

CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning

Pranav Rajpurkar^{*1} Jeremy Irvin^{*1} Kaylie Zhu¹ Brandon Yang¹ Hershel Mehta¹
 Tony Duan¹ Daisy Ding¹ Aarti Bagul¹ Curtis Langlotz² Katie Shpanskaya²
 Matthew P. Lungren² Andrew Y. Ng¹

Abstract

We develop an algorithm that can detect pneumonia from chest X-rays at a level exceeding practicing radiologists. Our algorithm, CheXNet, is a 121-layer convolutional neural network trained on ChestX-ray14, currently the largest publicly available chest X-ray dataset, containing over 100,000 frontal-view X-ray images with 14 diseases. Four practicing academic radiologists annotate a test set, on which we compare the performance of CheXNet to that of radiologists. We find that CheXNet exceeds average radiologist performance on pneumonia detection on both sensitivity and specificity. We extend CheXNet to detect all 14 diseases in ChestX-ray14 and achieve state of the art results on all 14 diseases.

1. Introduction

More than 1 million adults are hospitalized with pneumonia and around 50,000 die from the disease every year in the US alone (CDC, 2017). Chest X-rays are currently the best available method for diagnosing pneumonia (WHO, 2001), playing a crucial role in clinical care (Franquet, 2001) and epidemiological studies (Cherian et al., 2005). However, detecting pneumonia in chest X-rays is a challenging task that relies on the availability of expert radiologists. In this work, we present a model that can automatically detect pneumonia from chest X-rays at a level exceeding practicing

^{*}Equal contribution ¹Stanford University, Computer Science Department ²Stanford University, Medical School. Correspondence to: Pranav Rajpurkar <pranavs@cs.stanford.edu>, Jeremy Irvin <jirvin16@cs.stanford.edu>.

Project website at <https://stanfordmlgroup.github.io/projects/chexnet>

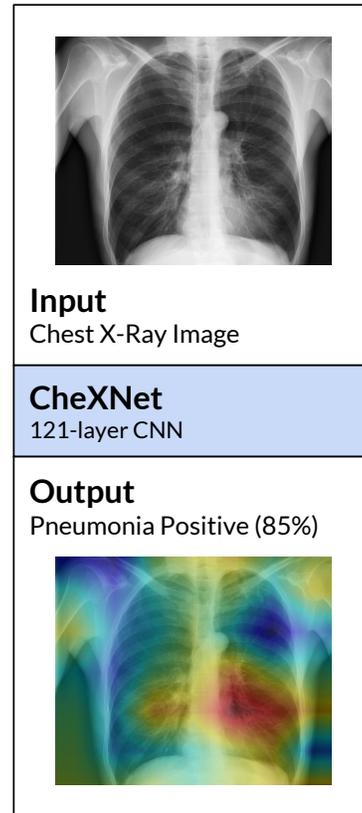


Figure 1. CheXNet is a 121-layer convolutional neural network that takes a chest X-ray image as input, and outputs the probability of a pathology. On this example, CheXNet correctly detects pneumonia and also localizes areas in the image most indicative of the pathology.

radiologists.

Our model, CheXNet (shown in Figure 1), is a 121-layer convolutional neural network that inputs a chest X-ray image and outputs the probability of pneumonia along with a heatmap localizing the areas of the image most indicative of pneumonia. We train CheXNet on the recently released ChestX-ray14 dataset (Wang

