

# Robust Classification from Noisy Labels: Integrating Additional Knowledge for Chest Radiography Abnormality Assessment

Sebastian Gundel<sup>a,c,\*</sup>, Arnaud A. A. Setio<sup>a</sup>, Florin C. Ghesu<sup>b</sup>, Sasa Grbic<sup>b</sup>, Bogdan Georgescu<sup>b</sup>, Andreas Maier<sup>c</sup>, Dorin Comaniciu<sup>b</sup>

<sup>a</sup>Digital Technology and Innovation, Siemens Healthineers, 91052 Erlangen, Germany

<sup>b</sup>Digital Technology and Innovation, Siemens Healthineers, Princeton, NJ 08540, USA

<sup>c</sup>Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg, 91058 Erlangen, Germany

---

## Abstract

Chest radiography is the most common radiographic examination performed in daily clinical practice for the detection of various heart and lung abnormalities. The large amount of data to be read and reported, with more than 100 studies per day for a single radiologist, poses a challenge in consistently maintaining high interpretation accuracy. The introduction of large-scale public datasets has led to a series of novel systems for automated abnormality classification. However, the labels of these datasets were obtained using natural language processed medical reports, yielding a large degree of label noise that can impact the performance. In this study, we propose novel training strategies that handle label noise from such suboptimal data. Prior label probabilities were measured on a subset of training data re-read by 4 board-certified radiologists and were used during training to increase the robustness of the training model to the label noise. Furthermore, we exploit the high comorbidity of abnormalities observed in chest radiography and incorporate this information to further reduce the impact of label noise. Additionally, anatomical knowledge is incorporated by training the system to predict lung and heart segmentation, as well as spatial knowledge labels. To deal with multiple datasets and images derived from various scanners that apply different post-processing techniques, we introduce a novel image normalization strategy. Experiments were performed on an extensive collection of 297,541 chest radiographs from 86,876 patients, leading to a state-of-the-art performance level for 17 abnormalities from 2 datasets. With an average AUC score of 0.880 across all abnormalities, our proposed training strategies can be used to significantly improve performance scores.

---

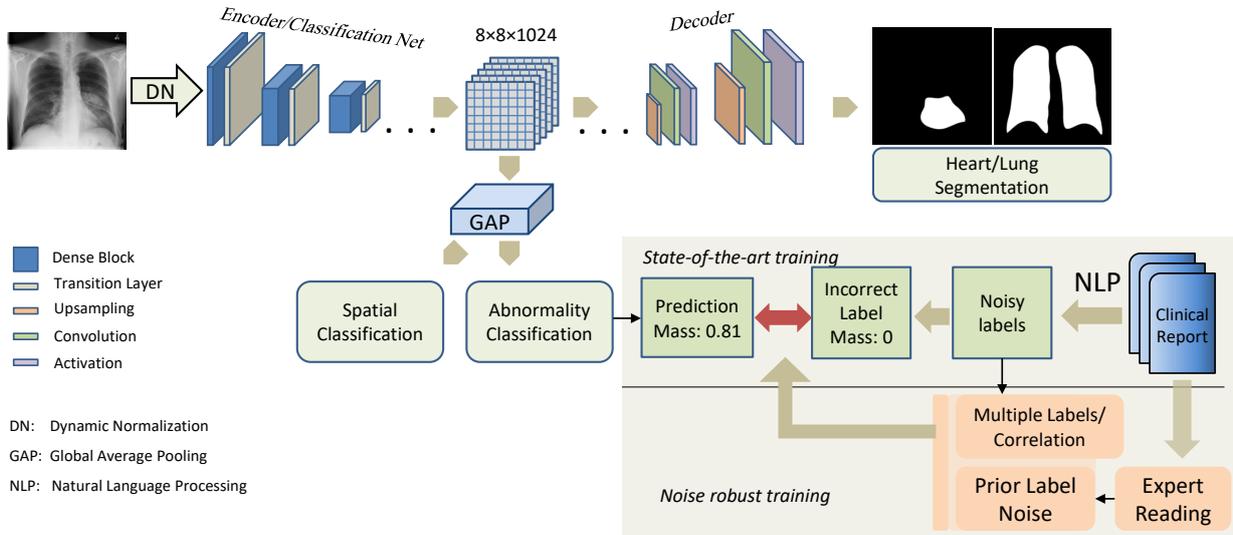
## 1. Introduction

Recent developments in the deep learning community combined with the availability of large labeled datasets have enabled the training of automated systems that can exceed human performance on a variety of classification, detection and segmentation tasks (Li et al., 2016; Ghesu et al., 2019a; Zhu et al., 2018). In different scenarios, such systems actively support humans, increasing the efficiency and accuracy of their workflow. In the medical domain, deep learning systems for image and data analysis and integration can potentially have an even greater impact, supporting the clinical workflow from patient admission to diagnosis, treatment and

follow-up investigations (Rajkomar et al., 2018; Ardila et al., 2019). In this paper, we focus on the problem of diagnosing multiple abnormalities based on chest radiography of the human body. In practice, this is a challenging problem reflected in a significant variability between different radiologists (Bruno et al., 2015). For example, a study on determining a diagnosis of tuberculosis between 25 radiologists on 50 chest radiographs leads to a moderate agreement with a kappa score of 0.448 (Balabanova et al., 2005). The most important cause for high reader variability is the lack of an established guideline of reading (like the Fleischner guideline for CT). As such the interpretation on radiographs is often subjective leading to a large degree of variability. Other factors include the complex appearance of pathologies in radiographs and the large number of scans that need to be read and analyzed daily under time pressure (Brady, 2016). The average time to read and

---

\*Corresponding author at: Digital Technology and Innovation, Siemens Healthineers, 91052 Erlangen, Germany and Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg, 91058 Erlangen, Germany  
E-mail address: sebastian.guendel@fau.de



**Fig. 1.** Additional Knowledge Integration: We integrate several additional features to handle label noise during training and to increase the abnormality classification performance. Our objective is adapted with prior label noise and label correlation information. In addition to the classification of the abnormalities, spatial classification and lung/heart segmentation are trained.

report a plain film is as little as 1.4 minutes (Fleishon et al., 2006).

Aiming to reduce the user variability, algorithms for automatic classification of abnormalities visible in chest radiographs have been proposed. The recent introduction of large-scale public datasets (Wang et al., 2017; Gohagan et al., 2000) has led to series of automated solutions using deep learning for abnormality classification. However, training on these datasets can be sub-optimal, as labels were obtained using natural language processed medical reports, which may introduce incorrect labels. As deep learning algorithms tend to overfit to provided (incorrect) labels, networks may be trained into the wrong direction, limiting the performance of the network on unseen data. Moreover, evaluating the performance on a test set with NLP-based labels can be misleading and may not be adequate in representing the performance needed for diagnostic interpretations in practice.

Furthermore, clinical reading reports typically include further relevant information in addition to a found abnormality, e.g., its approximate position. Certain abnormalities, e.g., nodules, appear as small objects in the radiograph. Thus, prior spatial knowledge about these abnormalities in shape of classification labels (e.g. nodule in the lower-left lung, yes or no?) may regularize the learning process when inductively transferring the information on a shared representation. The resulting learning process, the so-called multi-task learning procedure

has been demonstrated in several clinical applications in recent years, e.g. in Amyar et al. (2020), which leads to improved results on the classification of abnormalities.

In this study, we propose novel training strategies to handle label noise. First, we evaluate the performance of the NLP labels by performing an observer study, where a subset of training data was re-read by 4 board-certified radiologists.

Thereafter, a novel objective function that takes into account label noise and label correlation is proposed. Label noise was measured from our observer study and label correlation is incorporated to exploit possible comorbidity between abnormalities.

To regularize the network further, additional tasks are included. Anatomical knowledge is incorporated by adding lung and heart segmentation into the system. Moreover, we also ask the network to not only predict the abnormality, but also the location of the abnormality within the lung. We refer this task as spatial classification. To deal with multiple datasets and images derived from various scanners applying different post-processing techniques, a novel image normalization strategy is introduced.

The performance of the proposed algorithms was evaluated on a dataset with radiologists-based reference standard. For comparison purposes, performance scores on datasets with NLP-based reference standard is provided. The global architecture and the proposed methods can be seen in Figure 1.

### The contributions of this paper are as follows:

- We evaluate and measure label noise by comparing the original NLP-derived labels with the consensus labels from expert readers.
- We formulate a novel objective function that takes into account label noise and label correlation
- We propose a novel multi-task deep neural network for multi-abnormality classification, spatial classification and lung/heart segmentation based on frontal chest radiography images.
- We propose a novel image normalization strategy that robustly handles brightness and contrast variability of images from multiple datasets and scanners.
- We demonstrate that by using the proposed training strategies, additional anatomical knowledge, and novel normalization technique one can significantly increase the accuracy of the abnormality classification.

## 2. Related Work

### 2.1. Multi-Abnormality Classification

The publication of the ChestX-ray14 (NIH) dataset (Wang et al., 2017) has led to a series of recent publications that propose automatic systems for abnormality classification. At first, Wang et al. (2017) evaluated several state-of-the-art convolutional neural network architectures, reporting an area under the ROC curve (AUC) of 0.75 on average. Islam et al. (2017) defined an ensemble of multiple state-of-the-art network architectures to increase the classification performance. Subsequently, Rajpurkar et al. (2017) demonstrated that a common DenseNet architecture (Huang et al., 2017) can surpass the accuracy of radiologists in detecting pneumonia. In addition to a DenseNet, Yao et al. (2017) implemented a Long-short Term Memory (LSTM) model to exploit dependencies between the abnormalities. An attention guided convolutional neural network architecture was used by Guan et al. (2018) to specifically focus on the region of interest which is provided in a second network branch as a cropped image with higher resolution. While Rubin et al. (2018) designed a Dual-Network to extract the image information of both frontal and lateral views, Yan et al. (2018) used a DenseNet architecture integrated with "Squeeze-and-Excitation" blocks (Hu et al., 2018) to improve the performance.

