# Cross Chest Graph for Disease Diagnosis with Structural Relational Reasoning

Gangming Zhao, Baolian Qi, Jinpeng Li, [*†‡]

## Abstract

*Locating lesions is important in the computer-aided diagnosis of X-ray images. However, box-level annotation is time-consuming and laborious. How to locate lesions accurately with few, or even without careful annotations is an urgent problem. Although several works have approached this problem with weakly-supervised methods, the performance needs to be improved. One obstacle is that general weakly-supervised methods have failed to consider the characteristics of X-ray images, such as the highly-structural attribute. We therefore propose the Cross-chest Graph (CCG), which improves the performance of automatic lesion detection by imitating doctor's training and decision-making process. CCG models the intra-image relationship between different anatomical areas by leveraging the structural information to simulate the doctor's habit of observing different areas. Meanwhile, the relationship between any pair of images is modeled by a knowledge-reasoning module to simulate the doctor's habit of comparing multiple images. We integrate intra-image and inter-image information into a unified end-to-end framework. Experimental results on the NIH Chest-14 database (112,120 frontal-view X-ray images with 14 diseases) demonstrate that the proposed method achieves state-of-the-art performance in weakly-supervised localization of lesions by absorbing professional knowledge in the medical field.*

## 1. Introduction

Chest radiographs are a type of medical images that can be conveniently acquired for disease diagnosis. With the rapid development of deep learning, automatic disease detection in chest X-ray images has become an important task in the computer-aided diagnosis. Deep convolutional neural networks (DCNN) have been widely applied in many computer vision tasks, such as image classification [6, 18]

---
[*] *Gangming Zhao and Baolian Qi contributed equally to this work.*

[†] Gangming Zhao is with the Department of Computer Science, The University of Hong Kong, Hong Kong.

[‡] Qibao Lian and Jinpeng Li are with University of Chinese Academy of Sciences, Beijing, China
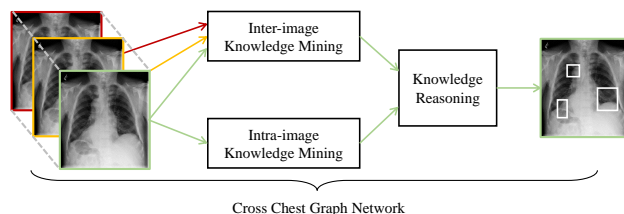
Figure 1. CCG network models the intra-image relationship between different anatomical areas by leveraging the structural information to simulate the doctor's habit of observing different areas. Meanwhile, the relationship between any pair of images is modeled by a knowledge-reasoning module to simulate the doctor's habit of comparing multiple images.

, object detection [3, 5, 16, 15, 10] and semantic segmentation [12, 17]. To achieve good performance in these tasks, substantial images with careful annotations are needed. Encouraged by the success of DCNN in computer vision, some researches have directly applied DCNN models to analyze the medical images but cannot achieve the same performance as in the natural images. The reasons lie in two folds: 1. it is expensive to acquire accurate localization or classification labels in chest X-ray images. 2. there exists much professional knowledge in medical images that DCNN cannot exploit well. Therefore, how to exploit the professional knowledge into DCNN models for solving these two questions still opens a fully challenging problem. Our work transfers the knowledge into DCNN models to reduce the problem of shortage of carefully annotated images.

Recent work paid much attention to utilize professional knowledge of chest X-ray images into DCNN frameworks. However, they just proposed a simple fused strategy to embed low-level information of chest X-ray into models, such as Liu et al. [9] utilized contrastive learning to provide more localization information with the help of healthy images. Zhao et al. [22] proposed to exploit the contralateral information of chest X-ray via a simple fusion module. These methods only exploit the apparent information of chest-Xray images. They all overlooked the inner structure information of chest X-rays. Therefore, they cannot apply their methods into real applications.

In this paper, we propose a Cross Chest Graph Network (CCG-Net) as shown in Fig 1, which firstly utilizes deep expert knowledge to automatical detect disease in chest X-ray images. We have known that medical experts have much experience in finding out disease and how to treat patients. In fact, the actions of medical experts consist of two phases: training and decision-making processes. They pay much time to learn distinguish disease and embed their experience into the decision process. During the training process, experts would like to observe different areas and find out the relationship between any pair of images. Our CCG-Net aims to model the observation way by a knowledge-reasoning module to simulate the doctor's habit of comparing multiple images. Then we integrate intra-image and inter-image information into a unified end-to-end framework.

Inspired from the experience of medical experts, our proposed CCG-Net consists of four modules, 1. an end-to-end framework for deciding where and what is a disease, 2. a inter-image relation module, which formulates the training process of medical experts, to compare multiple images, 3. a intra-image knowledge learning module, which builds the local relation graph for different patches of chest X-ray images. Due to their highly structured property, every chest X-ray image can be divided into several patches, we build a patch-wise relation graph on them, 4. a knowledge reasoning module, which excavates the inner knowledge from cross-image structural features. The last three operations (2, 3, and 4) are similar to medical experts' training process, which learn intra-image and inter-image information to gain professional knowledge. The first operation embeds the learned knowledge into DCNN frameworks leading to better disease diagnosis models. Above all, our contribution consists of three folds:

- We propose CCG-Net, which is the first to formulate the medical experts' training process by building relation graphs in the intra-image and inter-image information of chest X-ray images. More generally, it provides inspiration to address medical vision tasks with much professional knowledge like in chest X-ray images.

- We divide the experts' professional actions into two stages including training and decision-making processes. In addition, we utilize intra-image and inter-image relation to learn much professional knowledge that would be embedded in an end-to-end detection framework.

