-level labeled positive set (\mathcal{P}) and image-level labeled negative set (\mathcal{N}) are all used to train the teacher-student model. The training mechanism is illustrated in Fig. 1. For CXRs in \mathcal{R} and \mathcal{N} , the same supervised loss \mathcal{L}_{sup} is used. For CXRs in \mathcal{P} , the teacher model is applied to produce a pseudo GT map, which is further processed by an AALS operator. The sharpened pseudo GT of image x is denoted as

$$y' = S(f_{\theta'_t}(x)), \quad (3)$$

where $f_{\theta'_t}$ denotes the teacher model at the t -th step, $S(\cdot)$ denotes AALS. The KL divergence between the sharpened pseudo GT y' and the student model's prediction $f_{\theta_t}(x)$ is calculated as an additional loss:

$$\mathcal{L}_{semi} = \sum_{x \in \mathcal{P}} \text{KLDiv}(S(f_{\theta'_t}(x)), f_{\theta_t}(x)). \quad (4)$$

The total loss used to train the student network is

$$\mathcal{L} = \mathcal{L}_{sup} + \mathcal{L}_{semi}. \quad (5)$$

2.2 Adaptive Asymmetric Label Sharpening

In previous knowledge distillation models, the pseudo GTs are produced on unlabeled data to supervise the student model. Since no prior knowledge is given for the unlabeled data, the pseudo GTs are either directly used [13], or processed with symmetric softening [4] or sharpening [1]. In our problem setup, we have important prior knowledge that can be exploited: 1) image-level positive CXRs contain visible fracture sites, 2) due to the imbalanced positive/negative ratio, the pseudo GT tends to have low sensitivity (*i.e.*, low probabilities at fracture sites). Therefore, when the maximum value of the pseudo GT map is low, we are motivated to enhance the activation via AALS:

$$S(y') = \text{expit}(a \cdot \text{logit}(y') + (1 - a) \cdot \text{logit}(t)), \quad (6)$$

where $\text{expit}(\cdot)$ and $\text{logit}(\cdot)$ denote Sigmoid function and its inverse. a and t control the strength and center of the sharpening operator, respectively. The effects of a and t are shown in Fig. 2. Since the asymmetric sharpening aims to enhance the low probabilities in the pseudo GT, $t < 0.5$ should be used ($t = 0.4$ is used in our experiments). Since there are still many fracture sites missed in

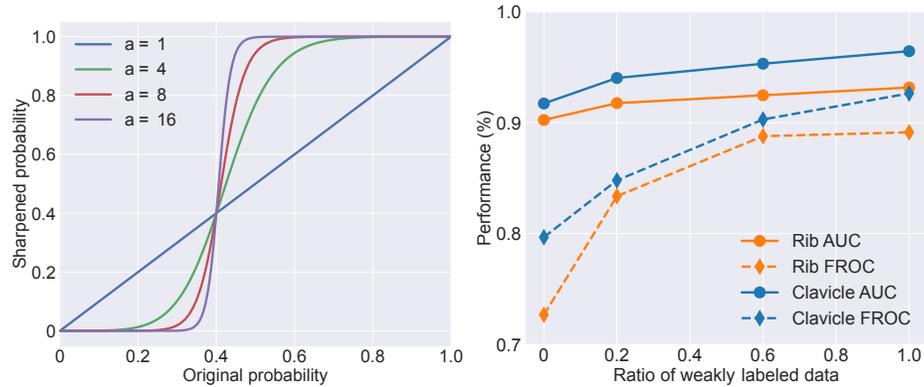


Fig. 2. Asymmetric label sharpening function at $t = 0.4$ with different a . **Fig. 3.** Model performance using a subset of \mathcal{P} .

y' (*i.e.* with low probability values) due to the imbalanced training data, we use $\max(S(y'), y')$ as the label sharpening function in our final solution to avoid over penalization of the student model’s activation on fracture sites with low probability values in y' .

The sharpening strength a is dynamically selected based on the maximum probability in the pseudo GT map, written as:

$$a = a_0 - (a_0 - 1)y'_{max}, \quad (7)$$

where y'_{max} is the maximum probability in the pseudo GT map, a_0 is a hyperparameter that controls the largest sharpening strength allowed. The sharpening strength a is negatively correlated with the maximum probability y'_{max} . When y'_{max} approaches 1, a approaches its minimum value 1, making $S(\cdot)$ an identity mapping. When y'_{max} decreases, a increases toward a_0 , leading to stronger sharpening of the pseudo GT. A dynamic a is required because the sharpening operator is asymmetric. If a constant $a > 1$ is used, the sharpening operation will always enlarge the activation area in the pseudo GT map, which drives the model to produce probability maps with overly large activation areas. With the adaptive sharpening strength, when a fracture site is confidently detected in a CXR (*i.e.*, y'_{max} approaches 1), the sharpening operation degenerates to identity mapping to avoid consistently expanding the activation area.