Moreover, Cai et al. (2018) used a multi-scale aggregation and, additionally, they implemented an attention mining strategy to find accurately diseased regions. Tang et al. (2018) integrated curriculum learning to specifically train based on severity-level attributes of the radiology reports. Attention guidance with heatmaps supported the learning process to increase the performance. Shen and Gao (2018) modified the DenseNet architecture to introduce dynamic routing between capsules. Furthermore, Guan and Huang (2018) developed a category-wise residual attention framework with a feature embedding and an attention learning module. Wang et al. (2020) developed a new network architecture with both a classification and an attention branch, where the latter calculates activation maps with gradient-weighted class activation mapping (Ba et al., 2014) which is subsequently concatenated with the classification branch. The same group developed a system for channel-wise, element-wise, and scale-wise attention learning to classify 14 abnormalities (Wang et al., 2021). Based on very limited location labels of the abnormalities, Li et al. (2018) trained a neural network to predict both classification and localization of the abnormalities. Liu et al. (2019) designed a network architecture with two branches, similar to Guan et al. (2018), where the second branch used a cropped image input based on existing lung masks. Yao et al. (2018) defined a network architecture which can be trained on different resolutions.

In contrast, Rajpurkar et al. (2018) used multiple radiologists to re-read the images: One subgroup of radiologists defined the ground truth of a set where the other subgroup and the neural network was evaluated on. In this way, performance of both the radiologists and the deep learning algorithm could be compared. Irvin et al. (2019) trained on a dataset where the ground truth consists of an additional uncertainty class. Different approaches were applied during training to increase the performance with the uncertainty information. In Huang et al. (2020), an ensemble of several state-of-the-art architectures is introduced. Additionally, the network is composed of more sub-networks with different resolutions. Finally, Ghesu et al. (2019b) and Ghesu et al. (2021) proposed a model that can implicitly learn to estimate predictive uncertainty as an orthogonal measure to the predicted abnormality probability using principles of subjective logic (Jøssang, 2016).

In an attempt to reduce label uncertainty and increase accuracy, the same group proposed an innovative framework for segmentation of airspace opacities from chest radiographs, trained on digitally reconstructed radiographs with precise annotations derived from computed tomography (CT) scans (Barbosa Jr et al., 2020)

In general, we emphasize that most of the published work report classification results by splitting the data completely randomly for training, validation and testing (Yao et al., 2017; Guan et al., 2018; Li et al., 2018; Islam et al., 2017). With this splitting strategy, images from the same patient may be located in both training and testing set. For a fair performance evaluation, the splitting should always be performed at patient level. Moreover, to ensure a fair comparison between different methods, the same data split needs to be applied.

### 2.2. Lung and Heart Segmentation

The segmentation of lungs and heart based on chest radiographs has been of high interest in the last years. Gómez et al. (2020) proposed a method to segment the lungs and the heart including a novel network architecture based on a U-net as backbone. Another method detects the lung boundaries with training features of the gray level co-occurrence matrix (Zotin et al., 2019). Gaál et al. (2020) developed an architecture for lung segmentation using an adversarial critic network, in addition to the segmentation network. In the work of Kholiavchenko et al. (2020), different state-of-the-art segmentation networks are used as backbone. In this method, the contours of the lungs were determined and used for training. Selvan et al. (2020) specifically focused on chest radiographs with high intensity change because of the abnormalities resulting in extensively dense regions. Larrazabal et al. (2020) demonstrated that an additional variational autoencoder fulfilled anatomical plausibility which improved the segmentation of lungs and heart.

### 2.3. Spatial Classification

For spatial classification there is currently no research available as the chest radiography datasets typically include either no additional spatial information or precise detection information which can be converted into segmentation maps.

### 2.4. Noisy Label Classification

While many publications highlight the high fraction of label noise in the ChestX-ray14 dataset, few solutions have been proposed to address this limitation. Rolnick et al. (2017) showed that, to some extent, deep learning models can be robust against label noise during training, in general. Therefore, many groups kept standard training procedures and focused solely on the re-definition of labels in the test set (Rajpurkar

et al., 2018). An overview of possibilities to deal with label noise on medical image data was demonstrated by Dou et al. (2020). The methods for noise robustness were defined in different categories, e.g., network architecture and training procedures. In terms of loss adaption, the categorical cross-entropy was changed by Rusiecki (2019), adding what they called the least trimmed absolute value estimator which makes the training more robust to outliers. Most publications in the field of label noise robustness, however, show their approaches based on standard datasets (Molchanov et al., 2017; Lakshminarayanan et al., 2017; Gal and Ghahramani, 2016), e.g., MNIST and CIFAR, remaining unclear how well these methods can generalize on medical datasets.

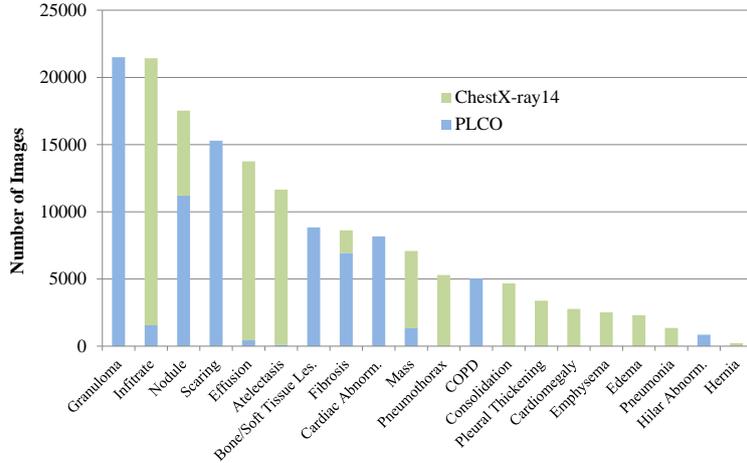
## 3. Problem Definition and Dataset Analysis

In this Section, the datasets that are used to train and validate our method are presented. The observer study aimed to evaluate the label noise on these datasets is described. Thereafter, correlation between labeled abnormalities is evaluated.

### 3.1. Dataset

Our data collection is composed of two different datasets, the ChestX-ray14 (Wang et al., 2017) and PLCO (Gohagan et al., 2000) dataset. The ChestX-ray14 database contains 67,310 posterior-anterior (PA) and 44,810 anterior-posterior (AP) images whereas the PLCO dataset is based on a screening study and solely contains PA images. Table 1 gives an overview of the datasets. By combining both datasets, 297,541 frontal chest radiographs from 86,876 patients can be used. Including follow-up scans, there is an average of 3-4 images per patient. The PLCO dataset includes spatial information, i.e. approximate knowledge about where abnormalities are located. Figure 2 shows the number of images that contain each abnormality. One image can also show multiple abnormalities. Additionally, the collection contains 178,319 images where none of the mentioned abnormalities appear, these images are not counted in Figure 2. Note that the high imbalance of the data collection with respect to different abnormalities represents a challenge in ensuring training stability and performance.

**Assessing the Quality of Labels:** As the ChestX-Ray14 dataset contains 112,120 images, re-reading entire dataset is not feasible. Therefore, we selected a subset of samples using random sampling. We limit the



**Fig. 2.** The number of images associated with abnormalities in ChestX-ray14 and PLCO datasets. The number of images where none of these pathologies appear is excluded.

**Table 1**

Overview of the 2 considered datasets. The combined data collection consists of 297,541 images from 86,876 patients. The last column describes the subset which is read by our expert radiologists. Note that the image number (first row) denotes the entire number of images in the dataset, however, patients used for the expert read are excluded in the ChestX-ray14 dataset (first column) for training and validation. The last two rows show the split statistics for training (train), validation (valid), calibration (calib), and testing (test).

Name	ChestX-ray14		PLCO			Consensus Expert Read (ChestX-ray14)	
Image / Patient number (#)	112,120 / 30,805		185,421 / 56,071			689 / 689	
Abnormality number	14		12			5	
Image size	1024 × 1024		~ 2500 × 2100			1024 × 1024	
Split	train	valid	train	valid	test	calib	test
Image number (%)	90	10	70	10	20	30	70
Image number (#)	99,462	8,338	129,658	18,773	36,990	207	482

process to 5 abnormalities (effusion, cardiomegaly, consolidation, atelectasis, mass) which were chosen based on clinical importance. We take only one image per patient. The expert reading subset contains 689 images. The remaining images of each patient are removed from the remaining data to avoid a patient overlap. The reading process was based on consensus decision-making: First, 4 board-certified radiologists blindly re-labeled the images. In a second stage, for all cases where consensus was not reached on all labels through the independent read, an open discussion was carried out to establish consensus labels. The original dataset labels were not provided during the re-reading process to avoid a biased decision towards the original labels.

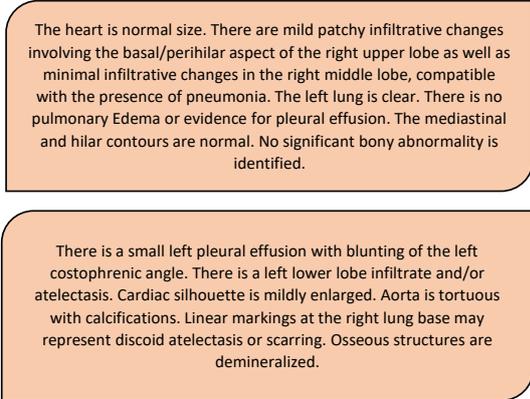
We hypothesize that the re-reading process following

of consensus of experts is imperative to understand and address the limited quality of the labels of the ChestX-Ray14 dataset. Thus, we assume that

$$Y^{True} \approx Y^{Rad}, \quad (1)$$

where  $Y^{True}$  are the true (unknown) labels, and  $Y^{Rad}$  the generated labels by consensus of expert radiologists. An overview of the subset can be seen in the last column of Table 1.

**Dataset Splits:** As Table 1 shows, 2 datasets are used and a subset of ChestX-ray14 images are re-read to obtain higher quality reference labels. The PLCO dataset is split into 70% for training, 10% for validation, and 20% for testing. The ChestX-ray14 dataset (excluding



**Fig. 3.** Two randomly chosen, abnormal clinical reports based on chest radiographs derived from the same reader.

the re-read samples) is split into 90% training and 10% validation. A subset of the ChestX-ray14 set is re-read (see last column in Table 1) where 70% is used for the evaluation. The last 30% serve as calibration of our proposed methods. *Patient-wise* splits are considered for all experiments to separate the patients into training, validation, and test set.