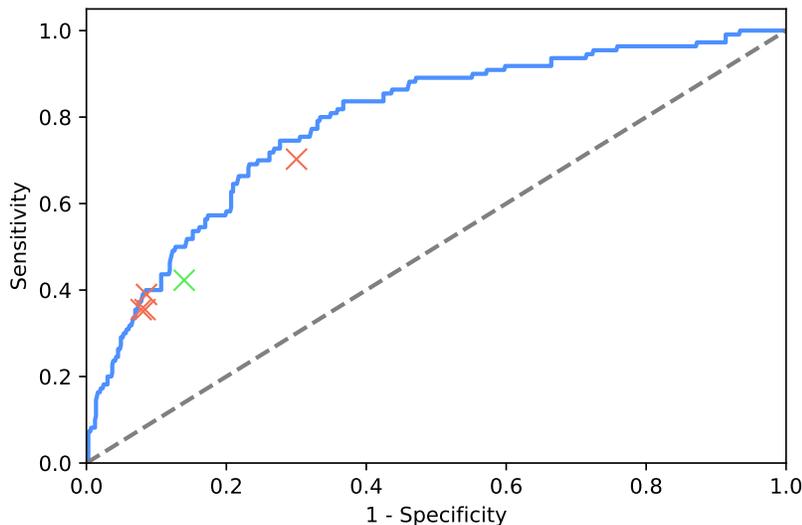


Figure 2. CheXNet outperforms the average of the radiologists at pneumonia detection using X-ray images. CheXNet is tested against 4 practicing radiologists on sensitivity (which measures the proportion of positives that are correctly identified as such) and specificity (which measures the proportion of negatives that are correctly identified as such). A single radiologist’s performance is represented by an orange marker, while the average is represented by green. CheXNet outputs the probability of detecting pneumonia in a Chest X-ray, and the blue curve is generated by varying the thresholds used for the classification boundary. The sensitivity-specificity point for each radiologist and for the average lie below the blue curve, signifying that CheXNet is able to detect pneumonia at a level matching or exceeding radiologists.

et al., 2017), which contains 112,120 frontal-view chest X-ray images individually labeled with up to 14 different thoracic diseases, including pneumonia. We use dense connections (Huang et al., 2016) and batch normalization (Ioffe & Szegedy, 2015) to make the optimization of such a deep network tractable.

Detecting pneumonia in chest radiography can be difficult for radiologists. The appearance of pneumonia in X-ray images is often vague, can overlap with other diagnoses, and can mimic many other benign abnormalities. These discrepancies cause considerable variability among radiologists in the diagnosis of pneumonia (Neuman et al., 2012; Davies et al., 1996; Hopstaken et al., 2004). To estimate radiologist performance, we collect annotations from four practicing academic radiologists on a subset of 420 images from ChestX-ray14. On these 420 images, we measure performance of individual radiologists using the majority vote of other radiologists as ground truth, and similarly measure model performance.

We find that the model exceeds the average radiologist performance at the pneumonia detection task on both sensitivity and specificity. To compare CheXNet against previous work using ChestX-ray14, we make simple modifications to CheXNet to detect all 14 dis-

eases in ChestX-ray14, and find that we outperform best published results on all 14 diseases. Automated detection of diseases from chest X-rays at the level of expert radiologists would not only have tremendous benefit in clinical settings, it would be invaluable in delivery of health care to populations with inadequate access to diagnostic imaging specialists.

2. CheXNet

2.1. Problem Formulation

The pneumonia detection task is a binary classification problem, where the input is a frontal-view chest X-ray image X and the output is a binary label $t \in \{0, 1\}$ indicating the absence or presence of pneumonia respectively. For a single example in the training set, we optimize the binary cross entropy loss

$$L(X, t) = -t \log p(T = 1|X) - (1 - t) \log p(T = 0|X)$$

where $p(T = i|X)$ is the probability that the network assigns to the label i .

2.2. Model Architecture and Training

CheXNet is a 121-layer Dense Convolutional Network (DenseNet) (Huang et al., 2016) trained on the

ChestX-ray 14 dataset. DenseNets improve flow of information and gradients through the network, making the optimization of very deep networks tractable. We replace the final fully connected layer with one that has a single output, after which we apply a sigmoid nonlinearity, outputting the probability that the image contains pneumonia.

The weights of the network are randomly initialized and trained end-to-end using Adam with standard parameters ($\beta_1 = 0.9$ and $\beta_2 = 0.999$) (Kingma & Ba, 2014). We train the model using minibatches of size 16, where we oversample the minority (positive) class (Buda et al., 2017). We use an initial learning rate of 0.01 that is decayed by a factor of 10 each time the validation loss plateaus after an epoch, and pick the model with the lowest validation loss.

3. Data

3.1. Training

We use the ChestX-ray14 dataset released by Wang et al. (2017) which contains 112,120 frontal-view X-ray images of 30,805 unique patients. Wang et al. (2017) annotate each image with up to 14 different thoracic pathology labels using automatic extraction methods on radiology reports. We label images that have pneumonia as one of the annotated pathologies as positive examples and label all other images as negative examples for the pneumonia detection task. We randomly split the entire dataset into 80% training, and 20% validation.