2.3 Implementation Details

We trained our model on a workstation with a single Intel Xeon E5-2650 v4 CPU @ 2.2 GHz, 128 GB RAM, 4 NVIDIA Quadro RTX 8000 GPUs. All methods are implemented in Python 3.6 and PyTorch v1.6. We use the ImageNet pre-trained weights to initialize the backbone network. Adam optimizer is employed in all methods. A learning rate of $4e - 5$, a weight decay of 0.0001 and a batch

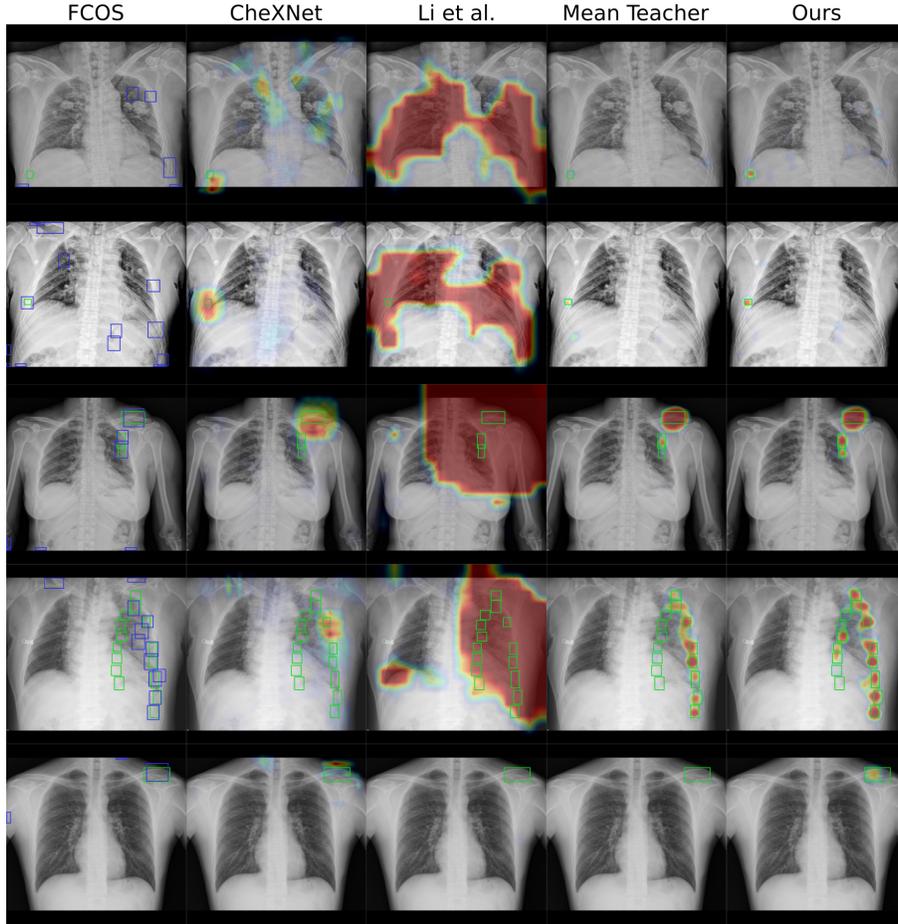


Fig. 4. Examples of the fracture detection results. GT and FCOS detected fracture bounding-boxes are shown in green and blue colors.

size of 48 are used to train the model for 25 epochs. All images are padded to square and resized to 1024×1024 for network training and inference. We randomly perform rotation, horizontal flipping, intensity and contrast jittering to augment the training data. The trained model is evaluated on the validation set after every training epoch, and the one with the highest validation AUROC is selected as the best model for inference.

3 Experiments

3.1 Experimental Settings

Dataset We collected 65,843 CXRs of unique patients that were admitted to the trauma center of Chang Gung Memorial Hospital from year 2008 to 2016. Based

Table 1. Fracture classification and localization performance comparison with state-of-the-art models. AUROC is reported for classification performance. FROC score is reported for localization performance.