### 3.2. Errors in the Abnormality Labeling Process

The misinterpretation of medical images can be caused by different reasons, e.g, radiologist fatigue or lack of attention (Brady, 2016). Recently, many publications highlight the significant label error in the ChestX-ray14 dataset (Oakden-Rayner, 2017). Beside the clinical reader, another major error domain is derived from the natural language processing (NLP) algorithm that transfers the clinical reports into binary abnormality labels.

Based on two reports from the same reader in Figure 3, the complexity and variability of the content can be observed. As medical reports for chest radiography do not follow well-established guidelines, the structure of the content is typically built up in arbitrary order. Furthermore, we see ambiguous terms, e.g., “*infiltrate and/or atelectasis*” or “*may represent*”. We hypothesize that this ambiguity combined with the unstructured complexity and additional user variance may lead to errors in the natural language processing algorithm and, thus, the noise in the dataset labels (Jusoh, 2018).

This error leads to an incorrect reference standard on which the networks are trained and evaluated. Several publications show that this label error negatively affects

the performance, e.g., in classifying abnormalities correctly (Rajpurkar et al., 2018). However, the degree of noise that is present in these datasets remains unclear. In order to measure the degree of label noise in our datasets, we perform an observer study on a subset of data and quantitatively assess the label noise probabilities.

### 3.3. Exploiting Abnormality Comorbidity: A Correlation Perspective

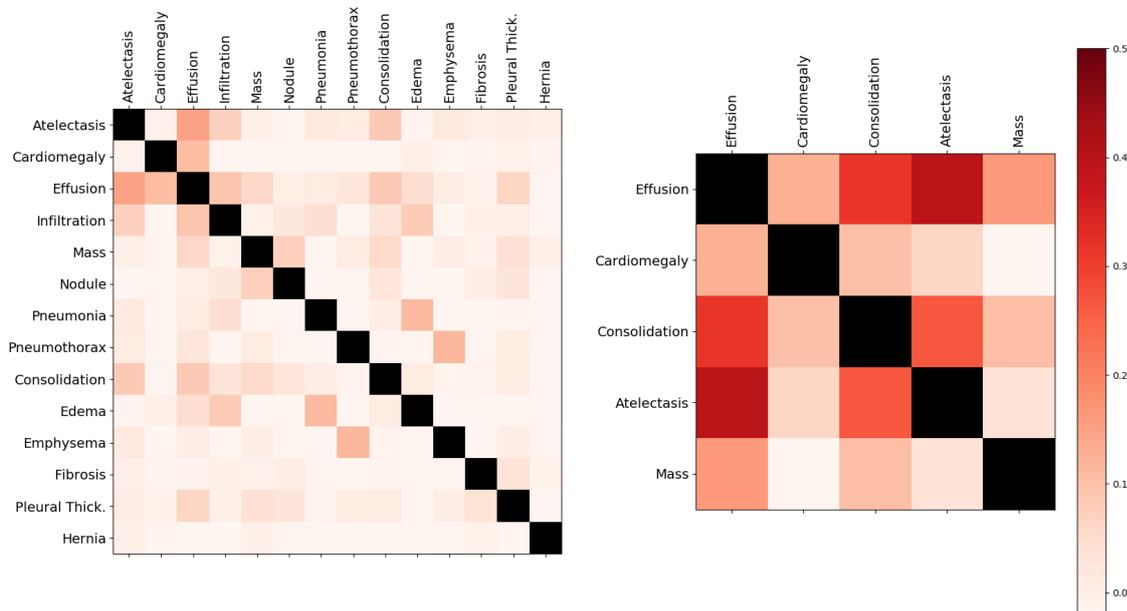
The comorbidity of abnormalities have been explored in many fields, e.g., in Guan et al. (2020). Comorbidity indicates simultaneous presence of multiple abnormalities. The knowledge of comorbidity may assist clinicians in making a more precise assessment of diseases or abnormalities (Bayliss et al., 2005). In this section we present a principled strategy of how to exploit the comorbidity of different abnormalities to improve the robustness of the model.

In the context of multi-label classification, the correlation of different labels is an important piece of information (Ghamrawi and McCallum, 2005; Zhu et al., 2005). Conceptually, we analyse how strong a set of class labels  $c^{(n)}$  for abnormality  $n$  correlate with a set of class labels  $c^{(r)}$  for abnormality  $r$  where  $r \in \{1 \dots D\} \setminus \{n\}$  and  $D$  denotes the number of abnormalities. We use the Pearson correlation coefficient which is sensitive to class imbalance to measure the correlation on the ChestX-Ray14 training set between the abnormalities:

$$\text{corr}_{\text{Pearson}}(c^{(n)}, c^{(r)}) = \frac{\text{cov}(c^{(n)}, c^{(r)})}{\sigma^{(n)}\sigma^{(r)}}, \quad (2)$$

where  $\text{cov}$  denotes the covariance and  $\sigma$  the standard deviation. In Figure 4 (left), we color the strength of correlation between the abnormalities of the original dataset labels. The correlation of the dataset with labels based on a consensus of radiologists can be seen in Figure 4 (right) which shows stronger correlation between the abnormalities, e.g., the Pearson correlation coefficient between atelectasis and effusion equals 0.395.

There are different reasons why the dataset labels have significant correlation: Based on the paper by Irvin et al. (2019), abnormalities can be represented hierarchically. Thus, an abnormality can consist of a subgroup of abnormalities, e.g., lung opacity includes pneumonia or edema. This stratification into different groups is also addressed in Oakden-Rayner (2019). Further existing publications analyse the relationships between abnormalities, e.g., lung inflammation and cardiovascular disease (Van Eeden et al., 2012).



**Fig. 4.** Pearson correlation of abnormalities; Left: Correlation of the original ChestX-Ray14 dataset labels; Right: Correlation of redefined labels based on a consensus of radiologists.

Overall, we hypothesize that the given correlation between abnormalities may be related to some degree of existing comorbidity. In practice, this correlation information can be exploited to serve an additional learning support during training.

#### 4. Methodology

Given an arbitrary AP or PA chest radiograph  $\mathbf{I}$  with size  $N \times N$  pixels, we design a deep learning based system parametrized by  $\theta$  which outputs the probability of different abnormalities being present in the image:  $\vec{\delta} = p(\mathbf{I}; \theta)$ , where  $\vec{\delta} \in [0, 1]^D$  and  $D$  is the number of considered abnormalities.

We exploit the correlation between abnormalities and regularize our loss function using information from both label noise and label correlation. The system is also designed to predict approximate spatial location of the abnormalities and a probabilistic segmentation map  $\mathbf{S} \in [0, 1]^{2 \times N \times N}$  for both lungs and the heart. The proposed normalization strategy wraps up our novel multi-task system focusing on the improvement of abnormality classification. Figure 1 shows the global architecture including our proposed methods.

##### 4.1. Deep Neural Network Design

The backbone architecture of our model is inspired from the DenseNet architecture (Huang et al., 2017). We adopt this network architecture with 5 dense blocks and a total of 121 convolutional layers. Each dense block consist of several dense layers which include batch normalization, rectified linear units, and convolution. The novelty of the DenseNet are the skip connections, meaning that within a block, each layer is connected to all subsequent layers. Between each dense block, a so-called transition layer is added, which includes batch normalization, convolution, and pooling, to reduce the dimensions.

The single grayscale input image  $\mathbf{I}$  is rescaled to  $N \times N$  pixels (in our experiments  $N=256$ ) using bilinear interpolation and fed into the network. The network is initialized with the pre-trained ImageNet model (Russakovsky et al., 2015). To maintain all pre-trained weights, we replicate the image to 3 channels. The global average pooling (GAP) layer is applied to obtain the representation to be used for classification. The number of output units is set to the number of abnormality classes  $D$ . We use sigmoid activation functions for each class to map the output to a probability interval  $[0, 1]$ .

**Table 2**

Hyperparameters used to train our model.

Hyperparameter	Value
Initial Learning Rate	$10^{-3}$
Image Input Size	256
Batch Size	128
Optimizer	Adam
Loss Function	Binary Cross Entropy
Number of Epochs	30 (early stopping)

#### 4.2. Model Training

The algorithm for multi-label classification problem is trained using the following approach: The training process is modified such that each class can be trained individually. We create  $D$  binary cross-entropy loss functions. The corresponding labels  $[c^{(1)}, c^{(2)} \dots c^{(D)}] \in \{0, 1\}$  (absence or presence of the abnormality, respectively) are compared with the network output  $[p^{(1)}, p^{(2)} \dots p^{(D)}] \in [0, 1]$  and the loss is measured. Due to the highly imbalanced problem we introduce additional weight constants  $w_P^{(n)}$  and  $w_N^{(n)}$  for each abnormality indexed by  $n$ , to the cross-entropy function:

$$\mathcal{L}_{Abn} = - \sum_{n=1}^D \sum_{i=1}^F \left[ w_P^{(n)} c_i^{(n)} \ln(p_i^{(n)}) + w_N^{(n)} (1 - c_i^{(n)}) \ln(1 - p_i^{(n)}) \right], \quad (3)$$

where  $w_P^{(n)} = \frac{P^{(n)} + N^{(n)}}{P^{(n)}}$  and  $w_N^{(n)} = \frac{P^{(n)} + N^{(n)}}{N^{(n)}}$ , with  $P^{(n)}$  and  $N^{(n)}$  indicating the number of positive and negative cases for the entire training dataset, respectively. The same weighting strategy is also applied in Wang et al. (2017) and Rajpurkar et al. (2018). Equation 3 calculates the loss sum over all images indexed by  $i$ , where  $F$  denotes the total number of images in the set. For all experiments, we train with 128 samples in each batch. The Adam optimizer (Kingma and Ba, 2015) ( $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ ) is used with an adaptive learning rate: the learning rate is initialized with  $10^{-3}$  and reduced by a factor of 10 when the validation loss plateaus. All parameters can be found in Table 2.

An important issue encountered with abnormality labeling is the varying and overlapping definition and interpretation between radiologists as seen in Section 3.2. Therefore, we treat the corresponding abnormalities of both datasets separately for the experiments in

**Table 3**

Sensitivity/Specificity scores of original dataset labels versus re-read labels.