Before inputting the images into the network, we downscale the images to 224×224 and normalize based on the mean and standard deviation of images in the ImageNet training set. We also augment the training data with random horizontal flipping.

3.2. Test

We collected a test set of 420 frontal chest X-rays. Annotations were obtained independently from four practicing radiologists at Stanford University, who were asked to label all 14 pathologies in Wang et al. (2017). The radiologists had 4, 7, 25, and 28 years of experience, and one of the radiologists is a sub-specialty fellowship trained thoracic radiologist. Radiologists did not have access to any patient information or knowledge of disease prevalence in the data. Labels were entered into a standardized data entry program.

4. CheXNet vs. Radiologist Performance

We assess radiologist performance on the test set on the pneumonia detection task. Recall that each of the images in test420 has a ground truth label from 4 practicing radiologists. We evaluate the performance of an individual radiologist by using the majority vote of the other 3 radiologists as ground truth. Similarly, we evaluate CheXNet using the majority vote of 3 of 4 radiologists, repeated four times to cover all groups of 3.

We compare CheXNet against radiologists on the Receiver Operating Characteristic (ROC) curve, which plots model sensitivity against 1 - specificity. Figure 2 illustrates the model ROC curve as well as the four radiologist and average radiologist operating points: a single radiologist’s performance is represented by an orange marker, while the average is represented by green. CheXNet outputs the probability of detecting pneumonia in a Chest X-ray, and the ROC curve is generated by varying the thresholds used for the classification boundary. CheXNet has an AUROC of 0.788. The sensitivity-specificity point for each radiologist and for the average lie below the ROC curve, signifying that CheXNet is able to detect pneumonia at a level matching or exceeding radiologists.

We identify two limitations of this setup. First, neither the model nor the radiologists were permitted to use prior examinations or patient history, which have been shown to decrease radiologist performance (Berbaum et al., 1985; Potchen et al., 1979). Second, only frontal radiographs were presented to the radiologists and model during diagnosis, but it has been shown that up to 15% of accurate diagnoses require the lateral view (Raouf et al., 2012). We thus expect that this setup provides a conservative estimate of human radiologist performance.

5. ChexNet vs. Previous State of the Art on the ChestX-ray14 Dataset

We extend the algorithm to classify multiple thoracic pathologies by making three changes. First, instead of outputting one binary label, ChexNet outputs a vector t of binary labels indicating the absence or presence of each of the following 14 pathology classes: Atelectasis, Cardiomegaly, Consolidation, Edema, Effusion, Emphysema, Fibrosis, Hernia, Infiltration, Mass, Nodule, Pleural Thickening, Pneumonia, and Pneumothorax. Second, we replace the final fully connected layer in CheXNet with a fully connected layer producing a 14-dimensional output, after which we apply an ele-

Pathology	Wang et al. (2017)	Yao et al. (2017)	CheXNet (ours)
Atelectasis	0.716	0.772	0.8209
Cardiomegaly	0.807	0.904	0.9048
Effusion	0.784	0.859	0.8831
Infiltration	0.609	0.695	0.7204
Mass	0.706	0.792	0.8618
Nodule	0.671	0.717	0.7766
Pneumonia	0.633	0.713	0.7632
Pneumothorax	0.806	0.841	0.8932
Consolidation	0.708	0.788	0.7939
Edema	0.835	0.882	0.8932
Emphysema	0.815	0.829	0.9260
Fibrosis	0.769	0.767	0.8044
Pleural Thickening	0.708	0.765	0.8138
Hernia	0.767	0.914	0.9387

Table 1. CheXNet outperforms the best published results on all 14 pathologies in the ChestX-ray14 dataset. In detecting Mass, Nodule, Pneumonia, Pneumothorax, and Emphysema, CheXNet has a margin of >0.05 AUROC over previous state of the art results.

mentwise sigmoid nonlinearity. The final output is the predicted probability of the presence of each pathology class. Third, we modify the loss function to optimize the sum of binary cross entropy losses

$$L(X, t) = \sum_{c=1}^{14} [-t_c \log p(T_c = 1|X) - (1 - t_c) \log p(T_c = 0|X)],$$

where $p(T_c = 1|X)$ is the predicted probability that the image contains the pathology c and $p(T_c = 0|X)$ is the predicted probability that the image does not contain the pathology c .