- We achieve state-of-the-art results on the localization of NIH ChestX-ray14.

## 2. Related Work

### 2.1. Disease Detection

Object detection is one of the most important computer vision tasks, aiming to localize and classify. Due to their strong feature representation ability, DCNN achieved much progress in object detection tasks. For detection tasks, DCNN methods consist of two style framework: 1. two-stage models, such as RCNN series [16], 2. one-stage models, such as YOLO [15] and SSD [10]. However, for disease detection, because of the shortage in careful annotations, traditional detection framework cannot directly be applied in chest X-ray images. Besides, since there is much distortion caused by other chest X-ray tissues, such low contrast also causes the difficulty of disease finding.

Weakly supervised object detection (WSOD) can be considered as an effective method to solve these problems. Based on CAM [23], researchers proposed many techniques to use only image-level labels to detect objects. Although there is no enough detection supervision, WSOD still achieved much progress. However, researchers still face a big challenge when it comes to disease detection in medical images. the existence of much professional knowledge greatly limits the development of the applications of DCNN in medical fields. Therefore, in this paper, we are inspired by the experts' learning and decision processes to propose CCG-Net, which not only exploits a larger amount of knowledge in chest X-ray images but also builds a unified framework to detect disease in an end-to-end style.

### 2.2. Knowledge-based Disease Diagnosis

Automatical disease diagnosis is a key problem in medical fields. However, due to the shortage of careful annotations and the existence of much professional knowledge, DCNN methods cannot achieve a good performance in medical tasks, especially such a tough problem: disease detection in chest X-ray images. To exploit medical knowledge and embed it into DCNN frameworks, researchers paid much effort to utilize medical experts' experience for disease diagnosis. Wang et al. [20] firstly proposed a carefully annotated chest X-ray dataset and led to a series of work that focuses on using image-level labels to localize the disease. Li et al. [8] integrated classification and localization in a whole framework with two multi instance-level losses and performed better. Liu et al. [9] improved their work to propose contrastive learning of paired samples, which utilizes healthy images to provide more localization information for disease detection. Zhao et al. [22] proposed to utilize the symmetry information in a chest X-ray to improve the disease localization performance. Besides, many works applied relation knowledge models to chest X-ray diagnosis. Ypsilantis et al. [21], Pesce et al. [14], and Guan et al. [4] proposed to build a relation attention model fusing DCNN
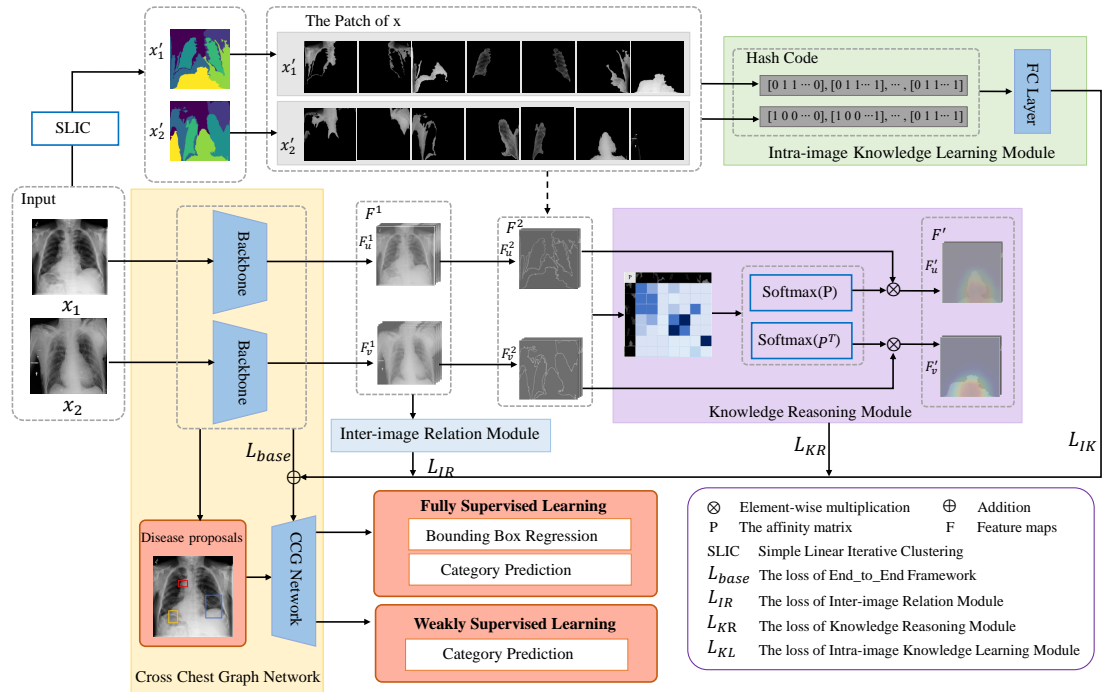
Figure 2. The network consists of four modules: 1. an end-to-end framework for disease detection under weakly-supervised settings, 2. the inter-image relation module among different samples, 3. the intra-image knowledge learning module based on the thoracic spatial structure, 4. the knowledge reasoning module mining cross-image structural features. Our four modules are tightly related and can be easily integrated into an end-to-end framework.

models achieved much progress. Li et al. [7] proposed a knowledge-graph based on medical reports and images to determine the dependencies among chest X-ray images. Cheng et al. [11] also proposed a new total strongly supervised dataset for tuberculosis detection. However, they all overlooked the structural relation among chest X-ray images. In this paper, we propose to build a structural, relational graph for disease detection under weakly supervised scenarios in chest X-ray images. Specifically, we build the global and local graph in chest X-ray via three modules: 1. a inter-image relation module, 2. a intra-image knowledge learning module, 3. knowledge reasoning module. Furthermore, we integrate three modules into an end-to-end framework to jointly train our network. Our proposed three relational modules provide better supervision since we exploit the local structural knowledge and global relation among different samples.

