This CVPR 2020 workshop paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version;

the final published version of the proceedings is available on IEEE Xplore.

Deep Learning for Automatic Pneumonia Detection

Tatiana Gabruseva Independent researcher tatigabru.com tatigabru@gmail.com Dmytro Poplavskiy Topcon Positioning Systems Brisbane, Queensland, Australia

Alexandr Kalinin University of Michigan Ann Arbor, MI 48109 USA, and Shenzhen Research Institute of Big Data, Shenzhen 518172, Guangdong, China

akalinin@umich.edu

Abstract

Pneumonia is the leading cause of death among young children and one of the top mortality causes worldwide. The pneumonia detection is usually performed through examine of chest X-Ray radiograph by highly-trained specialists. This process is tedious and often leads to a disagreement between radiologists. Computer-aided diagnosis systems showed the potential for improving diagnostic accuracy. In this work, we develop the computational approach for pneumonia regions detection based on single-shot detectors, squeeze-and-extinction deep convolution neural networks, augmentations and multi-task learning. The proposed approach was evaluated in the context of the Radiological Society of North America Pneumonia Detection Challenge, achieving one of the best results in the challenge.

1. Introduction

Pneumonia accounts for around 16% of all deaths of children under five years worldwide [4], being the world's leading cause of death among young children [1]. In the United States only, about 1 million adults seek care in a hospital due to pneumonia every year, and 50,000 die from this disease [1]. The pneumonia complicating recent coronavirus disease 2019 (COVID-19) is a life-threatening condition claiming thousands of lives in 2020 [10, 12, 6]. Pneumonia caused by COVID-19 is of huge global concern, with confirmed cases in 185 countries across five continents at the time of writing this paper [6].

The pneumonia detection is commonly performed

through examine of chest X-Ray radiograph (CXR) by highly-trained specialists. It usually manifests as an area or areas of increased opacity on CXR [11], the diagnosis is further confirmed through clinical history, vital signs and laboratory exams. The diagnosis of pneumonia on CXR is complicated because of a number of other conditions in the lungs, such as fluid overload, bleeding, volume loss, lung cancer, post-radiation or surgical changes. When available, comparison of CXRs of the patient taken at different time points and correlation with clinical symptoms and history is helpful in making the diagnosis. A number of factors such as positioning of the patient and depth of inspiration can alter the appearance of the CXR [17], complicating interpretation even further.

There is a known variability between radiologists in the interpretation of chest radiographs [19]. To improve the efficiency and accuracy of diagnostic services computer-aided diagnosis systems for pneumonia detection has been widely exploited in the last decade [21, 20, 26, 31, 23]. Deep learning approaches outperformed conventional machine learning methods in many computer vision and medical imaging tasks, including detection [23], classification [24] and segmentation [25]. Here, we present the solution of the Radiological Society of North America (RSNA) Pneumonia Detection Challenge for pneumonia regions detection hosted on Kaggle platform [3]. Our approach uses a single-shot detector (SSD), squeeze-and-extinction deep convolution neural networks (CNNs) [16], augmentations and multi-task learning. The algorithm automatically locates lung opacities on chest radiographs and demonstrated one of the best performance in the challenges. The source code is available at https://github.com/tatigabru/kaggle-rsna.

2. Dataset

The labelled dataset of the chest X-Ray images and patients metadata was publicly provided for the challenge by the US National Institutes of Health Clinical Center [30]. The database comprises frontal-view X-ray images from 26684 unique patients. Each image is labelled with one of three different classes from the associated radiological reports: "Normal", "No Lung Opacity / Not Normal", "Lung Opacity".

Usually, the lungs are full of air. When someone has pneumonia, the air in the lungs is replaced by other material, i.e. fluids, bacteria, immune system cells, etc. The lung opacities refers to the areas that preferentially attenuate the x-ray beam and therefore appear more opaque on CXR than they should, indicating that the lung tissue in that area is probably not healthy.

The "Normal" class contains data of healthy patients without any pathologies found (including, but not limited to pneumonia, pneumothorax, atelectasis, etc.). The "Lung Opacity" class has images with the presence of fuzzy clouds of white in the lungs, associated with pneumonia. The regions of lung opacities are labelled with bounding boxes. Any given patient can have multiple boxes if more than one area with pneumonia was detected. There are different kinds of lung opacities, some are related to pneumonia and some are not. The class "No Lung Opacity / Not Normal" illustrates data for patients with visible on CXR lung opacity regions, but without diagnosed pneumonia. Fig. 1 shows examples of CXRs for all three classes labeled with bounding boxes for unhealthy patients.