Abnormality	$s_{sens}$	$s_{spec}$
Effusion	0.300	0.966
Cardiomegaly	0.342	0.986
Consolidation	0.129	0.949
Atelectasis	0.221	0.970
Mass	0.364	0.972
Average	0.271	0.969

Section 5 and create different output classes. Given  $D_1 = 14$  abnormalities of the ChestX-ray14 dataset and  $D_2 = 12$  abnormalities of the PLCO dataset, we define  $D = D_1 + D_2 = 26$  classes for our network. Furthermore, we only compute gradients for labels of one dataset where the current image is derived from. This strategy avoids a class categorization step beforehand and ensures that each network layer (except the last) receives information of all images.

#### 4.3. Prior Label Noise Probability

One can compute sensitivity and specificity of the original dataset labels based on the subset of re-read labels. Assuming the corrected labels are the true labels as discussed in Section 3.1, we compute the performance scores which reflect the noise probabilities of positives and negatives on the original ChestX-Ray14 labels. Table 3 shows sensitivity  $s_{sens} = \frac{TP}{P}$  and specificity  $s_{spec} = \frac{TN}{N}$  scores of the selected abnormalities which are used to regularize the loss function. TP and TN respectively denote the number of original positive and negative labels which are correctly labeled based on the re-read subset. Parameter P and N stand for the total number of positive and negative cases in the re-read subset. Thus, low scores indicate strong label noise. We adapt our original loss function with a second term which is based on the inverse binary cross entropy function:

$$\begin{aligned} \mathcal{L}_{Noise} &= \mathcal{L}_{Abn} + r_{noise} \\ &= - \sum_{n=1}^D \sum_{i=1}^F \left[ w_P^{(n)} c_i^{(n)} \ln p_i^{(n)} + w_N^{(n)} (1 - c_i^{(n)}) \ln (1 - p_i^{(n)}) + \lambda_{Noise} \left[ f_P^{(n)} w_N^{(n)} (1 - c_i^{(n)}) \ln p_i^{(n)} + f_N^{(n)} w_P^{(n)} c_i^{(n)} \ln (1 - p_i^{(n)}) \right] \right], \quad (4) \end{aligned}$$

where  $f_P$  and  $f_N$  are the individual regularization weights for positive and negative examples. In our experiments, we set  $f_P^{(n)} = 1 - s_{sens}^n$  and  $f_N^{(n)} = 1 - s_{spec}^n$ . The parameters are derived based on the calibration set which is introduced in Section 3.1. Parameter  $\lambda_{Noise}$  is another weight to define the overall influence of the regularization term. Strong weights mean that a higher error is computed when predicting the original label. In our experiments, we set  $\lambda_{Noise} = 0.1$  (value with best performance on validation set). The regularization term is disregarded for abnormalities excluded in the calibration set.

#### 4.4. Prior Label Correlation

In Figure 4, we see that the labels of many classes are correlated. This correlation information can also be used as another prior in the loss function. We adapt the original loss function with a term to consider the information across all abnormality labels:

$$\begin{aligned} \mathcal{L}_{Corr} &= \mathcal{L}_{Abn} + r_{corr} \\ &= - \sum_{i=1}^F \sum_{n=1}^D \left[ w_P^{(n)} c_i^{(n)} \ln p_i^{(n)} + w_N^{(n)} (1 - c_i^{(n)}) \ln (1 - p_i^{(n)}) + \right. \\ &\quad \left. \sum_{r \in \{1 \dots D\} \setminus \{n\}} \left[ cov^{(n,r)} \left[ w_P^{(r)} c_i^{(r)} \ln p_i^{(r)} + \right. \right. \right. \\ &\quad \left. \left. \left. w_N^{(r)} (1 - c_i^{(r)}) \ln (1 - p_i^{(r)}) \right] \right] \right], \end{aligned} \quad (5)$$

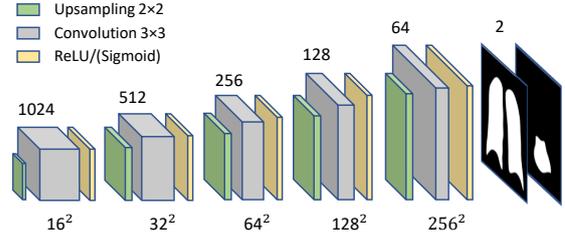
where  $cov^{(n,r)}$  with element  $(n, r)$  measures the covariance between label indexed as  $n$  and the label indexed as  $r$ . Thus, label correlation between 2 abnormalities influence on the overall loss function. Depending on the covariance matrix, all abnormality labels may influence on a specific abnormality. The covariance matrix is calculated based on the calibration subset. The regularization term is disregarded for abnormalities excluded in the calibration set.

#### 4.5. Integrating Additional Knowledge

Additional anatomical knowledge related to individual pathologies as well as the underlying anatomy, i.e., the heart and the lungs, can be exploited to increase the classification performance.

##### 4.5.1. Lung and Heart Segmentation

First, one can focus the learning task to the heart and lung regions. The image information outside of these regions may be regarded as irrelevant for the diagnosis of



**Fig. 5.** Architecture of the decoder to predict segmentation masks. The top number shows the channel number, the bottom value indicates the feature map size. The network is connected to the classification network after the last dense block (left). The final sigmoid layer predicts the lung and heart masks in 2 channels (right).

these lung/heart abnormalities. Lung and heart segmentation masks are available for the entire data collection.

Instead of providing the masks as input for the classification network, we extend the classification network with a decoder branch and predict the masks (see Figure 1). In this way, the additional knowledge about the shape of the heart and lungs is integrated in an implicit way, i.e., during learning through the flow of gradients. As such, in the encoder part, the network learns features that are not only relevant for the abnormality classification, but also for the isolation/segmentation of the relevant image regions.

The DenseNet model described in subsection 4.1 is extended to solve the segmentation task. Therefore, we add a decoder network whose input is the returning feature maps of the last dense block (see Figure 1). The decoder architecture is visualized in Figure 5. For the segmentation task, we use the mean squared error loss function:

$$\mathcal{L}_{Seg} = \sum_{i=1}^F \left[ \frac{1}{t} \sum_{z=1}^t (s_{i,z} - p_{i,z})^2 \right], \quad (6)$$

where  $t = 2 \times N \times N$  and  $p_{i,z} \in \mathbf{S}$  denotes the output prediction of the current pixel  $z$  and  $s_{i,z} \in \{0, 1\}$  the corresponding pixel label.

##### 4.5.2. Spatial Classification

We propose to add additional supervision during learning using several approximate spatial labels provided with the PLCO data. For five abnormalities (Nodule, Mass, Infiltrate, Atelectasis, Hilar Abnormality) coarse location information is made available (see Table 4).

The location information is incorporated in the cross-entropy loss using location-specific classes. The spatial

**Table 4**  
Spatial Class Labels for the PLCO data

No.	Region	No.	Region
1	Left lung	6	Upper-middle part
2	Right lung	7	Upper part
3	Lower part	8	Diffused
4	Lower-middle part	9	Multiple (more independent parts)
5	Middle part		

labels  $[b_1, b_2 \dots b_L] \in \{0, 1\}$ , where  $L$  is the total number of spatial classes listed in Table 4, are compared with the network prediction and the loss is calculated (Equation 7).

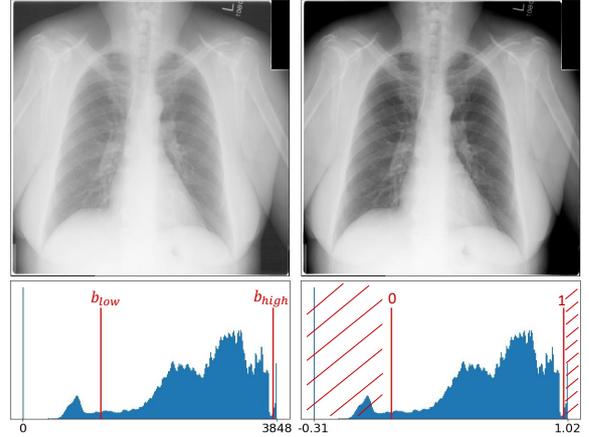
$$\mathcal{L}_{Loc} = - \sum_{m=1}^L \sum_{i=1}^F \left[ w_P^{(m)} b^{(m)} \log(p^{(m)}) + w_N^{(m)} (1 - b^{(m)}) \log(1 - p^{(m)}) \right], \quad (7)$$

where  $w_P^{(m)} = \frac{P_m + N_m}{P_m}$  and  $w_N^{(m)} = \frac{P_m + N_m}{N_m}$ , with  $P_m$  and  $N_m$  indicating, respectively, the number of presence and absence cases of spatial class  $m$  in the training set. The individual localization loss  $\mathcal{L}_{Loc}$  is activated/deactivated dynamically: If spatial labels are not available for abnormality  $n$ , all spatial labels are disregarded and no gradients are computed. Otherwise, the loss is calculated as Equation 7 shows.

#### 4.6. Dynamic Normalization

One challenge in processing chest radiographs is accounting for the large variability of the image appearance, depending on the acquisition source, radiation dose as well as proprietary non-linear postprocessing. In practice, one cannot systematically address this variation due to missing meta-information (e.g., unknown maximum high voltage for images from the ChestX-ray14 dataset). In this context, generic solutions have been proposed for the normalization of radiographs using multi-scale contrast enhancement/leveling techniques (Philipsen et al., 2015; Dippel et al., 2002).

For our diagnostic application, we propose to explicitly avoid altering the image appearance using one of these methods. Instead, we propose an efficient method for dynamically windowing each image, i.e., adjust the brightness and contrast via a linear transformation of the image intensities. Given an arbitrary chest radiograph  $\mathbf{I}$ , let us denote its pixel value histogram function as  $h(x; \mathbf{I})$ . We use a bandwidth of 256 for the image intensity histogram function. Using Gaussian smoothing and median filtering, one can significantly reduce



**Fig. 6.** An original image of the dataset and the corresponding intensity histogram ( $x$ -axis: intensity value;  $y$ -axis: number of pixels) is displayed (left). The preprocessed image where the described normalization technique based on the bounds  $b_{low}$  and  $b_{high}$  is applied (right). The high signal spike in the left histogram at value 0 represents the black anonymization box.

the noise of  $h$  (visible as, e.g., signal spikes due to black background or white text overlay) as well as account for long function tails that affect the windowing of the image. As such, based on the processed function  $h$ , we determine two bounds  $b_{low}$  and  $b_{high}$  which represent a tight intensity window for image  $\mathbf{I}$ . The value  $b_{low}$  indicates the lowest bin and  $b_{high}$  the highest bin along the image intensity histogram. The normalization is applied as follows,  $\mathbf{I} = (\mathbf{I} - b_{low}) / (b_{high} - b_{low})$ . A visual example is shown in Figure 6.