Method	Rib fracture		Clavicle fracture	
	AUROC	FROC	AUROC	FROC
CheXNet [11]	0.8867	-	0.9555	-
RetinaNet [9]	0.8609	0.4654	0.8610	0.7985
FCOS [14]	0.8646	0.5684	0.8847	0.8471
Li <i>et al.</i> [7]	0.8446	-	0.9560	-
II-Model [6]	0.8880	0.7703	0.9193	0.8536
Temporal Ensemble [6]	0.8924	0.7915	0.9132	0.8204
Mean Teacher [13]	0.9155	0.8540	0.9474	0.8884
Supervised pre-training	0.9025	0.7267	0.9174	0.7967
Ours	0.9318	0.8914	0.9646	0.9265
	(+1.63%)	(+3.74%)	(+0.86%)	(+3.81%)

on the clinical diagnosis records, the CXRs are assigned image-level labels for rib and clavicle fractures. In total, we obtained 6,792 ($\mathcal{R} \cup \mathcal{P}$) CXRs with positive label for at least one type of fracture and 59,051 (\mathcal{N}) CXRs with negative label for both fracture types. 808 (\mathcal{R}) image-level positive CXRs are randomly selected for expert annotation by two experienced trauma surgeons. The annotations are confirmed by the best available information, including the original CXR images, radiologist reports, clinical diagnoses, and advanced imaging modality findings (if available). All experiments are conducted using five-fold cross-validation with a 70%/10%/20% training, validation, and testing split, respectively.

Evaluation metrics We evaluate both fracture classification and localization performances. The widely used classification metric AUROC is used to assess classification performance. For object detection methods, the maximum classification score of all predicted bounding-boxes is taken as the classification score. For methods producing probability map, the maximum value of the probability map is taken as the classification score. We also assess the fracture localization performance of different methods. Since our method only produces probability map, standard FROC metric based on bounding-box predictions is infeasible. Thus, we report a modified FROC metric to evaluate the localization performance of all compared methods. A fracture site is considered as recalled if the center of its bounding-box is activated. And the activated pixels outside bounding-boxes are regarded as false positives. Thus, the modified FROC measures the fracture recall and the average ratio of false positive pixels per image. To calculate the modified FROC for object detection methods, we convert their predicted bounding-boxes into a binary mask using different thresholds, with the pixels within the predicted box as positive, and the pixels outside the box as

negative. To quantify the localization performance, we calculate an FROC score as an average of recalls at ten false positive ratios from 1% to 10%.

Compared methods We compare the proposed method with baseline methods in the following three categories. 1) **Weakly-supervised methods:** We evaluate *CheXNet* [11], a representative state-of-the-art X-ray CAD method trained purely using image-level labels. 2) **Object detection methods:** We evaluate two state-of-the-art object detection methods: an anchor-based detector *RetinaNet* [9] and an anchor-free detector *FCOS* [14]. 3) **Semi-supervised methods:** We evaluate three popular knowledge distillation methods, *II-Model* [6], *Temporal Ensemble* [6] and *Mean Teacher* [13], and a state-of-the-art medical image SSL method by Li *et al.* [7]. For all evaluated methods, ResNet-50 is employed as the backbone network. FPN is employed in the two detection methods, RetinaNet and FCOS.

3.2 Comparison with Baseline Methods

Table 1 summarizes the quantitative results of all compared methods and the proposed method. On the more challenging rib fracture detection task, Mean Teacher is the most competitive baseline method, measuring an AUROC of 0.9155 and an FROC score of 0.8540. Our proposed method measures an AUROC of 0.9318 and an FROC score of 0.8914, which significantly outperforms Mean Teacher by a 1.63% gap on the AUROC, and a 3.74% gap on the FROC score. The ROC and FROC curves of the evaluated methods are shown in Fig. 5. On the easier clavicle fracture detection task, CheXNet and Li *et al.* [7] report the highest AUROCs (*i.e.*, above 0.95) among the baseline methods. Mean Teacher delivers the strongest FROC score of 0.8884 among the baseline methods. Our proposed method also outperforms all baseline methods on the clavicle fracture detection task, reporting an AUROC of 0.9646 and an FROC of 0.9265.