The dataset is well-balanced with the distribution of classes as in Table 1.

Class	Target	Patients
Lung Opacity	1	9555
No Lung Opacity / Not Normal	0	11821
Normal	0	8851

Table 1. Classes distribution in the dataset. Target 1 or 0 indicates weather pneumonia is diagnosed or not, respectively.

3. Evaluation

The models were evaluated using the mean average precision (mAP) at different intersection-over-union (IoU) thresholds [2]. The threshold values range from 0.4 to 0.75 with a step size of 0.05: (0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75). A predicted object is considered a "hit" if its intersection over union with a ground truth object is greater than 0.4. The average precision (AP) of a single image is calculated as the mean of the precision values at each IoU

threshold:

$$AP = \frac{1}{|thresholds|} \sum_{t} \frac{TP(t)}{TP(t) + FP(t) + FN(t)} \quad (1)$$

Lastly, the score returned by the competition metric, mAP, is the mean taken over the individual average precisions of each image in the test dataset.

4. Model

Often, the solutions in machine learning competitions are based on large and diverse ensembles, test-time augmentation, and pseudo labelling, which is not always possible and feasible in real-life applications. At test-time, we often want to minimize a memory footprint and inference time. Here, we propose a solution based on a single model, ensembled over several checkpoints and 4 folds. The model utilises an SSD RetinaNet [29] with SE-ResNext101 encoder pre-trained on ImageNet [9].

4.1. Base model

The model is based on RetinaNet [29] implementation on Pytorch [14] with the following modifications:

- 1. Images with empty boxes were added to the model and contributed to the loss calculation/optimisation (the original Pytorch RetinaNet implementation [14] ignored images with no boxes).
- 2. An extra output for small anchors was added to the CNN to handle smaller boxes.
- 3. Extra output for global image classification with one of the classes ('No Lung Opacity / Not Normal', 'Normal', 'Lung Opacity') was added to the model. This output was not used directly to classify the images, however, making the model predict the other related function improved the result.
- 4. We added dropout to the global classification output to reduce overfitting. In addition to extra regularisation, it helped to achieve the optimal classification and regression results around the same epoch.

4.2. Model training

The training dataset included data for 25684 patients and the test set had data for 1000 patients. We used a range of base models pre-trained on ImageNet dataset [9]. The models without pre-train on the ImageNet performed well on classification, but worse on regression task. The following hyper-parameters were used for all training experiments (Table 2): As the training dataset was reasonably balanced (see Table 1), there was no need for extra balancing techniques. For learning rate scheduler we used available in Pytorch ReduceLROnPlateau with the patience of 4



Figure 1. Examples of the chest X-Ray images for (a) "Normal", (b) "No Lung Opacity / Not Normal", and (c) "Lung Opacity" cases. The lung opacities regions are shown on (c) with red bounding boxes.

Parameter	Description	
Optimizer	Adam	
Initial learning rate	1e-5	
Learning rate scheduler	ReduceLROnPlateau	
Patience	4	
Image size	512 x 512	

Table 2. Common models hyper-parameters.

and learning rate decrease factor of 0.2. The losses of whole image classification, individual boxes classification and anchors regression were combined with weights and used as a total loss.

4.3. Model encoders

A number of different encoder architectures has been tested: Xception [8], NASNet-A-Mobile [32], ResNet-34, -50, -101 [13], SE-ResNext-50, -101 [16], and DualPathNet-92 [7], Inception-ResNet-v2 [28], PNASNet-5-Large [18]. To enable reasonably fast experiments and model iterations, we considered architectures with good trade-offs betweens accuracy and complexity/parameters number and hence training time [5]. In this regard, VGG nets [27] and MobileNets [15] do not provide optimal accuracy on ImageNet dataset [9], while SeNet-154 [16] and NasNet-A-Large [32] have the largest number of parameters and require the most floating-point operations. Fig. 2 shows validation loss during training for various encoders used in the RetinaNet SSD. The SE-ResNext architectures demonstrated optimal performance on this dataset with a good trade-off between accuracy and complexity [5].

4.4. Multi-task learning

Extra output for global image classification with one of the classes ('No Lung Opacity / Not Normal', 'Normal',



Figure 2. Evolution of the validation loss during training for the RetinaNet model with various encoders.

'Lung Opacity') was added to the model. The total loss was combined of this global classification output with regression loss and individual boxes classification loss.