## 5. Experimental Results

In our experiments, we measured the performance of our baseline system (**baseline**) presented in Section 4.1 and 4.2, and quantified the improvement achieved by 1) including the image normalization component (**norm**), 2) segmenting the heart/lung region (**seg**), 3) approximately localizing pathologies within the image (**loc**), and 4) adding the regularized loss terms to compete with the noisy labels (**noise/corr**). All results can be seen in Table 7. Finally, we show the performance of 2 joint models: one including normalization, anatomy information and noise regularization (**baseline + norm + seg + loc + noise**); and the other including normalization, anatomy information and correlation regularization (**baseline + norm + seg + loc + corr**). We refer to the first joint model as **all-noise** and the to the second as **all-corr**. The column **all** includes all previous fea-



**Fig. 7.** Left: Example chest radiograph. Right: Predicted segmentation masks of lungs (blue) and heart (red).

**Table 5**  
Segmentation scores of the multi-task network

	Dice	IoU
Heart	96.8	96.9
Lung	98.2	94.7

**Table 6**  
Spatial classification scores of the multi-task network

Class	AUC
Left Lung	0.861
Right Lung	0.826
Lower	0.835
Lower-middle	0.734
Middle	0.748
Upper-middle	0.781
Upper	0.847
Diffuse	0.846
Multiple	0.709

tures and both regularization schemes. Please recall that we use 70% of the consensus expert reading subset (ChestX-ray14) and 30% of the PLCO dataset. All experiments were measured with the area under the curve (AUC). Corresponding statistical tests (DeLong et al., 1988) were conducted to measure whether the improvement between the baseline and the best model is significant. The proposed values (p-values) are defined to show the significant difference of 2 performance measures, low p-values indicate that they are statistically different. The first column of Table 7 shows our baseline model.

### 5.1. Normalization

Including the normalization step based on dynamic windowing had a two-fold benefit. First, the training time was reduced on average 2-3 times (in terms of average number of epochs). We hypothesize that this is because the normalization ensures images to be more aligned in terms of brightness and contrast, which in some sense simplifies the learning task. Second, this also improved the generalization of the model, and led to a performance increase across all abnormalities, e.g., Effusion (ChestX-ray14) could be improved by 0.026

to 0.949 ( $p = 0.019$ ) as can be seen in the Column 2 (**norm**).

### 5.2. Lung and Heart Segmentation

Column 3 in Table 7 shows improved classification scores when the network was additionally trained to generate lung and heart segmentation masks (**seg**). Some abnormalities were significantly improved, e.g., cardiomegaly from 0.929 to 0.950 ( $p = 0.042$ ).

An example image is visualized in Figure 7 (left). The probabilistic segmentation map is thresholded and overlaid with the image. The red mask defines the heart area, the blue mask indicates the two lungs (right).

### 5.3. Spatial Knowledge

We measured the impact of the location labels on the performance of the classification. A performance gain of all abnormalities supported with spatial information could be observed, as can be seen in the fourth column of Table 7 (**loc**). Atelectasis (PLCO) improved significantly by 0.040 to an AUC score of 0.884 ( $p = 0.047$ ).

In addition, a slight improvement across all abnormalities, can be seen. Thus, we hypothesize that the partial integration of spatial knowledge may also help to improve abnormalities which were not directly supported with spatial information, e.g., due to the abnormality correlation aspect.

### 5.4. Regularization

We evaluated the noise and correlation experiments on the re-read Chest-Xray14 test set only as the calibration set is derived from this data. In Table 7, we can see the performance of both experiments (**noise/corr**). Adding the regularization terms significantly improved the performance on individual abnormalities. We highlight the results for the consolidation class as the abnormality shows the highest label noise (see Table 3) and one of the stronger correlation classes (see Figure 4). The performance for this class improved by 0.024 to 0.836 ( $p = 0.037$ ) with label noise regularization and the consolidation class by 0.019 to 0.831 ( $p = 0.044$ ) with regularization based on label correlation.

### 5.5. All-in-One Joint Model

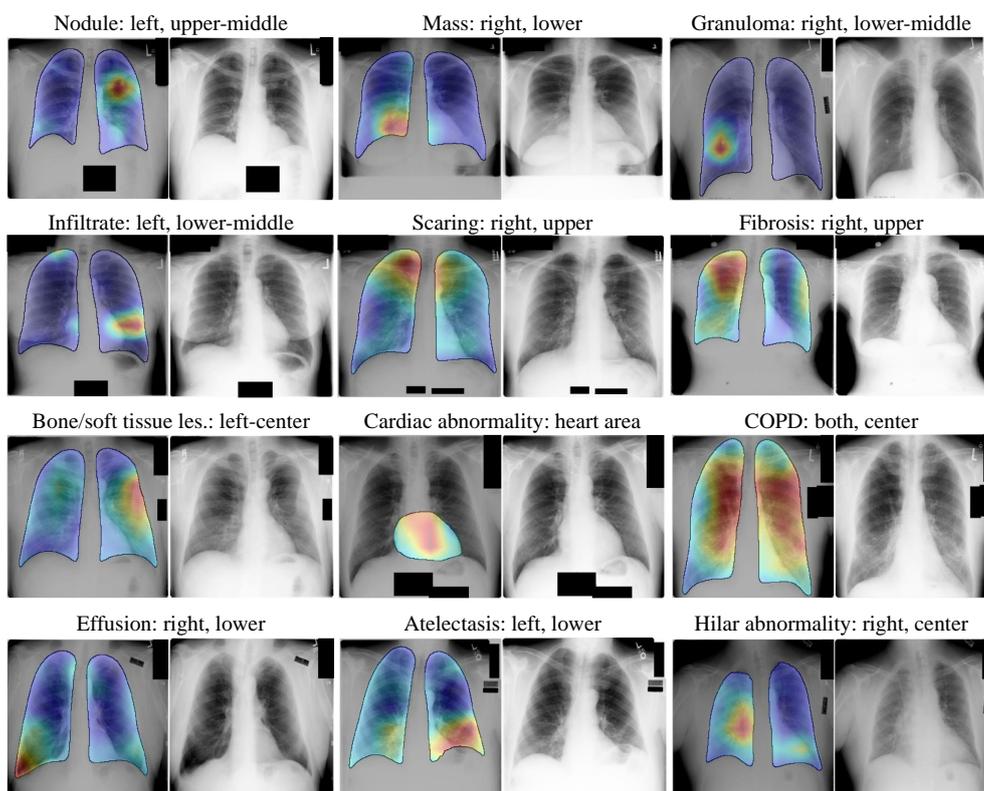
For the multi-task learning network, we combined all previous tasks and features in addition to the baseline model:

- Normalization of images
- Segmentation of lungs and heart

**Table 7**

AUC classification scores for experiments. As test set we use 70% of the consensus expert reading ChestX-ray14 data and 20% of the PLCO data. With + we denote the combination between the baseline system and one feature, e.g., “+norm” denotes the baseline system with the normalization component. The columns “all-noise”, “all-corr”, and “all” denote a combined system with all previous features including noise, correlation regularization, and both, respectively. The last column shows p-values between the best model (bold values) and the baseline model for each class. The p-values were computed using DeLong’s test (DeLong et al., 1988)

Abnormalities	baseline	+norm	anatomy		regularization		combination			p-value	
			+seg	+loc	+noise	+corr	all-noise	all-corr	all		
Expert Read ChestX-ray14	Effusion	0.923	0.949	0.938	0.936	0.940	0.915	<b>0.951</b>	0.933	0.950	0.016
	Cardiomegaly	0.926	0.943	0.950	0.942	0.927	0.940	0.932	0.955	<b>0.957</b>	0.031
	Consolidation	0.812	0.847	0.823	0.837	0.836	0.831	0.838	0.850	<b>0.852</b>	0.024
	Atelectasis	0.821	0.847	0.834	0.819	0.845	0.831	<b>0.851</b>	0.840	0.848	0.045
	Mass	0.804	0.815	0.805	0.842	0.829	0.815	0.830	<b>0.838</b>	<b>0.838</b>	0.021
PLCO	Nodule	0.809	0.817	0.815	0.821	0.817	0.816	0.821	0.825	<b>0.826</b>	< 0.001
	Mass	0.837	0.851	0.858	0.844	0.846	0.845	<b>0.862</b>	0.857	0.859	0.011
	Granuloma	0.881	0.885	0.881	0.883	0.883	0.884	0.883	0.888	<b>0.889</b>	< 0.001
	Infiltrate	0.871	0.882	0.873	<b>0.885</b>	0.878	0.876	0.878	0.879	0.877	0.036
	Scarring	0.842	0.845	0.846	0.845	0.847	0.847	0.847	0.850	<b>0.851</b>	0.001
	Fibrosis	0.870	0.874	0.875	0.870	0.875	0.875	<b>0.877</b>	0.875	<b>0.877</b>	0.015
	B./S. Tissue L.	0.842	0.832	<b>0.844</b>	0.839	0.837	0.834	0.837	0.841	0.838	0.288
	Cardiac Abn.	0.925	0.925	0.927	0.925	0.927	0.926	0.923	0.928	<b>0.929</b>	0.201
	COPD	0.870	0.883	0.876	0.873	0.879	0.877	0.872	<b>0.889</b>	0.884	0.034
	Effusion	0.931	0.950	0.946	0.942	0.953	0.955	0.946	<b>0.956</b>	0.955	0.263
	Atelectasis	0.844	0.853	0.853	0.884	0.875	0.856	0.860	<b>0.889</b>	0.888	0.043
	Hilar Abn.	0.813	0.810	0.826	0.818	0.812	0.823	0.829	0.817	<b>0.830</b>	0.089



**Fig. 8.** Heatmap prediction with GradCAM++ (Chattopadhyay et al., 2018) : Image pairs (left: GradCam++ prediction filtered with lung/heart segmentation masks; right: original image with image normalization) for all 12 PLCO abnormalities from left to right and from top to bottom. All images can be defined as true positive based on a prediction threshold of 0.5.