## 3. Method

### 3.1. Overview

Given the images $X = \{x_1, x_2, ..., x_n\}$. Our proposed framework consists of four modules:

- The End to End Framework is to localize and classify

the disease in chest X-ray images. In our paper, we utilize the same multi-instance level losses used in [9] and [8].

- Inter-image Relation Module, which includes a learnable matrix $G \in R^{n \times n}$. We also use a contrast-constrained loss to share similar information of $X$ and exploit their contrasted structural knowledge. We build a cross-sample graph for them to exploit the dependencies among different samples. The graph $G \in R^{n \times n}$ is to build the inter-image relation among sampled samples, which is a learnable matrix, and every element is initialized by $\frac{1}{n}$. The element $g_{ij}$ of $G$, $i, j \in \{1, 2, ..., n\}$, represents the similarity wight of images $x_i$ and $x_j$.

- Intra-image Knowledge Learning, which firstly acquires patch-wise features of different images. Then the network can achieve a new image graph via building a structural knowledge-based module. We denote this graph as $G_k \in R^{n \times n}$. Assumed that the number of patches are $|p_i|$ and $|p_j|$ of images $x_i$ and $x_j$. The graph $G_k$ would be calculated on using the graph $G_l \in R^{|p_i| \times |p_j|}$, which learns the relationship between different paired patches of images.

- Knowledge Reasoning Module, which is based on cross image structural knowledge. When we get the whole structural information of different images, we will utilize it to reason the inner structural dependencies among different patches in different images.

## 3.2. End to End Framework

The end to end framework is to localize and classify the disease in chest X-ray images in a coarse-grained style. More specifically, the input images $X = \{x_1, x_2, ..., x_n\}$ of the module are resized to $512 \times 512$. ResNet-50 pre-trained from the ImageNet dataset is adopted as the backbone for this module. We use the feature map $F$ after C5 (last convolutional output of 5th-stage), which is 32 times downsampled from the input image, and of size $2048 \times 16 \times 16$. Each grid in the feature map denotes the existent probability of disease. We pass $F$ through two $1 \times 1$ convolutional layers and a sigmoid layer to obtain the class-aware feature map $P$ of size $C \times H \times W$, where $C$ is the number of classes. Then we follow the paradigm used in [9], computing losses and making predictions in each channel for the corresponding class. For images with box-level annotations, if the grid in the feature map overlaps with the projected ground truth box, we assign label 1 to the grid. Otherwise, we assign 0 to it. Therefore, we use the binary cross-entropy loss as used in [9] for each grid:

$$L_i^k(P) = \sum_j -y_{ij}^k \log(p_{ij}^k) - \sum_j (1-y_{ij}^k)\log(1-p_{ij}^k) \quad (1)$$

where $k$, $i$, and $j$ are the index of classes, samples, and grids respectively. $y_{ij}^k$ denotes the target label of the grid and $p_{ij}^k$ denotes the predicted probability of the grid.

For images with only image-level annotations, we use the MIL loss used in [8].

$$\begin{aligned} L_i^k(P) = &-y_i^k \log(1 - \prod_j(1-p_{ij}^k)) \\ &-(1-y_i^k)\log(\prod_j(1-p_{ij}^k)) \end{aligned} \quad (2)$$

where $y_i^k$ denotes the target label of the image. For this end to end framework, the whole loss $L_{base}$ as shown in Fig. 2, is formulated as follows.

$$L_{base} = \sum_i \sum_k \lambda_i^k \beta_B L_i^k(P) + (1 - \lambda_i^k)L_i^k(P) \quad (3)$$

where $\lambda_i^k \in 0, 1$ denotes if the $k_{th}$ class in the $i_{th}$ sample has box annotation, and $\beta_B$ is the balance weight of the two losses and is set to 4.

## 3.3. Inter-image Relation Module

Inter-image relation is formulated as a learnable matrix $G \in R^{n \times n}$. A contrast-constrained loss is used to share

similar information of $X$ and exploit their contrasted structural knowledge, as following equation.

$$L_{IR} = \frac{\sum_{(u,v) \in G} G(u,v)D(F_u, F_v)}{n \times n} \quad (4)$$

$D(\cdot)$ is the distance metric function, where it is a Euclidean distance. $F_u$ and $F_v$ means the feature map after C5 of the image $x_u$ and $x_v$. We build a cross-sample graph for them to exploit the dependencies among different samples. The graph $G \in R^{n \times n}$ is to build the inter-image relation among sampled samples, which is a learnable matrix, and every element is initialized by $\frac{1}{n}$. The element $g_{ij}$ of $G$, $i, j \in \{1, 2, ..., n\}$, represents the similarity wight of images $x_i$ and $x_j$. G is adaptively adjusted during training processes and changes with diverse inputs to exploit the relationship fully.