For ablation study, we trained the RetinaNet model with SE-ResNext-101 encoder and fixed augmentations with and without global classification output. The training dynamics is shown in Fig. 3. The output of global classification was not used directly to classify the images, however, making the model predict the other related function improved the result compared to training the regression-only output of the model.

As the classification output overfits faster than the detected anchors' positions/size regression, we added a dropout for the global image classification output. Besides regularization, it helped to achieve the optimal classification and regression results around the same epoch. Various dropout probabilities have been tested. Fig. 4 shows examples of training curves for SE-ResNext-101 with differ-



Figure 3. Evolution of the validation loss during training of RetinaNet with SE-ResNext-101 encoder with (red) and without (black) multi-task learning.

ent dropouts and pre-train. Without a pre-train, the models took a much longer time to converge. RetinaNet SSD with SE-ResNext-101 encoder pre-trained in Imagenet and with dropouts of 0.5 and 0.75 for the global classification output showed the best test metrics on this dataset.



Figure 4. Evolution of the validation loss during training for different versions of RetinaNet with SE-ResNext-101 encoders.

5. Images preprocessing and augmentations

The original images were scaled to 512×512 px resolution. The 256 resolution yield degradation of the results, while the full original resolution (typically, over 2000 x 2000 px) was not practical with heavier base models. Since the original challenge dataset is not very large the following images augmentations were used to reduce overfitting:

- mild rotations (up to 6 degrees)
- · shift, scale, shear

- · horizontal flip
- for some images random level of blur, noise and gamma changes
- a limited the amount of brightness / gamma augmentations

An example of a patient X-ray scan with heavy augmentations is shown in Fig. 5.

5.1. Ablation study

To examine experimentally the effect of image augmentations, we conducted an ablation study with different augmentation sets. In the ablation study, we ran training sessions on the same model with fixed hyper-parameters and changed only the sets of image augmentations. We used the following augmentation sets:

- 1. No augmentations: after resizing and normalisation, no changes were applied to the images
- 2. Light augmentations: affine and perspective changes (scale=0.1, shear=2.5), and rotations (angle=5.0)
- 3. Heavy augmentations: random horizontal flips, affine and perspective changes (scale=0.15, shear=4.0), rotations (angle=6.0), occasional Gaussian noise, Gaussian blur, and additive noise
- 4. Heavy augmentations without rotation: heavy augmentations described above without rotations
- 5. Heavy augmentations with custom rotation: heavy augmentations described above with mild rotations of 6 deg, customised as shown in Fig. 6

The dynamics of the training with different sets of augmentations is shown in Fig. 7.

The results for all experiments are presented in Table 3.

Augmentations	Best validation mAP
no augmentations	0.246127
light augmentations	0.254429
heavy augmentations	0.250230
heavy augmentations	
custom rotation	0.255617
heavy augmentations,	
no rotation	0.260971

Table 3. Pneumonia detection mean average precision results achieved with various augmentations sets on validation.

Without enough image augmentations the model showed signs of overfitting when the validation loss stopped improving (see Fig. 7). With light and heavy augmentations, the same model showed better validation loss and mAP



Figure 5. The example of a patient chest X-ray scan with heavy augmentations and rotations.



Figure 6. The diagram illustrating custom rotation of bounding boxes.

scores. The image rotations had a measurable effect on the results, as the rotation of the bounding boxes around corners modifies the original annotated regions significantly. To reduce the impact of the rotation on bounding box sizes, instead of rotating the corners we rotated two points at each edge, at 1/3 and 2/3 edge length from the corner (8 points in total), and calculated the new bounding box as min/max of the rotated points, as illustrated in Fig. 6. We tested the same model with usual rotation, custom rotation and no rotation at all. The custom rotation improved the results, but the heavy augmentations without any rotation gave the best metrics on the validation.



Figure 7. Evolution of the validation loss during training for different sets of augmentations.

6. Postprocessing

There was a difference in train and test the labelling process of the dataset provided. The train set was labelled by a single expert, while the test set was labelled by three independent radiologists and the intersection of their labels was used for the ground truth. This would yield to a smaller labelled boxes size, especially in the complicated cases. This process can be simulated using outputs from 4 folds and/or predictions from multiple checkpoints. The 20 percentile was used instead of the mean output of anchor sizes, and then it was reduced even more, proportionally to the difference between 80 and 20 percentiles for individual models (with the scale of 1.6 optimised as a hyper-parameter).