- Classification of spatial location of abnormalities

In addition, we used either the label noise regularization (**all-noise**) or the label correlation regularization (**all-corr**). The bold values denote the experiment with highest score for each abnormality. The 2 joint experiments share the best scores of most abnormalities. The highest improvement could be achieved on atelectasis. The ChestX-ray14 class improved by 0.030 to 0.851 (**all-noise**) and the PLCO class by 0.045 to 0.889 (**all-corr**). If we consider the best values for each abnormality, an average AUC score of 0.880 was reached.

**Statistical significance:** In order to measure whether the performance scores are statistically significant, we calculated the p-values. In the last column of Table 7, p-values can be seen between the reference and the model leading to the highest AUC score (bold values). Except of “Bone/Soft Tissue Lesion”, “Cardiac Abnormality”, and “Effusion” (PLCO), all abnormalities improved significantly (i.e.  $p < 0.05$ ) including the proposed methods. However, we hypothesize that we reach better performance scores evaluating the models on a clean PLCO test set.

**Performance of other tasks:** We also evaluated the segmentation of heart and lungs. Table 5 shows the dice score and Intersection-over-Union (IoU). The evaluation of the spatial classes can be seen in Table 6.

**Visual Interpretation:** We used GradCAM++ (Chattopadhyay et al., 2018) to show attention maps. In Figure 8, we show image pairs (GradCam++ prediction and original image) for each abnormality of the PLCO dataset. The attention maps are additionally filtered with the help of the lung/heart segmentation masks. For all lung abnormalities, we show attention limited to both lungs and for heart abnormalities, i.e., cardiac abnormality, we show the heart region. The images containing that abnormality were chosen randomly conditioned by a prediction threshold greater than 0.5. All examples show local attention in correspondence to the given abnormality position.

## 6. Discussion

In this study we proposed a multi-task convolutional neural network which outperforms the baseline model in classifying a wide range of abnormalities in chest radiographs. The loss function was designed to deal with noisy labels during training and exploit the correlation between abnormalities. In addition to the classification of the abnormalities, segmentation masks of lungs and heart and approximate spatial classification prediction of the abnormalities shape the architecture

of our network. All images were processed with our dynamic normalization strategy as an unknown degree of variation in the chest radiographs exists.

We showed that training with prior information helped to achieve a performance gain under label noise. By including label noise ratios and label correlation information the performance increased on the 5 analyzed abnormalities of the ChestX-ray14 dataset. A reading process including 4 board-certified radiologists was performed to calibrate and evaluate our proposed strategies on higher quality reference labels. The calibration of the loss regularization was based on a small portion, i.e.  $< 1\%$  of the whole ChestX-ray14 dataset, and therefore, it is expected that using larger re-read dataset can calibrate the labels more accurately. The regularization predominantly improved abnormalities with stronger noise and correlation, e.g., consolidation from 0.812 AUC to 0.836 (noise) and to 0.831 AUC (correlation). We further show performance scores with noise regularization on the PLCO set. However, the improvement on the ChestX-ray14 dataset exceeded as the PLCO data was not re-read. We hypothesize that a re-reading process of the PLCO may help to increase the performance of the corresponding classes.

Information about anatomical structures in chest radiographs and location of abnormalities were included. Certain classes significantly improved, e.g., hilar abnormality from 0.813 to 0.826 by including segmentation masks. We hypothesize that the network interpreted the location of the hilar region better as the abnormality is located close to the lung border. The provided nodule location helps to improve the performance from 0.809 to 0.821 as nodules are often hard to find because they appear as a tiny fraction in the images.

Chest radiographs typically appear with high variation in image quality (e.g., contrast ratio, level of exposure, noise level). This is an important challenge, as algorithms should ideally be applicable for images acquired in many institutions. Our dynamic normalization technique was applied to compensate for the strong intensity variation by applying noise filters and linear transformation of image intensities. Evaluated on two different datasets (ChestX-ray14 and PLCO), we showed that our dynamic normalization can be applied to improve the classification performance of the proposed algorithm. We specifically see on the ChestX-ray14 dataset that dynamic normalization contributes mainly with respect to the performance improvement. Oakden-Rayner (2019) mentioned that there is a substantially low image quality without any standards, i.e. the image intensities vary significantly

	Pat.	Re.	Atel.	Card.	Eff.	Inf.	Mass	Nod.	Pneu.	PTX	Cons.	Ed.	Emph.	Fib.	Pleu.	Hern.
Other Split			Gündel et al. (2019) - DNet													
	×		0.826	0.911	0.885	0.716	0.854	0.774	0.765	0.872	0.806	0.892	0.925	0.820	0.785	0.941
			Huang et al. (2020) - Fusion High-Resolution Network													
	×		0.794	0.902	0.839	0.714	0.827	0.727	0.703	0.848	0.773	0.834	0.911	0.824	0.752	0.916
			Guan and Huang (2018) - Baseline vs. Attention Guided CNN													
			0.832	0.906	0.887	0.717	0.870	0.791	0.732	0.891	0.808	0.905	0.912	0.823	0.802	0.883
		0.853 0.939 0.903 0.754 0.902 0.828 0.774 0.921 0.842 0.924 0.932 0.864 0.837 0.921														
		Rajpurkar et al. (2018) - Radiologist vs. Network														
	×	×	0.808	0.888	0.900	0.734	0.886	0.899	0.823	0.940	0.841	0.910	0.911	0.897	0.779	0.985
	×	×	0.862	0.831	0.901	0.721	0.909	0.894	0.851	0.944	0.893	0.924	0.704	0.806	0.798	0.851
	×	×	Ours - experiment <b>all</b> based on the re-read ChestX-ray14 subset													
	×	×	0.848	0.957	0.950	-	0.838	-	-	-	0.852	-	-	-	-	-
Official Split	×		Gündel et al. (2019) - DNet													
	×		0.767	0.883	0.828	0.709	0.821	0.758	0.731	0.846	0.745	0.835	0.895	0.818	0.761	0.896
			Liu et al. (2019) - Baseline vs. SDFN													
	×		0.762	0.878	0.822	0.693	0.791	0.744	0.707	0.855	0.737	0.837	0.912	0.826	0.760	0.902
	×		0.781	0.885	0.832	0.700	0.815	0.765	0.719	0.866	0.743	0.842	0.921	0.835	0.791	0.911
	×		Ours - experiment <b>all</b> based on the official ChestX-ray14 split													
	×		0.785	0.892	0.836	0.710	0.826	0.755	0.735	0.847	0.747	0.837	0.925	0.838	0.785	0.905

**Table 8**

Performance scores (AUC) of different methods based on the ChestX-ray14 dataset. The abbreviated abnormalities are placed in the same order as visualized in Figure 4 (left). The overall table is separated into methods where the official ChestX-ray14 split and other splits were applied. Furthermore, we highlight key factors of the methods: We separate the methods into random split and patient-wise split (Pat.), additionally, we denote if a subset is re-read by radiologists (Re.). Note that reported performance is influenced by different factors (e.g., data selection, evaluation strategy, reference standard) and therefore, direct comparison between different methods for other splitting strategies is not trivial.

across the images. Therefore, the benefit of using the dynamic normalization may be more pronounced on several abnormalities than the benefit of using the proposed regularization. However, there are also other abnormalities where our regularization technique predominated, e.g. mass with 0.815 (normalization) and 0.829 (noise regularization).

The overall aim of the paper is to build an automated method that can simultaneously predict multiple abnormalities. Based on the complexity of classifying multiple abnormalities with different characteristics, one integrated method cannot help to improve all classes. For that reason, our investigation were focused in searching strategies to improve most of the given abnormalities. In combination of the proposed methods we achieved a significant improvement of the performance for 14 of 17 abnormality classes.

An observer study was performed to evaluate the performance of the NLP labels against radiologists-based labels. While the NLP-based labels maintain high specificity, we observed low sensitivity scores in all relevant abnormalities. This means that using only NLP-based labels may not be ideal for developing deep learning algorithms. The strong label difference be-

tween original labels and our radiologist study disclose the high bias that needs to be taken into account in interpreting the reported performance of algorithms trained using NLP-based labels. Large scale radiologist-based labels remain crucial in this field and, if available, a re-labeling process is necessary in order to uncover label noise. While obtaining radiologists-based labels is crucial, inter and intra-observer variability can still impact the performance. To improve this process, more than 4 radiologists may be considered to further reduce label noise. In this work, the main limitation was the small subset of re-read labels (689 cases), for which 30% were used to find the parameters for regularization. A higher number of cases for the calibration may further increase the performance. Moreover, we used the original labels for the validation set due to the limited amount of re-read examples. For a more precise validation, the validation set labels also need to be re-read.

As we proposed several sub-strategies to increase the performance, our proposed system required a lot of prior work. On the one hand, the data for the auxiliary tasks, i.e., lung/heart segmentation masks and spatial classification labels need to be available

to train our multi-task system. On the other hand, the prior re-reading process is expensive, time-consuming, and therefore limited in terms of scalability. However, we show in this work that a small subset is sufficient to obtain a performance gain. Nevertheless, to entirely prevent a prior re-reading process, unsupervised approaches may be considered which adapt the loss depending on noise level probabilities. For example Arazo et al. (2019) applied a mixture model to find noisy labels. Some other approaches claim that the mean absolute error function shows more robustness against label noise if weighing training examples with uncertainty differently. (Wang et al., 2019; Zhang and Sabuncu, 2018). However, these methods have solely been tested on standard, non-medical datasets.

In the field of chest radiographic abnormality classification, most studies focus on the performance evaluation with the AUC score. The high publication rate helps to analyse proposed methods by comparing the performance scores. However, these performance scores are often directly compared without prior knowledge about strongly influencing factors. In Table 8, we show different studies and highlight essential factors in order to avoid a direct performance comparison between methods. As highlighted in Section 2, most groups split data randomly, which leads to unique patients in both training and testing set. For example the ChestX-ray14 dataset has an average of 3.6 images per patient. Therefore, the split should be generated patient-wise which is denoted with “Pat.” in Table 8.