## 3.4. Intra-image Knowledge Learning

Intra-image Knowledge Learning, which firstly utilizes Simple linear iterative clustering (SLIC) [1], a superpixel method to generate the patches for different images. Assumed that the patches of the image $x_i$ is $p_i = p_1^i, p_2^i, ..., p_m^i$. Then the network can achieve a new image graph via building a structural knowledge-based module with the help of $p_i, i \in 1, 2, ..., n$. We denote the graph as $G_k \in R^{n \times n}$, which is the intra-image graph between paired images $x_i$ and $x_j$. The graph $G_k$ is calculated on using the graph $G_l \in R^{|p_i| \times |p_j|}$, which learns the dependencies among different paired patches of images. Then the same contrast-constrained loss using this graph to provide more structural knowledge for the whole framework.

$$L_{IK} = \frac{\sum_{(u,v) \in G_k} G_k(u,v)D(F_u, F_v)}{n \times n} \quad (5)$$

$$G_k = W_l(G_l) \quad (6)$$

Where, $W_l$ is a fully connected layer and

$$G_l(l,p) = D^{'}(H_l, H_p^{'}), l \in 1, 2, ..., |p_i|, p \in 1, 2, ..., |p_j| \quad (7)$$

$H_l$ is the hash code [19] of the patch $p_l^i$ in the image $x_i$ and $H_p^{'}$ is the hash code of the patch $p_p^j$ in the image $x_j$. $D^{'}(\cdot)$ is the Hamming distance.

## 3.5. Knowledge Reasoning Module

In addition to previous efforts to focus on information in a whole image, we also explored the value of cross-image semantic relations in the medical object. The correlations between patches across images are emphasized, especially, the correlations between corresponding patches in two images.

Knowledge Reasoning Module focuses on the correlations of two images. After getting the feature map $F_u$ and $F_v$ of the images, the affinity matrix $P$ is firstly calculated between $F_u$ and $F_v$.

$$P = F_u^{\mathrm{T}} W_P F_v \in \mathbb{R}^{HW \times HW}$$

where the feature map $F_u \in \mathbb{R}^{C \times HW}$ and $F_v \in \mathbb{R}^{C \times HW}$ are flattened into matrix formats, and $W_P \in \mathbb{R}^{C \times C}$ is a learnable matrix. The affinity matrix $P$ represents the similarity of all pairs of patches in $F_u$ and $F_v$.

Then $P$ is normalized column-wise to get the attention map of $F_u$ for each patch in $F_v$ and row-wise to get the attention map of $F_v$ for each patch in $F_u$.

$$F_u^{'} = F_u softmax(P) \in \mathbb{R}^{C \times HW}$$
$$F_v^{'} = F_v softmax(P^{\mathrm{T}}) \in \mathbb{R}^{C \times HW}$$

where $softmax(P)$ and $softmax(P^{\mathrm{T}})$ pay attention to the similar patches of the feature map $F_u$ and $F_v$ respectively. Therefore, they can be used to enhance $F_u$ and $F_v$ respectively, so that similar patches in $F_u$ and $F_v$ are highlighted.

The cross-image method can extract more contextual information between images than using a single image. This module exploits the context of other related images to improve the reasoning ability of the feature map, which is beneficial to the localization and classification of disease in chest X-ray images. Furthermore, we exploit the enhanced feature map to calculate the new similarity between the paired images to gain a more strong supervisor.

$$L_{KR} = \frac{\sum_{(u,v) \in G_k^{'}} G_k^{'}(u,v) D(F_u^{'}, F_v^{'})}{n \times n} \quad (8)$$

The graph $G_k^{'}$ is calculated on using the graph $G_l^{'} \in R^{|p_i| \times |p_j|}$.

$$G_k^{'} = W_l^{'}(G_l^{'}) \quad (9)$$

where $W_l^{'}$ is a fully connected layer and

$$G_l^{'}(l,p) = D^{'}(P_l, P_p),$$
$$l \in \{1, 2, ..., |p_i|\}, p \in \{1, 2, ..., |p_j|\} \quad (10)$$

$P_l$ is the $l$-th feature patch of $F_u^{'}$ and $P_p$ is the $p$-th feature patch of $F_v^{'}$, respectively.

### 3.6. Training Loss

The overall loss function during the training is a weighted combination of four loss functions,

$$L_{all} = w_1 L_{base} + w_2 L_{IR} + w_3 L_{IK} + w_4 L_{KR} \quad (11)$$

where $\sum_{i=1}^{4} w_i = 1$. In our experiments, we always set $w_i = 0.25, i \in 1, 2, .., 4$.

### 3.7. Training and Test

**Training** All the models are trained on NIH chest X-ray dataset using the SGD algorithm with the Nesterov momentum. With a total of 9 epochs, the learning rate starts from 0.001 and decreases by 10 times after every 4 epochs. Additionally, the weight decay and the momentum is 0.0001 and 0.9, respectively. All the weights are initialized by pre-trained ResNet [6] models on ImageNet [2]. The mini batch size is set to 2 with the NVIDIA 1080Ti GPU. All models proposed in this paper are implemented based on PyTorch [13].

**Testing** We also use the threshold of 0.5 to distinguish positive grids from negative grids in the class-wise feature map as described in [8] and [9]. All test setting is same as [9], we also up-sampled the feature map before two last fully convolutional layers to gain a more accurate localization result.

## 4. Experiments

### 4.1. Dataset and Evaluation Metrics

**Dataset.** NIH chest X-ray dataset [20] include 112,120 frontal-view X-ray images of 14 classes of diseases. There are different diseases in each image. Furthermore, the dataset contains 880 images with 984 labeled bounding boxes. We follow the terms in [8] and [9] to call 880 images as "annotated" and the remaining 111,240 images as "unannotated". Following the setting in [9], we also resize the original 3-channel images from resolution of $1024 \times 1024$ to $512 \times 512$ without any data augmentation techniques.