The optimal threshold for the non-maximum suppression (NMS) algorithm was also different for the train and test sets due to different labelling process. The test set true labels were available after the challenge. The NMS thresholds had a dramatic impact on the mAP metric values. Fig. 8 shows the validation mAP metrics evolution for different training epochs and NMS thresholds. The optimal NMS thresholds on validation set varied significantly from epoch to epoch with the optimum between 0.45 and 1 depending on the model.



Figure 8. The validation mAP metric versus epochs and NMS thresholds.

The other approach is re-scaling the predicted boxes sizes for the test set to 87.5% of the original sizes to reflect the difference between test and train set labelling process. The coefficient of 87.5% was chosen to approximately match the sizes to the previous approach. These differences between the train and test sets reflect differences in the annotation process for these datasets, with a consensus of expert radiologists used as ground truth in the test sets.

7. Results

The results of detection models can change significantly between epochs and depend largely on thresholds. Therefore, it is beneficial to ensemble models from different checkpoints to achieve a more stable and reliable solution. The outputs from the same model for 4 cross-validation folds and several checkpoints were combined before applying NMS algorithms and optimizing thresholds (see the diagram of the ensemble in Fig. 9.



Figure 9. The diagram of the same model ensemble technique.

The final top results of the challenge are shown in Table 4.

Team Name	Test set, mAP
Ian Pan and Alexandre Cadrin-Chênevert	0.25475
Dmytro Poplavskiy	0.24781
Phillip Cheng	0.23908

Table 4. The final leader board results in Pneumonia detection challenge showing mAP metric calculated on the private test set.

The method described in this paper took second place in the challenge. The model was based on RetineNet SSD with Se-ResNext101 encoders pre-trained on ImageNet dataset, heavy augmentations with custom rotation as described in Section 6, multi-task learning with global classification output (see Section 5) and postprocessing as in Section 7. For the final ensemble, the outputs from the same model for 4 cross-validation folds and several checkpoints were combined before applying NMS algorithms (as shown in Fig. 9). The postprocessing with re-scaling predictions was applied to compensate for the difference between the train and test sets labelling processes.

8. Discussion

The other winner's solutions were also based on the ensemble of RetinaNet models with various inputs and encoders[22]. Remarkably, all top teams made similar discoveries regarding the differences between the training and test sets. All three teams found that lowering threshold for the NMS algorithm for the test predictions compared to the validation set improved the test set scores.

In addition, systematic size reductions of the predicted bounding boxes have been also applied by the other winning teams [22]. These difference between the train and test set reflect differences in the datasets labelling process. The train set was labelled by a single expert, while the test set was labelled by three independent radiologists and the intersection of their labels was used for the ground truth.

9. Conclusions

In this paper, we propose a simple and effective algorithm for the localization of lung opacities regions. The model was based on single-shot detector RetinaNet with Se-ResNext101 encoders, pre-trained on ImageNet dataset. The number of improvements was implemented to increase the accuracy of the model. In particular, the global classification output added to the model, heavy augmentations were applied to the data, the ensemble of 4 folds and several checkpoints was unitised to generalise the model. Ablation studies have shown the improvements by the proposed approaches for the model accuracy. This method purposely does not involve test-time augmentation and provides a good trade-off between accuracy and resources. The reported method achieved one of the best results in the challenge.

10. Acknowledgements

We thank the National Institutes for Health Clinical Center for providing the chest x-ray images used in the competition, Kaggle, Inc. for hosting the challenge. The authors thank Google Cloud Platform and Dutch internet service provider HOSTKEY B.V. (hostkey.com) for access to GPU servers and technical assistance. We also acknowledge the Radiological Society of North America, the Society of Thoracic Radiology, and Kaggle, Inc. for annotating the images and organizing the competition. The authors thank the Open Data Science community (ods.ai) for useful suggestions.

References

- [1] White paper: Top 20 pneumonia facts. www.thoracic.org/patients/patient-resources/resources/toppneumonia-facts.pdf.
- [2] Evaluation metric. www.kaggle.com/c/rsna-pneumoniadetection-challenge/overview/evaluation, 2018.
- [3] Rsna challenge. www.kaggle.com/c/rsna-pneumoniadetection-challenge/overview, 2018.
- [4] World health organization: World pneumonia day 2018. www.who.int/maternal_child_adolescent/child/worldpneumonia-day-2018/en/, 2018.
- [5] Simone Bianco, Remi Cadene, Luigi Celona, and Paolo Napoletano. Benchmark analysis of representative deep neural network architectures. *IEEE Access*, 6:64270–64277, 2018.
- [6] Johns Hopkins Coronavirus Resource Center. Covid-19 dashboard by the center for systems science and engineering (csse) at johns hopkins university (jhu). https://www.arcgis.com/apps/opsdashboard/index.html, 2020.

