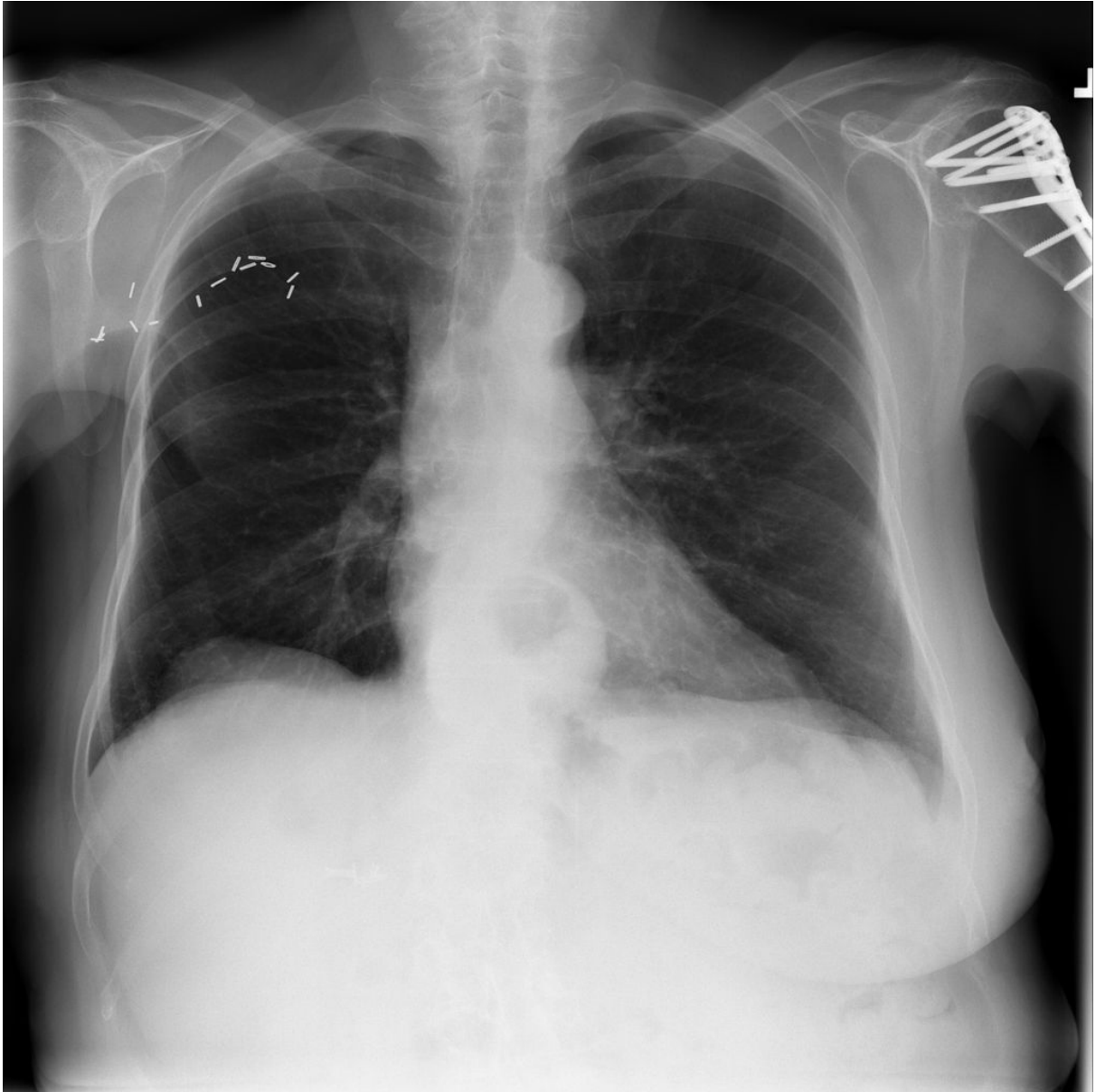# Image analysis of lung diseases



**By:**    Sander, Stephan, Joshua & Michelle

**Date:**    04/10/2020

# The Classifier

The model created by this project is a convolutional neural network (CNN), which is common practice when analyzing visual data. The first four scripts are part of preprocessing steps whilst the last does the training of the actual model. The first script, filter_diseases.py, filters the data by selecting the data belonging to three diseases. We chose for the diseases effusion, cardiomegaly and pneumothorax due to the fact that effusion and pneumothorax look alike on the scans, and cardiomegaly for being more divergent. The second script, resize_images, crops and resizes images with the goal of discarding irrelevant content. The third script, reconcile_labels.py, utilizes the csv file to "tag" the resized images with their corresponding diseases. The fourth script, image_to_array.py, very straightforwardly creates a numpy array from these images. The last script splits the data in testing and training sets, creates a model using the training set and tests the model on the test set. A table of these different modules and their description can be found below in **Table 1**.