**Evaluation Metrics.** We follow the metrics used in [8]. The localization accuracy is calculated by the IoU (Intersection over Union) between predictions and ground truths. Since it is a coarse-grained task, our localization predictions are discrete small rectangles. The eight diseases with ground truth boxes is reported in our paper. The localization result is regarded as correct when $IoU > T(IoU)$, where T(*) is the threshold.

### 4.2. Comparison with the State-of-the-art

In order to evaluate the effectiveness of our models for weakly supervised disease detection, we design the experiments on three sets of data and conduct a 5-fold cross-validation. In the first experiment, we use the 50% unannotated images and 80% annotated images for training, and test the models with the remaining 20% annotated images. In the second experiment, we use the 100% unannotated images and no any annotated images for training, and test the models with all annotated images. In the third experiment, we use the 100% unannotated images and 40% annotated images for training, and test the models with remaining 60% annotated images. Additionally, our experimental results are mainly compared with four methods. The

| T (IoU) | Models | Atelectasis | Cardiomegaly | Effusion | Infiltration | Mass | Nodule | Pneumonia | Pneumothorax | Mean |
|---|---|---|---|---|---|---|---|---|---|---|
| | X, Wang [20] | 0.24 | 0.46 | 0.30 | 0.28 | 0.15 | 0.04 | 0.17 | 0.13 | 0.22 |
| 0.3 | Z, Li [8] | 0.36 | **0.94** | 0.56 | 0.66 | 0.45 | 0.17 | 0.39 | **0.44** | 0.49 |
| | J, Liu [9] | **0.53** | 0.88 | 0.57 | 0.73 | **0.48** | 0.10 | 0.49 | 0.40 | 0.53 |
| | Ours | 0.44 | 0.86 | **0.68** | **0.84** | 0.47 | **0.29** | **0.67** | 0.40 | **0.60** |
| | X, Wang [20] | 0.05 | 0.18 | 0.11 | 0.07 | 0.01 | 0.01 | 0.03 | 0.03 | 0.06 |
| 0.5 | Z, Li [8] | 0.14 | 0.84 | 0.22 | 0.30 | 0.22 | 0.07 | 0.17 | 0.19 | 0.27 |
| | J, Liu [9] | **0.32** | 0.78 | 0.40 | 0.61 | 0.33 | 0.05 | 0.37 | 0.23 | 0.39 |
| | Ours | 0.27 | **0.86** | **0.48** | **0.72** | **0.53** | **0.14** | **0.58** | **0.35** | **0.49** |
| | X, Wang [20] | 0.01 | 0.03 | 0.02 | 0.00 | 0.00 | 0.00 | 0.01 | 0.02 | 0.01 |
| 0.7 | Z, Li [8] | 0.04 | 0.52 | 0.07 | 0.09 | 0.11 | 0.01 | 0.05 | 0.05 | 0.12 |
| | J, Liu [9] | 0.18 | 0.70 | 0.28 | 0.41 | 0.27 | 0.04 | 0.25 | 0.18 | 0.29 |
| | Ours | **0.20** | **0.86** | **0.48** | **0.68** | **0.32** | **0.14** | **0.54** | **0.30** | **0.44** |

Table 1. The comparison results of disease localization among the models using 50% unannotated images and 80% annotated images. For each disease, the best results are bolded.

| T (IoU) | Models | Atelectasis | Cardiomegaly | Effusion | Infiltration | Mass | Nodule | Pneumonia | Pneumothorax | Mean |
|---|---|---|---|---|---|---|---|---|---|---|
| | Z, Li [8] | 0.59 | 0.81 | 0.72 | 0.84 | 0.68 | 0.28 | 0.22 | 0.37 | 0.57 |
| 0.1 | J, Liu [9] | 0.39 | **0.90** | 0.65 | **0.85** | **0.69** | **0.38** | 0.30 | 0.39 | 0.60 |
| | Ours | **0.66** | 0.88 | **0.79** | **0.85** | **0.69** | 0.28 | **0.40** | **0.47** | **0.63** |
| | J, Liu [9] | 0.34 | 0.71 | **0.39** | 0.65 | **0.48** | **0.09** | 0.16 | 0.20 | 0.38 |
| 0.3 | Baseline | **0.36** | 0.69 | 0.35 | 0.64 | 0.44 | 0.08 | 0.02 | 0.23 | 0.35 |
| | Ours | 0.31 | **0.79** | 0.37 | **0.75** | 0.40 | 0.06 | **0.24** | **0.27** | **0.40** |
| | J, Liu [9] | **0.19** | 0.53 | **0.19** | 0.47 | **0.33** | 0.03 | **0.08** | 0.11 | 0.24 |
| 0.5 | Baseline | 0.18 | 0.51 | 0.14 | 0.47 | 0.27 | 0.03 | 0.01 | 0.12 | 0.22 |
| | Ours | **0.19** | **0.71** | 0.14 | **0.52** | 0.31 | **0.08** | 0.05 | **0.13** | **0.27** |
| | J, Liu [9] | 0.08 | 0.30 | **0.09** | 0.25 | 0.19 | 0.01 | 0.04 | 0.07 | 0.13 |
| 0.7 | Baseline | **0.11** | 0.34 | 0.06 | 0.32 | **0.20** | 0.01 | 0.00 | 0.06 | 0.14 |
| | Ours | 0.06 | **0.64** | 0.08 | **0.38** | 0.19 | 0.01 | **0.08** | **0.09** | **0.19** |

Table 2. The comparison results of disease localization among the models using 100% unannotated images and no any annotated images. For each disease, the best results are bolded.

first method is X, Wang [20], which proposes a carefully annotated chest X-ray dataset and a unified weakly supervised multi-label image classification and disease localization framework. The second method is Z, Li [8], which uses fully convolutional neural network to localize and classify the disease in chest X-ray images. The third method is J, Liu [9], which proposes contrastive learning of paired samples to provide more localization information for disease detection. The last method is our baseline model, which is a unified end-to-end framework that doesn't use our approach to locate and classify the disease.