**Official Split:** In this table, we show the performance of the model proposed in Gündel et al. (2019) using the official patient-wise splits. In Liu et al. (2019), both a baseline model and a proposed method, described in Section 2, were evaluated. The (official) patient-wise split used for both models thus enabling a fair comparison and further enables a comparison with other work. In our approach (last row), we show performance scores of our proposed model **all**. However, to ensure a fair comparison, we changed to the official ChestX-ray14 split and re-trained the model, in contrast to Table 7. We observed improved AUC scores on most abnormalities compared to the other approaches using the official split.

**Other Split:** In Gündel et al. (2019), the same model was used both for the official and another split. Even if both splits are patient-wise, a significant performance difference can be seen across all abnormalities. Thus, a direct comparison of the performance should only be considered if the same split is used since there is a significant performance variability by using different

test sets. Another patient-wise split was applied in Huang et al. (2020). Because of the missing baseline, a comparison to other approaches should be avoided. Guan and Huang (2018) proposed an approach with an additional baseline model on another split, however, the splits were separated randomly.

Because of the high label noise ratio, Rajpurkar et al. (2018) evaluated their model on a re-read subset (“Re.” in Table 8). The radiologist performance was compared with the model performance across all 14 abnormalities. Finally, we show the performance scores of the 5 abnormalities (Table 7, **all**) based on our re-read subset.

In general, we highlight the importance of the applied splitting strategies to enable or avoid performance comparisons. We see a high performance variance when using different reference labels, e.g. cardiomegaly from 0.831 AUC (Rajpurkar et al., 2018) to 0.939 AUC (Guan and Huang, 2018). In addition to the increasing availability of chest radiograph datasets for abnormality classification, radiologist studies and different splitting strategies lead to many subsets with different labels. If no re-read subset can be established, a self-defined baseline model should therefore be a consistent reference to compare the proposed methods in the future.

## 7. Conclusion

We developed a novel method that achieves state-of-the-art performance for the classification of 17 abnormalities based on chest radiographs. A radiologist study enabled an analysis of the accuracy of NLP labels compared to the consensus of 4 board-certified radiologists indicated a large fraction of label noise. We introduced a novel training strategy to handle label noise in the training data. In addition, a regularizer was implemented to exploit the label correlation information during training. The network was expanded to predict segmentation masks of the underlying anatomy as well as the spatial classes for the location of the abnormalities. We implemented a dynamic image normalization technique as the scans are derived from different scanners and post-processing methods. We show that our proposed method can be used to significantly improve the performance of deep learning methods trained on datasets with noisy labels.

## Acknowledgement

The authors thank the National Cancer Institute (NCI) for access to their data collected by the Prostate, Lung, Colorectal and Ovarian (PLCO) Cancer Screening Trial. The authors thank the National Institutes of Health (NIH) for access to the ChestX-ray14 collection. The statements contained herein are solely those of the authors and do not represent or imply concurrence or endorsement by NCI or NIH.

**Disclaimer:** The concepts and information presented in this paper are based on research results that are not commercially available. This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

## References

- Amyar, A., Modzelewski, R., Ruan, S., 2020. Multi-task deep learning based CT imaging analysis for COVID-19: Classification and segmentation. doi:10.1101/2020.04.16.20064709.
- Arazo, E., Ortego, D., Albert, P., O'Connor, N., McGuinness, K., 2019. Unsupervised label noise modeling and loss correction.
- Ardila, D., Kiraly, A.P., Bharadwaj, S., Choi, B., Reicher, J.J., Peng, L., Tse, D., Etemadi, M., Ye, W., Corrado, G., Naidich, D.P., Shetty, S., 2019. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature Medicine* 25, 954–961.
- Ba, J., Mnih, V., Kavukcuoglu, K., 2014. Multiple object recognition with visual attention. *CoRR abs/1412.7755*. arXiv:1412.7755.
- Balabanova, Y., Coker, R., Fedorin, I., Zakharova, S., Plavinskij, S., Krukov, N., Atun, R., Drobniowski, F., 2005. Variability in interpretation of chest radiographs among russian clinicians and implications for screening programmes: Observational study. *BMJ (Clinical research ed.)* 331, 379–82. doi:10.1136/bmj.331.7513.379.
- Barbosa Jr, E., Gefter, W., Yang, R., Ghesu, F., Liu, S., Mailhe, B., Mansoor, A., Grbic, S., Piat, S., Chabin, G., Balachandran, A., Vogt, S., Ziebandt, V., Kappler, S., Comaniciu, D., 2020. Automated detection and quantification of COVID-19 airspace disease on chest radiographs: A novel approach achieving radiologist-level performance using a cnn trained on digital reconstructed radiographs (drrs) from ct-based ground-truth, *Investigative Radiology*.
- Bayliss, E., Ellis, J., Steiner, J., 2005. Subjective assessments of comorbidity correlate with quality of life health outcomes: Initial validation of a comorbidity assessment instrument. *Health and quality of life outcomes* 3, 51. doi:10.1186/1477-7525-3-51.
- Brady, A.P., 2016. Error and discrepancy in radiology: Inevitable or avoidable? *Insights into Imaging* 8, 171–182. doi:10.1007/s13244-016-0534-1.
- Bruno, M.A., Walker, E.A., Abujudeh, H.H., 2015. Understanding and confronting our mistakes: The epidemiology of error in radiology and strategies for error reduction. *RadioGraphics* 35, 1668–1676.
- Cai, J., Lu, L., Harrison, A.P., Shi, X., Chen, P., Yang, L., 2018. Iterative attention mining for weakly supervised thoracic disease pattern localization in chest X-Rays, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*, Springer International Publishing, Cham. pp. 589–598.
- Chattopadhyay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N., 2018. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks, in: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 839–847. doi:10.1109/WACV.2018.00097.
- DeLong, E., DeLong, D., Clarke-Pearson, D., 1988. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* 44, 837–845. doi:10.2307/2531595.
- Dippel, S., Stahl, M., Wiemker, R., Blaffert, T., 2002. Multiscale contrast enhancement for radiographies: Laplacian pyramid versus fast wavelet transform. *IEEE Transactions on Medical Imaging* 21, 343–353.
- Dou, H., Warfield, S., Gholipour, A., 2020. Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. *Medical Image Analysis* 65, 101759. doi:10.1016/j.media.2020.101759.
- Fleishon, H.B., Bhargavan, M., Meghea, C., 2006. Radiologists' reading times using PACS and using films: One practice's experience. *Academic Radiology* 13, 453–460.
- Gaál, G., Maga, B., Lukács, A., 2020. Attention U-net based adversarial architectures for chest X-ray lung segmentation.
- Gal, Y., Ghahramani, Z., 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, in: Balcan, M.F., Weinberger, K.Q. (Eds.), *Proceedings of The 33rd International Conference on Machine Learning, PMLR, New York, New York, USA*. pp. 1050–1059. URL: <http://proceedings.mlr.press/v48/gal16.html>.
- Ghamrawi, N., McCallum, A., 2005. Collective multi-label classification, in: *Proceedings of the 14th ACM International Conference on Information and Knowledge Management, ACM, New York, NY, USA*. pp. 195–200. doi:10.1145/1099554.1099591.
- Ghesu, F., Georgescu, B., Zheng, Y., Grbic, S., Maier, A., Hornegger, J., Comaniciu, D., 2019a. Multi-scale deep reinforcement learning for real-time 3D-landmark detection in CT scans. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 176–189. doi:10.1109/TPAMI.2017.2782687.
- Ghesu, F.C., Georgescu, B., Gibson, E., Guendel, S., Kalra, M.K., Singh, R., Digumarthy, S.R., Grbic, S., Comaniciu, D., 2019b. Quantifying and leveraging classification uncertainty for chest radiograph assessment, in: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, Springer International Publishing, Cham. pp. 676–684.
- Ghesu, F.C., Georgescu, B., Mansoor, A., Yoo, Y., Gibson, E., Vishwanath, R., Balachandran, A., Balter, J.M., Cao, Y., Singh, R., Digumarthy, S.R., Kalra, M.K., Grbic, S., Comaniciu, D., 2021. Quantifying and leveraging predictive uncertainty for medical image assessment. *Medical Image Analysis* 68, 101855. doi:https://doi.org/10.1016/j.media.2020.101855.
- Gómez, O., Mesejo, P., Ibáñez, O., Valsecchi, A., Cordon, O., 2020. Deep architectures for high-resolution multi-organ chest X-ray image segmentation. *Neural Computing and Applications* 32. doi:10.1007/s00521-019-04532-y.
- Gohagan, J.K., Prorok, P.C., Hayes, R.B., Kramer, B.S., 2000. The prostate, lung, colorectal and ovarian (PLCO) cancer screening trial of the national cancer institute: History, organization, and status. *Controlled clinical trials* 21, 251S–272S.
- Guan, Q., Huang, Y., 2018. Multi-label chest X-ray image classification via category-wise residual attention learning. *Pattern Recognition Letters* doi:https://doi.org/10.1016/j.patrec.2018.10.027.
- Guan, Q., Huang, Y., Zhong, Z., Zheng, Z., Zheng, L., Yang, Y., 2018. Diagnose like a radiologist: Attention guided convolutional neural network for thorax disease classification. *CoRR abs/1801.09927*.