Following previous work on ChestX-ray14 (Wang et al., 2017; Yao et al., 2017), we randomly split the dataset into training (70%), validation (10%), and test (20%) sets; the choice of split has been shown to have insignificant effect on performance. We compare the per-class Area Under the ROC Curve (AUROC) of the model against the previous state of the art held by Yao et al. (2017) on 13 classes and Wang et al. (2017) on the remaining 1 class.

We find that CheXNet achieves state of the art results on all 14 pathology classes. Table 1 illustrates the per-class AUROC comparison on the test set. On Atelectasis, Mass, Nodule, Pneumonia, and Emphysema, we outperform previous state of the art considerably (> 0.05 increase in AUROC).

6. Model Interpretation

To interpret the network predictions, we also produce heatmaps to visualize the areas of the image most indicative of the disease using class activation mappings (CAMs) (Zhou et al., 2016). To generate the CAMs, we feed an image into the fully trained network and extract the feature maps that are output by the final convolutional layer. Let f_k be the k th feature map and let $w_{c,k}$ be the weight in the final classification layer for feature map k leading to pathology c . We obtain a map M_c of the most salient features used in classifying the image as having pathology c by taking the weighted sum of the feature maps using their associated weights. Formally,

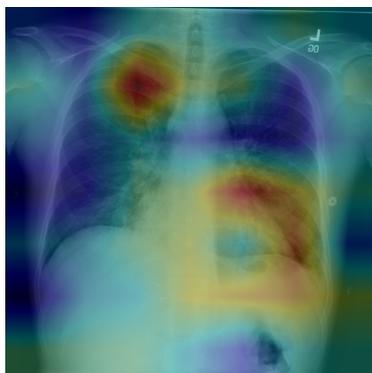
$$M_c = \sum_k w_{c,k} f_k.$$

We identify the most important features used by the model in its prediction of the pathology c by upscaling the map M_c to the dimensions of the image and overlaying the image.

Figure 3 shows several examples of CAMs on the pneumonia detection task as well as the 14-class pathology classification task.

7. Related Work

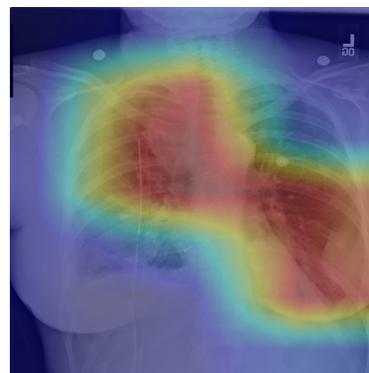
Recent advancements in deep learning and large datasets have enabled algorithms to surpass the performance of medical professionals in a wide variety of medical imaging tasks, including diabetic retinopathy



(a) Patient with multifocal community acquired pneumonia. The model correctly detects the airspace disease in the left lower and right upper lobes to arrive at the pneumonia diagnosis.



(b) Patient with a left lung nodule. The model identifies the left lower lobe lung nodule and correctly classifies the pathology.



(c) Patient with primary lung malignancy and two large masses, one in the left lower lobe and one in the right upper lobe adjacent to the mediastinum. The model correctly identifies both masses in the X-ray.



(d) Patient with a right-sided pneumothorax and chest tube. The model detects the abnormal lung to correctly predict the presence of pneumothorax (collapsed lung).



(e) Patient with a large right pleural effusion (fluid in the pleural space). The model correctly labels the effusion and focuses on the right lower chest.



(f) Patient with congestive heart failure and cardiomegaly (enlarged heart). The model correctly identifies the enlarged cardiac silhouette.

Figure 3. CheXNet localizes pathologies it identifies using Class Activation Maps, which highlight the areas of the X-ray that are most important for making a particular pathology classification.

detection (Gulshan et al., 2016), skin cancer classification (Esteva et al., 2017), arrhythmia detection (Rajpurkar et al., 2017), and hemorrhage identification (Grewal et al., 2017).

Automated diagnosis from chest radiographs has received increasing attention with algorithms for pulmonary tuberculosis classification (Lakhani & Sundaram, 2017) and lung nodule detection (Huang et al., 2017). Islam et al. (2017) studied the performance of various convolutional architectures on different abnormalities using the publicly available OpenI dataset (Demner-Fushman et al., 2015). Wang et al. (2017) released ChestX-ray-14, an order of magnitude larger than previous datasets of its kind, and also benchmarked different convolutional neural network archi-

tectures pre-trained on ImageNet. Recently Yao et al. (2017) exploited statistical dependencies between labels in order to make more accurate predictions, outperforming Wang et al. (2017) on 13 of 14 classes.