In the first experiment, we compare the localization results of our model with [20], [8] and [9]. We can observe that our model outperforms existing methods in most cases, as shown in Table 1. Particularly, with the increase of T(IoU), our model has greater advantages over the reference models. For example, when T(IoU) is 0.3, the mean accuracy of our model is 0.60, and outperforms [20], [8] and [9] by 0.38, 0.11 and 0.07 respectively. However, when T(IoU) is 0.7, the mean accuracy of our model is 0.44, and outperforms [20], [8] and [9] by 0.43, 0.32 and 0.15 respectively. Overall, the experimental results shown in Table 1 demonstrate that our method is more accurate for disease localization and classification, which provides a great role for clinical practices.

In the second experiment, we train our model without any annotated images comparing the first experiment. Since [8] only provides the results when T(IoU) = 0.1, in order to better show the performance of our model, we add an evaluation method of T(IoU) = 0.1. It can be seen that our model outperforms [8] and [9] in most cases, as shown in Table 2. For example, when T(IoU) is 0.1, the mean accuracy of our models is 0.63, which is 0.06 higher than [8], and 0.03 higher than [9]. Furthermore, when T(IoU) is 0.7, the mean localization result of our model is 0.19, which is 0.06 higher than [8] and 0.05 higher than [9]. Compared with the baseline model, our approach performs better in most classes except for "Atelectasis" and "Nodule". The trend stays the same that at higher T(IoU), our approach demonstrates more advantages over baseline methods. The added unannotated training samples contribute more than the removed annotated ones in those classes, which implies that our approach can better utilize the unannotated samples. The overall results show that even without annotated data used for training, our approach can achieve decent localization results.

In the third experiment, we use more annotated images comparing the second experiment. We compare the localization results of our model with [9] in same data setting. It can be seen that our model outperforms [9] in most cases, as shown in Table 3. With T(IoU) = 0.3 and 0.7, our model outperforms [9] by 0.02 and 0.05 respectively. Similar improvements are achieved comparing the second experiment. Overall, the experimental results demonstrate

| T (IoU) | Models | Atelectasis | Cardiomegaly | Effusion | Infiltration | Mass | Nodule | Pneumonia | Pneumothorax | Mean |
|---|---|---|---|---|---|---|---|---|---|---|
| | J, Liu [9] | **0.55** | 0.73 | 0.55 | 0.76 | 0.48 | **0.22** | 0.39 | **0.30** | 0.50 |
| 0.3 | Baseline | 0.47 | 0.84 | 0.65 | 0.82 | 0.33 | 0.04 | **0.57** | 0.29 | 0.50 |
| | Ours | 0.49 | **0.87** | **0.66** | **0.88** | **0.48** | 0.10 | 0.51 | 0.20 | **0.52** |
| | J, Liu [9] | **0.36** | 0.57 | 0.37 | 0.62 | **0.34** | **0.13** | 0.23 | 0.17 | 0.35 |
| 0.5 | Baseline | 0.27 | 0.76 | 0.39 | 0.58 | 0.24 | 0.02 | 0.39 | **0.21** | 0.36 |
| | Ours | 0.26 | **0.80** | **0.41** | **0.67** | 0.15 | 0.06 | **0.42** | 0.18 | **0.37** |
| | J, Liu [9] | **0.19** | 0.47 | 0.20 | 0.41 | **0.22** | **0.06** | 0.12 | **0.11** | 0.22 |
| 0.7 | Baseline | 0.14 | 0.62 | 0.20 | 0.42 | 0.07 | 0.00 | 0.23 | 0.08 | 0.22 |
| | Ours | 0.18 | **0.71** | 0.20 | **0.50** | 0.20 | 0.02 | **0.29** | 0.06 | **0.27** |

Table 3. The comparison results of disease localization among the models using 100% unannotated images and 40% annotated images. For each disease, the best results are bolded.
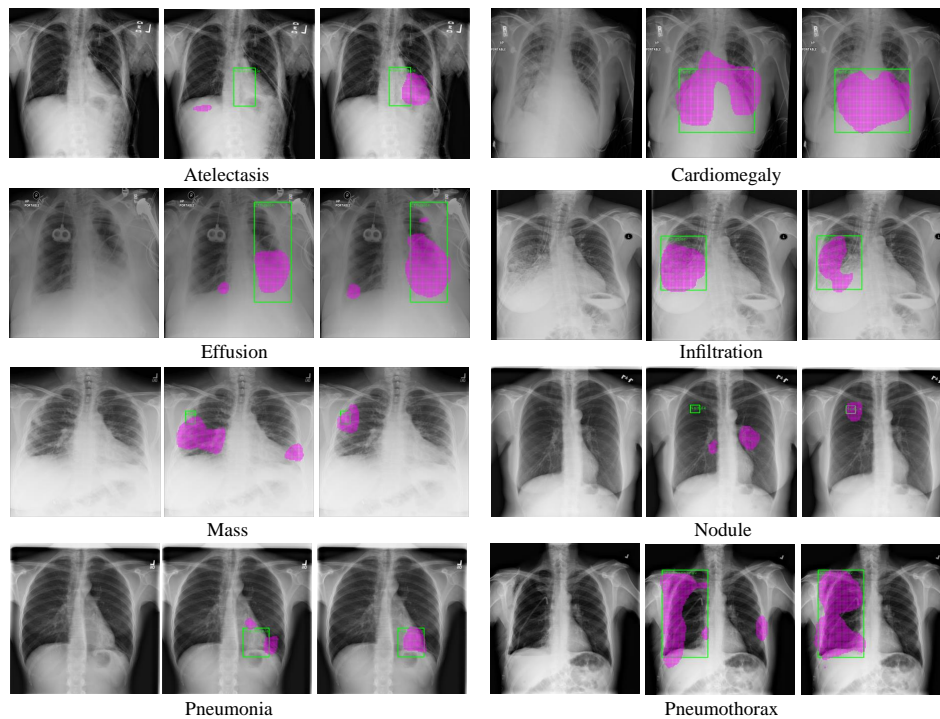