We note that the three knowledge distillation methods, II-Model, Temporal Ensemble and Mean Teacher, perform stronger than the supervised detection methods. The advantage is more significant on the easier clavicle fracture detection task. This is mainly because clavicle fractures have simpler geometric property and similar visual patterns, which knowledge distillation methods can effectively learn from the pseudo GT of unlabeled data. However, on the more complex rib fracture detection, the advantage of knowledge distillation methods is much less significant. Due to the complex visual patterns of rib fracture and the limited region-labeled positive data, the pseudo GT maps have a low sensitivity (*i.e.*, the supervised pre-trained model reports a low FROC score of 0.7267), which limits the knowledge transferred to the distilled model. Using the proposed AALS, our method effectively transfers more knowledge to the student model, hence achieving significantly improved performance compared to the previous knowledge distillation methods.

We observed that CheXNet and Li *et al.* [7] significantly outperform baseline knowledge distillation methods on the clavicle fracture AUROC metric, but

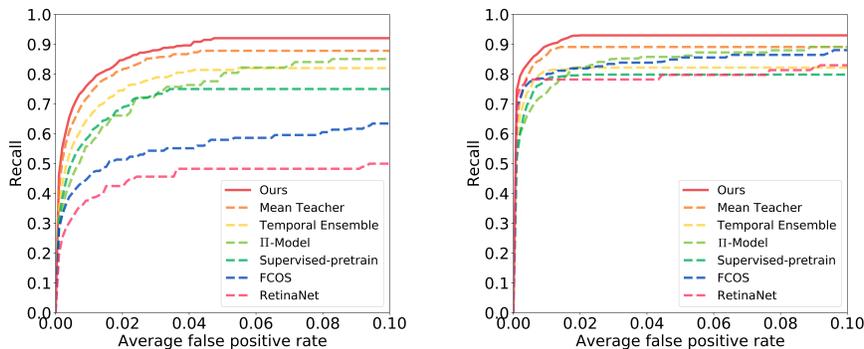


Fig. 5. FROC curves of rib fracture (left) and clavicle fracture (right) detection results using different methods.

no performance advantage is observed on the rib fracture AUROC. This is because CheXNet and Li *et al.* [7] specifically use the positive image-level label, while the baseline knowledge distillation methods do not. In particular, CheXNet is trained via weakly-supervised learning purely using image-level labels. Li *et al.* [7] exploits image-level positive labels in a multi-instance learning manner. In contrast, the baseline knowledge distillation methods treat the image-level positive images as unlabeled data. While weakly-supervised learning and multi-instance learning are effective on learning the simpler clavicle fractures, they are less effective on more complex rib fractures. In addition, CheXNet and Li *et al.* [7] also produce poor localization performances. CheXNet provides localization visualization via class activation maps (CAM). Since the CAM values are not comparable across images, the FROC cannot be calculated for CheXNet results. As Li *et al.* [7] consistently produces overly large activation areas, it does not report meaningful FROC scores. For both CheXNet and Li *et al.* [7], we qualitatively verified that their localization performances are worse than other methods, as demonstrated by the examples shown in Fig. 4.

3.3 Ablation Study

We validate our proposed AALS by conducting experiments with different sharpening strengths a_0 and centers t , respectively. First, to analyze the effect of the label sharpening center t , we evaluate AALS with $t = 0.2, 0.3, 0.4, 0.5$ and summarize the results in Table 2. Using $t = 0.4$ achieves the best detection performance, measuring the highest/second highest AUROC score of 0.9318/0.9646, and the highest FROC score of 0.8914/0.9265, on rib/clavicle fracture detection. Note that for clavicle fracture classification, the best AUROC score of 0.9661 achieved at $t = 0.2$ is only marginally better than that of $t = 0.4$. The sharpening center behaves as a trade-off between sensitivity and specificity. We note that our method consistently outperforms baseline methods using all four t values. Sec-

Table 2. Study of the sharpening bias. **Table 3.** Study of the sharpening strength.

t	Rib fracture		Clavicle fracture		a_0	Rib fracture		Clavicle fracture	
	AUROC	FROC	AUROC	FROC		AUROC	FROC	AUROC	FROC
0.2	0.9289	0.8902	0.9661	0.9236	1	0.9222	0.8783	0.9550	0.9036
0.3	0.9261	0.8888	0.9611	0.9168	4	0.9318	0.8914	0.9646	0.9265
0.4	0.9318	0.8914	0.9646	0.9265	8	0.9283	0.8884	0.9606	0.9090
0.5	0.9271	0.8848	0.9577	0.9106	16	0.9302	0.8911	0.9620	0.9185

ond, we fix the center $t = 0.4$ and evaluate $a_0 = 1, 4, 8, 16$ to study the impact of the sharpening strength. As summarized in Table 3, label sharpening with strength $a_0 = 4$ results in the best detection performance. For $a_0 = 1$, no label sharpening is applied, which results in degraded performance. For $a_0 = 8, 16$, the label sharpening becomes overly aggressive (as shown in Fig. 2), which also causes false positives in sharpened pseudo GT and hence slight performance degradation.