- [7] Yunpeng Chen, Jianan Li, Huaxin Xiao, Xiaojie Jin, Shuicheng Yan, and Jiashi Feng. Dual path networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4467–4475. Curran Associates, Inc., 2017.
- [8] François Chollet. Xception: Deep learning with depthwise separable convolutions, 2016.
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [10] Claire Duployez, Rémi Le Guern, Claire Tinez, Anne-Laure Lejeune, Laurent Robriquet, Sophie Six, Caroline Loïez, and Frédéric Wallet. Panton-valentine leukocidin–secreting staphylococcus aureus pneumonia complicating COVID-19. *Emerging Infectious Diseases*, 26(8), aug 2020.
- [11] Tomás Franquet. Imaging of community-acquired pneumonia. *Journal of Thoracic Imaging*, page 1, jul 2018.
- [12] Leiwen Fu, Bingyi Wang, Tanwei Yuan, Xiaoting Chen, Yunlong Ao, Thomas Fitzpatrick, Peiyang Li, Yiguo Zhou, Yi fan Lin, Qibin Duan, Ganfeng Luo, Song Fan, Yong Lu, Anping Feng, Yuewei Zhan, Bowen Liang, Weiping Cai, Lin Zhang, Xiangjun Du, Huachun Zou, Linghua Li, and Yuelong Shu. Clinical characteristics of coronavirus disease 2019 (COVID-19) in china: A systematic review and metaanalysis. *Journal of Infection*, apr 2020.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [14] Yann Henon. Retinanet, github repo. github.com/yhenon/pytorch-retinanet, 2018.
- [15] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications, 2017.
- [16] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [17] Barry Kelly. The chest radiograph. Ulster Med J, 2012.
- [18] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search, 2017.
- [19] Mark I. Neuman, Edward Y. Lee, Sarah Bixby, Stephanie Diperna, Jeffrey Hellinger, Richard Markowitz, Sabah Servaes, Michael C. Monuteaux, and Samir S. Shah. Variability in the interpretation of chest radiographs for the diagnosis of pneumonia in children. *Journal of Hospital Medicine*, 7(4):294–298, oct 2011.
- [20] Norliza Mohd. Noor, Omar Mohd. Rijal, Ashari Yunus, and S.A.R. Abu-Bakar. A discrimination method for the detection of pneumonia using chest radiograph. *Computerized Medical Imaging and Graphics*, 34(2):160–166, mar 2010.
- [21] Leandro Luís Galdino Oliveira, Simonne Almeida e Silva, Luiza Helena Vilela Ribeiro, Renato Maurício de Oliveira, Clarimar José Coelho, and Ana Lúcia S. S. Andrade. Computer-aided diagnosis in chest radiography for detection of childhood pneumonia. *International Journal of Medical Informatics*, 77(8):555–564, aug 2008.

- [22] Ian Pan, Alexandre Cadrin-Chênevert, and Phillip M. Cheng. Tackling the radiological society of north america pneumonia detection challenge. *American Journal of Roentgenol*ogy, 213(3):568–574, Mar. 2020.
- [23] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, Matthew P. Lungren, and Andrew Y. Ng. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. arXiv:1711.05225v1, 2017.
- [24] Alexander Rakhlin, Alexey Shvets, Vladimir Iglovikov, and Alexandr A. Kalinin. Deep convolutional neural networks for breast cancer histology image analysis. In *Lecture Notes* in *Computer Science*, pages 737–744. Springer International Publishing, 2018.
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In *Lecture Notes in Computer Science*, pages 234– 241. Springer International Publishing, 2015.
- [26] Parveen N. Ravia Shabnam and Sathik M. Mohamed. Detection of pneumonia in chest x-ray images. *Journal of X-Ray Science and Technology*, 19(4):423–428, 2011.
- [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- [28] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning, 2016.
- [29] Lin T.Y., Goyal P., Girshick R., He K., and Dollár P. Focal loss for dense object detection. *IEEE International Conference on Computer Vision*, page 2999–3007, 2017.
- [30] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R.M. Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *IEEE CVPR*, 2017.
- [31] John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and Eric Karl Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study. *PLOS Medicine*, 15(11):e1002683, nov 2018.
- [32] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition, 2017.