Figure 3. Visualization of the predicted results on both the baseline model and our method. The first column shows the original images, the second and third columns show baseline and our method. The green bounding box and red area mean the the ground truth and prediction.

that our method can improve the performance of models with limited annotated images.

To better demonstrate the final effect of our approach on disease localization and classification, we visualize some of typical predictions of both the baseline model and our method, as shown in Figure 3. The first column shows the original images, the second and third columns show baseline model and our method. The green bounding box and red area mean the ground truth and prediction. It can be seen that our models can predict more accurate in most cases comparing the baseline model. For example, the class "Atelectasis" and "Nodule", the localization reslut of the baseline model is completely inconsistent with the ground truth, but the localization reslut of our method is consistent with the ground truth. It shows that using the structural in-

formation of intra-image and inter-image can improve the performance of automatic lesion detection. Additionally, we also visualize the generated heatmap and ground truth of our model, as shown in Figure 4. It can be seen that the proposed method can effectively locate and classify medical images.

### 4.3. Ablation Studies

In this section, we explore the influence of different modules on our method for ablation studies. To evaluate our method more comprehensively, we build 6 models, including the model of the end to end framework (Baseline), the model with the intra-image knowledge learning (IK), the model with the inter-image relation module (IR), the model with the knowledge reasoning (KR), the model combining
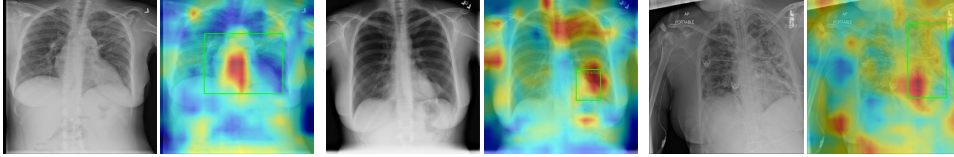
Figure 4. Visualization of the generated heatmap and ground truth of our method, where the green bounding box means the ground truth.

| Data | Models | Atelectasis | Cardiomegaly | Effusion | Infiltration | Mass | Nodule | Pneumonia | Pneumothorax | Mean |
|---|---|---|---|---|---|---|---|---|---|---|
| | X, Wang [20] | 0.01 | 0.03 | 0.02 | 0.00 | 0.00 | 0.00 | 0.01 | 0.02 | 0.01 |
| | Z, Li [8] | 0.04 | 0.52 | 0.07 | 0.09 | 0.11 | 0.01 | 0.05 | 0.05 | 0.12 |
| | J, Liu [9] | 0.18 | 0.70 | 0.28 | 0.41 | 0.27 | 0.04 | 0.25 | 0.18 | 0.29 |
| | Baseline | **0.34** | **1.00** | 0.40 | **0.68** | 0.11 | **0.14** | **0.65** | 0.00 | 0.41 |
| 0.5_0.8 | IK | 0.22 | 0.82 | 0.36 | 0.56 | 0.32 | **0.14** | 0.25 | **0.35** | 0.38 |
| | IR | 0.24 | 0.82 | 0.40 | 0.56 | 0.32 | 0.07 | 0.38 | 0.30 | 0.39 |
| | KR | 0.24 | 0.89 | 0.32 | **0.68** | 0.26 | **0.14** | 0.21 | 0.30 | 0.38 |
| | (IR+IK) | 0.20 | 0.86 | **0.48** | **0.68** | 0.32 | **0.14** | 0.54 | 0.30 | **0.44** |
| | IR+IK+KR | 0.27 | 0.86 | 0.40 | 0.56 | **0.37** | **0.14** | 0.13 | 0.30 | 0.38 |
| | J, Liu [9] | 0.08 | 0.30 | **0.09** | 0.25 | 0.19 | 0.01 | 0.04 | 0.07 | 0.13 |
| | Baseline | 0.11 | 0.34 | 0.06 | 0.32 | 0.20 | 0.01 | 0.00 | 0.06 | 0.14 |
| | IK | 0.10 | 0.59 | 0.07 | 0.37 | 0.20 | 0.00 | **0.13** | 0.06 | **0.19** |
| 1.0_0.0 | GR | 0.06 | 0.61 | 0.07 | 0.28 | 0.14 | 0.00 | 0.05 | 0.08 | 0.16 |
| | IK | 0.09 | 0.63 | 0.06 | 0.36 | **0.22** | 0.00 | 0.09 | 0.07 | **0.19** |
| | IR+IK | 0.06 | **0.64** | 0.08 | **0.38** | 0.19 | 0.01 | 0.08 | **0.09** | **0.19** |
| | IR+IK+KR | **0.12** | 0.51 | 0.07 | 0.36 | **0.22** | **0.03** | 0.02 | 0.07 | 0.17 |
| | J, Liu [9] | **0.19** | 0.47 | 0.20 | 0.41 | **0.22** | **0.06** | 0.12 | 0.11 | 0.22 |
| | Baseline | 0.14 | 0.62 | 0.20 | 0.42 | 0.07 | 0.00 | 0.23 | 0.08 | 0.22 |
| | IK | 0.14 | 0.66 | 0.09 | 0.47 | 0.15 | 0.00 | **0.30** | 0.06 | 0.23 |
| 1.0_0.4 | GR | 0.14 | **0.75** | **0.24** | 0.42 | 0.11 | 0.00 | 0.26 | **0.12** | 0.25 |
| | KR | 0.13 | 0.68 | 0.20 | 0.47 | 0.19 | **0.06** | 0.17 | 0.08 | 0.25 |
| | IR+IK | 0.13 | 0.72 | 0.13 | 0.43 | 0.20 | 0.00 | 0.23 | 0.06 | 0.24 |
| | IR+IK+KR | 0.18 | 0.71 | 0.20 | **0.50** | 0.20 | 0.02 | 0.29 | 0.06 | **0.27** |