8. Conclusion

Pneumonia accounts for a significant proportion of patient morbidity and mortality (Gonçalves-Pereira et al., 2013). Early diagnosis and treatment of pneumonia is critical to preventing complications including death (Aydogdu et al., 2010). With approximately 2 billion procedures per year, chest X-rays are the most common imaging examination tool used in practice, critical for screening, diagnosis, and management of a

variety of diseases including pneumonia (Raouf et al., 2012). However, two thirds of the global population lacks access to radiology diagnostics, according to an estimate by the World Health Organization (Mollura et al., 2010). There is a shortage of experts who can interpret X-rays, even when imaging equipment is available, leading to increased mortality from treatable diseases (Kesselman et al., 2016).

We develop an algorithm which exceeds the performance of radiologists in detecting pneumonia from frontal-view chest X-ray images. We also show that a simple extension of our algorithm to detect multiple diseases outperforms previous state of the art on ChestX-ray14, the largest publicly available chest X-ray dataset. With automation at the level of experts, we hope that this technology can improve healthcare delivery and increase access to medical imaging expertise in parts of the world where access to skilled radiologists is limited.

9. Acknowledgements

We would like to acknowledge the Stanford Center for Artificial Intelligence in Medicine and Imaging for clinical dataset infrastructure support (AIMI.stanford.edu).

References

- Aydogdu, M, Ozyilmaz, E, Aksoy, Handan, Gursel, G, and Ekim, Numan. Mortality prediction in community-acquired pneumonia requiring mechanical ventilation; values of pneumonia and intensive care unit severity scores. *Tuberk Toraks*, 58(1):25–34, 2010.
- Berbaum, K, Franken Jr, EA, and Smith, WL. The effect of comparison films upon resident interpretation of pediatric chest radiographs. *Investigative radiology*, 20(2):124–128, 1985.
- Buda, Mateusz, Maki, Atsuto, and Mazurowski, Maciej A. A systematic study of the class imbalance problem in convolutional neural networks. *arXiv preprint arXiv:1710.05381*, 2017.
- CDC, 2017. URL <https://www.cdc.gov/features/pneumonia/index.html>.
- Cherian, Thomas, Mulholland, E Kim, Carlin, John B, Ostensen, Harald, Amin, Ruhul, Campo, Margaret de, Greenberg, David, Lagos, Rosanna, Lucero, Marilla, Madhi, Shabir A, et al. Standardized interpretation of paediatric chest radiographs for the diagnosis of pneumonia in epidemiological studies. *Bulletin of the World Health Organization*, 83(5):353–359, 2005.
- Davies, H Dele, Wang, Elaine E-l, Manson, David, Babyn, Paul, and Shuckett, Bruce. Reliability of the chest radiograph in the diagnosis of lower respiratory infections in young children. *The Pediatric infectious disease journal*, 15(7):600–604, 1996.
- Demner-Fushman, Dina, Kohli, Marc D, Rosenman, Marc B, Shooshan, Sonya E, Rodriguez, Laritza, Antani, Sameer, Thoma, George R, and McDonald, Clement J. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310, 2015.
- Esteva, Andre, Kuprel, Brett, Novoa, Roberto A, Ko, Justin, Swetter, Susan M, Blau, Helen M, and Thrun, Sebastian. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, 2017.
- Franquet, T. Imaging of pneumonia: trends and algorithms. *European Respiratory Journal*, 18(1):196–208, 2001.
- Gonçalves-Pereira, João, Conceição, Catarina, and Póvoa, Pedro. Community-acquired pneumonia: identification and evaluation of nonresponders. *Therapeutic advances in infectious disease*, 1(1):5–17, 2013.
- Grewal, Monika, Srivastava, Muktabh Mayank, Kumar, Pulkit, and Varadarajan, Srikrishna. Radnet: Radiologist level accuracy using deep learning for hemorrhage detection in ct scans. *arXiv preprint arXiv:1710.04934*, 2017.
- Gulshan, Varun, Peng, Lily, Coram, Marc, Stumpe, Martin C, Wu, Derek, Narayanaswamy, Arunachalam, Venugopalan, Subhashini, Widner, Kasumi, Madams, Tom, Cuadros, Jorge, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22):2402–2410, 2016.
- Hopstaken, RM, Witbraad, T, Van Engelshoven, JMA, and Dinant, GJ. Inter-observer variation in the interpretation of chest radiographs for pneumonia in community-acquired lower respiratory tract infections. *Clinical radiology*, 59(8):743–752, 2004.
- Huang, Gao, Liu, Zhuang, Weinberger, Kilian Q, and van der Maaten, Laurens. Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993*, 2016.