- Guan, W.j., Liang, W.h., Zhao, Y., Liang, H.r., Chen, Z.s., Li, Y.m., Liu, X.q., Chen, R.c., Tang, C.l., Wang, T., Ou, C.q., Li, L., Chen, P.y., Sang, L., Wang, W., Li, J.f., Li, C.c., Ou, L.m., Cheng, B., Xiong, S., Ni, Z.y., Xiang, J., Hu, Y., Liu, L., Shan, H., Lei, C.l., Peng, Y.x., Wei, L., Liu, Y., Hu, Y.h., Peng, P., Wang, J.m., Liu, J.y., Chen, Z., Li, G., Zheng, Z.j., Qiu, S.q., Luo, J., Ye, C.j., Zhu, S.y., Cheng, L.l., Ye, F., Li, S.y., Zheng, J.p., Zhang, N.f., Zhong, N.s., He, J.x., 2020. Comorbidity and its impact on 1590 patients with COVID-19 in china: A nationwide analysis. *European Respiratory Journal* doi:10.1183/13993003.00547-2020.
- Gündel, S., Grbic, S., Georgescu, B., Liu, S., Maier, A., Comaniciu, D., 2019. Learning to recognize abnormalities in chest X-rays with location-aware dense networks, in: *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, pp. 757–765.
- Hu, J., Shen, L., Sun, G., 2018. Squeeze-and-excitation networks, in: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7132–7141. doi:10.1109/CVPR.2018.00745.
- Huang, G., Liu, Z., v. d. Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks, in: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269. doi:10.1109/CVPR.2017.243.
- Huang, Z., Lin, J., Xu, L., Wang, H., Bai, T., Pang, Y., Meen, T.H., 2020. Fusion high-resolution network for diagnosing ChestX-ray images. *Electronics* 9, 190. doi:10.3390/electronics9010190.
- Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Illcus, S., Chute, C., Marklund, H., Haghighi, B., Ball, R., Shpanskaya, K., et al., 2019. CheXpert: a large chest radiograph dataset with uncertainty labels and expert comparison, in: *Thirty-Third AAAI Conference on Artificial Intelligence*.
- Islam, M.T., Aowal, M.A., Minhaz, A.T., Ashraf, K., 2017. Abnormality detection and localization in chest X-rays using deep convolutional neural networks. *CoRR* abs/1705.09850. arXiv:1705.09850.
- Jøsang, A., 2016. *Subjective Logic: A Formalism for Reasoning Under Uncertainty*. 1st ed., Springer Publishing Company, Incorporated.
- Jusoh, S., 2018. A study on nlp applications and ambiguity problems. *Journal of Theoretical and Applied Information Technology* 96, 1486–1499.
- Kholiavchenko, M., Sirazitdinov, I., Kubrak, K., Badrutdinova, R., Kuleev, R., Yuan, Y., Vrtovec, T., Ibragimov, B., 2020. Contour-aware multi-label chest X-ray organ segmentation. *International Journal of Computer Assisted Radiology and Surgery* 15. doi:10.1007/s11548-019-02115-9.
- Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization, in: *Bengio, Y., LeCun, Y. (Eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Lakshminarayanan, B., Pritzel, A., Blundell, C., 2017. Simple and scalable predictive uncertainty estimation using deep ensembles, in: *Proceedings of the 31st International Conference on Neural Information Processing Systems, Curran Associates Inc., Red Hook, NY, USA*. p. 6405–6416.
- Larrazabal, A.J., Martínez, C., Glocker, B., Ferrante, E., 2020. Post-DAE: Anatomically plausible segmentation via post-processing with denoising autoencoders. *IEEE Transactions on Medical Imaging* 39, 3813–3820. doi:10.1109/TMI.2020.3005297.
- Li, W., Cao, P., Zhao, D., Wang, J., 2016. Pulmonary nodule classification with deep convolutional neural networks on computed tomography images. *Computational and Mathematical Methods in Medicine* 2016, 1–7.
- Li, Z., Wang, C., Han, M., Xue, Y., Wei, W., Li, L., Fei-Fei, L., 2018. Thoracic disease identification and localization with limited supervision, in: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8290–8299. doi:10.1109/CVPR.2018.00865.
- Liu, H., Wang, L., Nan, Y., Jin, F., Wang, Q., Pu, J., 2019. SDFN: Segmentation-based deep fusion network for thoracic disease classification in chest X-ray images. *Computerized Medical Imaging and Graphics* 75, 66–73. doi:https://doi.org/10.1016/j.compmedimag.2019.05.005.
- Molchanov, D., Ashukha, A., Vetrov, D., 2017. Variational dropout sparsifies deep neural networks, in: *Proceedings of the 34th International Conference on Machine Learning - Volume 70, JMLR.org*. p. 2498–2507.
- Oakden-Rayner, L., 2017. Exploring the chestxray14 dataset: problems. <https://lukeoakdenrayner.wordpress.com/2017/12/18/the-chestxray14-dataset-problems/>. Accessed: 2017-12-18.
- Oakden-Rayner, L., 2019. Half a million x-rays! first impressions of the Stanford and MIT chest x-ray datasets. <https://lukeoakdenrayner.wordpress.com/2019/02/25/half-a-million-x-rays-first-impressions-of-the-stanford-and-mit-chest-x-ray-datasets/>. Accessed: 2019-02-25.
- Philipsen, R.H.H.M., Maduskar, P., Hogeweg, L., Melendez, J., Sánchez, C.I., van Ginneken, B., 2015. Localized energy-based normalization of medical images: Application to chest radiography. *IEEE Transactions on Medical Imaging* 34, 1965–1975.
- Rajkomar, A., Oren, E., Chen, K., Dai, A.M., Hajaj, N., Hardt, M., Liu, P.J., Liu, X., Marcus, J., Sun, M., Sundberg, P., Yee, H., Zhang, K., Zhang, Y., Flores, G., Duggan, G.E., Irvine, J., Le, Q., Litsch, K., Mossin, A., Tansuwan, J., Wang, D., Wexler, J., Wilson, J., Ludwig, D., Volchenboum, S.L., Chou, K., Pearson, M., Madabushi, S., Shah, N.H., Butte, A.J., Howell, M.D., Cui, C., Corrado, G.S., Dean, J., 2018. Scalable and accurate deep learning with electronic health records. *npj Digital Medicine* 1.
- Rajpurkar, P., Irvin, J., Ball, R.L., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C.P., Patel, B.N., Yeom, K.W., Shpanskaya, K., Blankenberg, F.G., Seekins, J., Amrhein, T.J., Mong, D.A., Halabi, S.S., Zucker, E.J., Ng, A.Y., Lungren, M.P., 2018. Deep learning for chest radiograph diagnosis: A retrospective comparison of the CheXNeXt algorithm to practicing radiologists. *PLOS Medicine* 15, 1–17. doi:10.1371/journal.pmed.1002686.
- Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C., Shpanskaya, K., et al., 2017. CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning abs/1711.05225.
- Rolnick, D., Veit, A., Belongie, S., Shavit, N., 2017. Deep learning is robust to massive label noise .
- Rubin, J., Sanghavi, D., Zhao, C., Lee, K., Qadir, A., Xu-Wilson, M., 2018. Large scale automated reading of frontal and lateral chest X-Rays using dual convolutional neural networks. *CoRR* abs/1804.07839.
- Rusiecki, A., 2019. Trimmed Robust Loss Function for Training Deep Neural Networks with Label Noise. pp. 215–222. doi:10.1007/978-3-030-20912-4\_21.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L., 2015. ImageNet large scale visual recognition challenge. *IJCV* 115, 211–252.
- Selvan, R., Dam, E., Detlefsen, N., Rischel, S., Sheng, K., Nielsen, M., Pai, A., 2020. Lung segmentation from chest X-rays using variational data imputation.
- Shen, Y., Gao, M., 2018. Dynamic routing on deep neural network for thoracic disease classification and sensitive area localization, in: *Shi, Y., Suk, H.I., Liu, M. (Eds.), Machine Learning in Medical Imaging, Springer International Publishing, Cham*. pp. 389–397.

- Tang, Y., Wang, X., Harrison, A.P., Lu, L., Xiao, J., Summers, R.M., 2018. Attention-guided curriculum learning for weakly supervised classification and localization of thoracic diseases on chest radiographs, in: Shi, Y., Suk, H.I., Liu, M. (Eds.), *Machine Learning in Medical Imaging*, Springer International Publishing, Cham. pp. 249–258.
- Van Eeden, S., Leipsic, J., Paul Man, S.F., Sin, D.D., 2012. The relationship between lung inflammation and cardiovascular disease. *American Journal of Respiratory and Critical Care Medicine* 186, 11–16. doi:10.1164/rccm.201203-0455PP, arXiv:https://doi.org/10.1164/rccm.201203-0455PP. PMID: 22538803.
- Wang, H., Jia, H., Lu, L., Xia, Y., 2020. Thorax-Net: An attention regularized deep neural network for classification of thoracic diseases on chest radiography. *IEEE Journal of Biomedical and Health Informatics* 24, 475–485. doi:10.1109/JBHI.2019.2928369.
- Wang, H., Wang, S., Qin, Z., Zhang, Y., Li, R., Xia, Y., 2021. Triple attention learning for classification of 14 thoracic diseases using chest radiography. *Medical Image Analysis* 67, 101846. doi:10.1016/j.media.2020.101846.
- Wang, X., Kodirov, E., Hua, Y., Robertson, N., 2019. IMAE for noise-robust learning: Mean absolute error does not treat examples equally and gradient magnitude’s variance matters.
- Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M., 2017. ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases, in: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3462–3471.
- Yan, C., Yao, J., Li, R., Xu, Z., Huang, J., 2018. Weakly supervised deep learning for thoracic disease classification and localization on chest X-rays, in: *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, pp. 103–110. doi:10.1145/3233547.3233573.
- Yao, L., Poblenz, E., Dagunts, D., Covington, B., Bernard, D., Lyman, K., 2017. Learning to diagnose from scratch by exploiting dependencies among labels. CoRR abs/1710.10501.
- Yao, L., Prosky, J., Poblenz, E., Covington, B., Lyman, K., 2018. Weakly supervised medical diagnosis and localization from multiple resolutions. CoRR abs/1803.07703. arXiv:1803.07703.
- Zhang, Z., Sabuncu, M., 2018. Generalized cross entropy loss for training deep neural networks with noisy labels.
- Zhu, S., Ji, X., Xu, W., Gong, Y., 2005. Multi-labelled classification using maximum entropy method, pp. 274–281. doi:10.1145/1076034.1076082.
- Zhu, W., Liu, C., Fan, W., Xie, X., 2018. Deeplung: Deep 3D dual path nets for automated pulmonary nodule detection and classification. *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 673–681.
- Zotin, A., Hamad, Y., Simonov, K., Kurako, M., 2019. Lung boundary detection for chest X-ray images classification based on GLCM and probabilistic neural networks. *Procedia Computer Science* 159, 1439–1448. doi:10.1016/j.procs.2019.09.314.