Table 4. The comparison results of disease localization among the models using three sets of data at T(IoU)=0.7, including 50% unannotated and 80% annotated images (0.5_0.8), 100% unannotated and no any annotated images (1.0_0.0), and 100% unannotated and 40% unannotated images (1.0_0.4). For each disease, the best results are bolded.

the inter-image relation module and the intra-image knowledge learning (IR+IK), the model combining the inter-image relation module, the intra-image knowledge learning and the knowledge reasoning module (IR+IK+KR).

Table 4 shows the results of the three experiments mentioned in section 4.2 at T(IOU)=0.7. It can be seen that our method performs better in most classes except for "Atelectasis", "Effusion" and "Mass" comparing [20], [8] and [9]. Furthermore, comparing the baseline model, it can be observed that the performance of our other models are improved in most cases, which shows that our method is effective for improving model performance. However, a model does not always maintain the advantage in the three experiments, for example, the model (IR+IK) achieves the best performance in the data (0.5_0.8), the model (IK), the model (KR) and the model (IR+IK) achieve the best performance in the data (1.0_0.0), and the model (IR+IK+KR) achieves the best performance in the data (1.0_0.4). Overall, the experimental results demonstrate that using structural relational information can improve the performance of models. For different experimental data, our models can achieve different results. It is difficult for us to determine which model is the best, but we can be sure that our method is effective,

because no matter what kind of data we use, our models achieve great improvement. Particularly, the method can achieve good localization results even without any annotation images for training.

## 5. Conclusion

By imitating doctor's training and decision-making process, we propose the Cross-chest Graph (CCG) to improve the performance of automatic lesion detection under limited supervision. CCG models the intra-image relationship between different anatomical areas by leveraging the structural information to simulate the doctor's habit of observing different areas. Meanwhile, the relationship between any pair of images is modeled by a knowledge-reasoning module to simulate the doctor's habit of comparing multiple images. We integrate intra-image and inter-image information into a unified end-to-end framework. Experimental results on the NIH Chest-14 dataset demonstrate that the proposed method achieves state-of-the-art performance in diverse situations.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels. Technical report, 2010. 4

[2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5

[3] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 1

[4] Qingji Guan, Yaping Huang, Zhun Zhong, Zhedong Zheng, Liang Zheng, and Yi Yang. Diagnose like a radiologist: Attention guided convolutional neural network for thorax disease classification. *arXiv preprint arXiv:1801.09927*, 2018. 2

[5] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid poolling in deep convolutional networks for visual recognition. *In European Conference of Computer Vision (ECCV)*, 2014. 1

[6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 5

[7] Christy Y Li, Xiaodan Liang, Zhiting Hu, and Eric P Xing. Knowledge-driven encode, retrieve, paraphrase for medical image report generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6666–6673, 2019. 3

[8] Zhe Li, Chong Wang, Mei Han, Yuan Xue, Wei Wei, Li-Jia Li, and Li Fei-Fei. Thoracic disease identification and localization with limited supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8290–8299, 2018. 2, 3, 4, 5, 6, 8

[9] J. Liu, G. Zhao, Y. Fei, M. Zhang, Y. Wang, and Y. Yu. Align, attend and locate: Chest x-ray diagnosis via contrast induced attention network with limited supervision. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10631–10640, 2019. 1, 2, 3, 4, 5, 6, 7, 8

[10] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 1, 2

[11] Yun Liu, Yu-Huan Wu, Yunfeng Ban, Huifang Wang, and Ming-Ming Cheng. Rethinking computer-aided tuberculosis diagnosis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2646–2655, 2020. 3

[12] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1

[13] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5

[14] Emanuele Pesce, Petros-Pavlos Ypsilantis, Samuel Withey, Robert Bakewell, Vicky Goh, and Giovanni Montana. Learning to detect chest radiographs containing lung nodules using visual attention networks. *arXiv preprint arXiv:1712.00996*, 2017. 2

[15] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 1, 2

[16] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015. 1, 2

[17] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 1

[18] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 1

[19] R. Venkatesan, S. M. Koon, M. H. Jakubowski, and P. Moulin. Robust image hashing. In *International Conference on Image Processing*, 2000. 4

[20] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. 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*, pages 2097–2106, 2017. 2, 5, 6, 8

[21] Petros-Pavlos Ypsilantis and Giovanni Montana. Learning what to look in chest x-rays with a recurrent visual attention model. *arXiv preprint arXiv:1701.06452*, 2017. 2

[22] Gangming Zhao, Chaowei Fang, Guanbin Li, Licheng Jiao, and Yizhou Yu. Contralaterally enhanced networks for thoracic disease detection, 2020. 1, 2

[23] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 2