- Huang, Peng, Park, Seyoun, Yan, Rongkai, Lee, Junghoon, Chu, Linda C, Lin, Cheng T, Hussien, Amira, Rathmell, Joshua, Thomas, Brett, Chen, Chen, et al. Added value of computer-aided ct image features for early lung cancer diagnosis with small pulmonary nodules: A matched case-control study. *Radiology*, pp. 162725, 2017.
- Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456, 2015.
- Islam, Mohammad Tariqul, Aowal, Md Abdul, Minhaz, Ahmed Tahseen, and Ashraf, Khalid. Abnormality detection and localization in chest x-rays using deep convolutional neural networks. *arXiv preprint arXiv:1705.09850*, 2017.
- Kesselman, Andrew, Soroosh, Garshasb, Mollura, Daniel J, and Group, RAD-AID Conference Writing. 2015 rad-aid conference on international radiology for developing countries: The evolving global radiology landscape. *Journal of the American College of Radiology*, 13(9):1139–1144, 2016.
- Kingma, Diederik and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lakhani, Paras and Sundaram, Baskaran. Deep learning at chest radiography: Automated classification of pulmonary tuberculosis by using convolutional neural networks. *Radiology*, pp. 162326, 2017.
- Mollura, Daniel J, Azene, Ezana M, Starikovskiy, Anna, Thelwell, Aduke, Iosifescu, Sarah, Kimble, Cary, Polin, Ann, Garra, Brian S, DeStigter, Kristen K, Short, Brad, et al. White paper report of the rad-aid conference on international radiology for developing countries: identifying challenges, opportunities, and strategies for imaging services in the developing world. *Journal of the American College of Radiology*, 7(7):495–500, 2010.
- Neuman, Mark I, Lee, Edward Y, Bixby, Sarah, Diperna, Stephanie, Hellinger, Jeffrey, Markowitz, Richard, Servaes, Sabah, Monuteaux, Michael C, and Shah, Samir S. Variability in the interpretation of chest radiographs for the diagnosis of pneumonia in children. *Journal of hospital medicine*, 7(4): 294–298, 2012.
- Potchen, EJ, Gard, JW, Lazar, P, Lahaie, P, and Andary, M. Effect of clinical history data on chest film interpretation-direction or distraction. In *Investigative Radiology*, volume 14, pp. 404–404, 1979.
- Rajpurkar, Pranav, Hannun, Awni Y, Haghpanahi, Masoumeh, Bourn, Codie, and Ng, Andrew Y. Cardiologist-level arrhythmia detection with convolutional neural networks. *arXiv preprint arXiv:1707.01836*, 2017.
- Raof, Suhail, Feigin, David, Sung, Arthur, Raof, Sabiha, Irugulpati, Lavanya, and Rosenow, Edward C. Interpretation of plain chest roentgenogram. *CHEST Journal*, 141(2):545–558, 2012.
- Wang, Xiaosong, Peng, Yifan, Lu, Le, Lu, Zhiyong, Bagheri, Mohammadhadi, and Summers, Ronald M. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. *arXiv preprint arXiv:1705.02315*, 2017.
- WHO. Standardization of interpretation of chest radiographs for the diagnosis of pneumonia in children. 2001.
- Yao, Li, Poblenz, Eric, Dagunts, Dmitry, Covington, Ben, Bernard, Devon, and Lyman, Kevin. Learning to diagnose from scratch by exploiting dependencies among labels. *arXiv preprint arXiv:1710.10501*, 2017.
- Zhou, Bolei, Khosla, Aditya, Lapedriza, Agata, Oliva, Aude, and Torralba, Antonio. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2921–2929, 2016.