We further conduct an experiment to study the involvement of image-level positive set \mathcal{P} . Figure 3 shows the classification and detection performances for rib and clavicle using a subset of \mathcal{P} with different ratios (0%, 20%, 60%, 100%), where 0% and 100% correspond to the supervised pre-training student model and the proposed method, respectively. We observe that larger \mathcal{P} improves both the classification AUROC and detection FROC scores. This verifies the motivation of our method that CAD model training can benefit from utilizing image-level labels from clinical diagnoses. It also suggests a potential of our method to further improve its performance by incorporating more data with clinical diagnoses without additional annotation efforts.

4 Conclusion

In this paper, we introduced a specifically-designed SSL method to exploit both limited expert annotated region-level labels and large-scale image-level labels mined from the clinical diagnoses records for training a fracture detection model on CXR. We demonstrated that by accounting for the imbalanced data distribution and exploiting the clinical diagnoses, the proposed AALS scheme can effectively improve the effectiveness of knowledge distillation on only image-level labeled data. On a large-scale real-world scenario dataset, our method reports the state-of-the-art performance and outperforms previous methods by substantial margins. Our method offers a promising solution to exploit potentially unlimited and automatically mined clinical diagnosis data to facilitate CAD model training.

References

1. Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A.: Mixmatch: A holistic approach to semi-supervised learning. In: *Advances in Neural Information Processing Systems*. pp. 5049–5059 (2019) [2](#), [5](#)
2. Brady, A.P.: Error and discrepancy in radiology: inevitable or avoidable? Insights into imaging **8**(1), 171–182 (2017) [2](#)
3. Gulshan, V., Peng, L., Coram, M., Stumpe, M.C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., Cuadros, J., et al.: Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama* **316**(22), 2402–2410 (2016) [2](#)
4. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015) [5](#)
5. Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghighi, B., Ball, R., Shpanskaya, K., et al.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In: *AAAI*. vol. 33, pp. 590–597 (2019) [2](#)
6. Laine, S., Aila, T.: Temporal ensembling for semi-supervised learning. arXiv preprint arXiv:1610.02242 (2016) [2](#), [8](#), [9](#)
7. Li, Z., Wang, C., Han, M., Xue, Y., Wei, W., Li, L.J., Fei-Fei, L.: Thoracic disease identification and localization with limited supervision. In: *CVPR*. pp. 8290–8299 (2018) [2](#), [3](#), [8](#), [9](#), [10](#)
8. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2117–2125 (2017) [3](#)
9. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *ICCV*. pp. 2980–2988 (2017) [8](#), [9](#)
10. Liu, Y., Jain, A., Eng, C., Way, D.H., Lee, K., Bui, P., Kanada, K., de Oliveira Marinho, G., Gallegos, J., Gabriele, S., et al.: A deep learning system for differential diagnosis of skin diseases. *Nature Medicine* pp. 1–9 (2020) [2](#)
11. Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C., Shpanskaya, K., et al.: Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. arXiv preprint arXiv:1711.05225 (2017) [2](#), [8](#), [9](#)
12. Shen, L., Margolis, L.R., Rothstein, J.H., Fluder, E., McBride, R., Sieh, W.: Deep learning to improve breast cancer detection on screening mammography. *Scientific reports* **9**(1), 1–12 (2019) [2](#)
13. Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In: *NeurIPS*. pp. 1195–1204 (2017) [2](#), [4](#), [5](#), [8](#), [9](#)
14. Tian, Z., Shen, C., Chen, H., He, T.: Fcos: Fully convolutional one-stage object detection. In: *ICCV*. pp. 9627–9636 (2019) [8](#), [9](#)
15. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2097–2106 (2017) [2](#)
16. Wang, Y., Lu, L., Cheng, C.T., Jin, D., Harrison, A.P., Xiao, J., Liao, C.H., Miao, S.: Weakly supervised universal fracture detection in pelvic x-rays. In: *MICCAI*. pp. 459–467. Springer (2019) [2](#)

17. Xie, Y., Xia, Y., Zhang, J., Song, Y., Feng, D., Fulham, M., Cai, W.: Knowledge-based collaborative deep learning for benign-malignant lung nodule classification on chest ct. *IEEE TMI* **38**(4), 991–1004 (2018) [2](#)