*Table 1:* *Overview of scripts used in this project.*

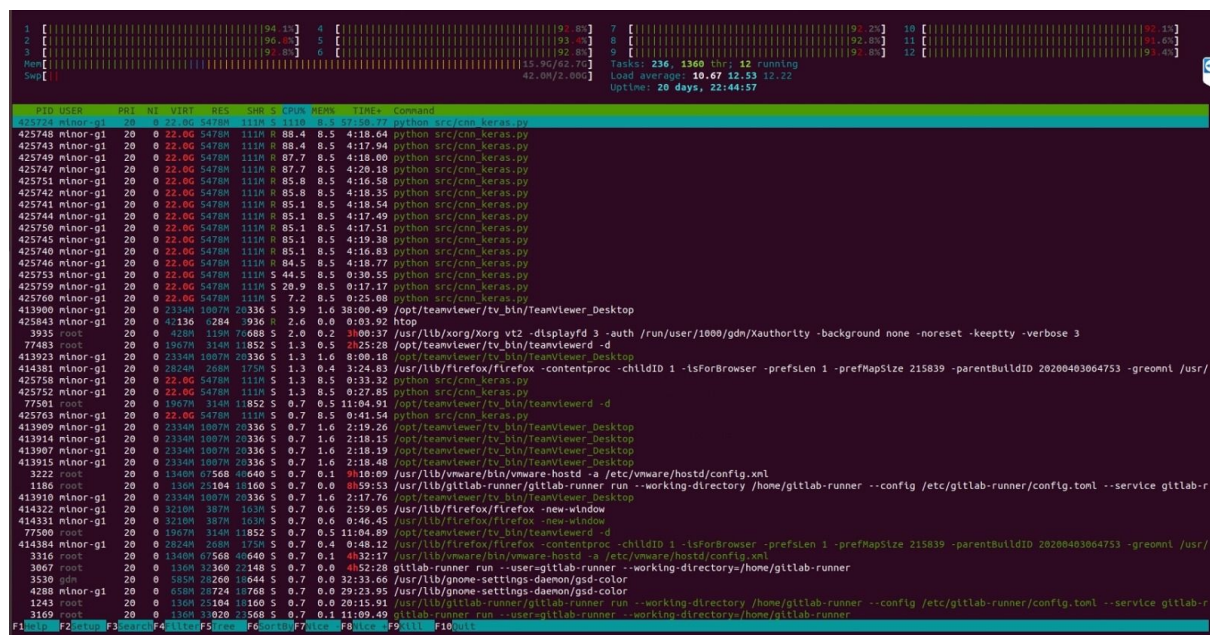| Step | Description | Script |
|------|-------------|--------|
| **1** | Split data for three diseases | filter_diseases.py |
| **2** | Resize Images | resize_images.py |
| **3** | Reconcile Labels | reconcile_labels.py |
| **4** | Convert Images to Arrays | image_to_array.py |
| **5** | CNN Model | cnn_keras.py |

# Cuda

A comparison was made between the performance of the model training when using a CPU and a GPU (**Table 2**). The runtime of the CNN training model without using GPU was just over 8 hours. To find out how much faster the prediction can be by using GPU, the GPU of Minor-g1 was used. With a runtime of 32 minutes, the speed of the prediction model was improved by approximately 15 times.

**Table 2:** *Comparison of the runtime, accuracy, precision, recall and F1 score of the CNN model between using CPU and GPU.*

|  | Runtime (hh:mm:ss) | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| **CPU** | 08:01:49 | 0.4008 | 0.16 | 0.40 | 0.23 |
| **GPU** | 00:33:23 | 0.4000 | 0.16 | 0.41 | 0.23 |

The results from the test dataset are viewed as validation of the model, considering testing on another validation set would result in very similar results. Also this helps keep the training set of substantial size for proper training of the model.

# Model Optimization



**Figure 1:** *Overview of CPU usage with htop command.*

**Figure 1** shows that the script utilizes around 1100% of the CPU when training and validating the model with tensorflow cpu. Individual cpu cores show an average usage of 91%.

***Figure 2:*** *Overview of GPU usage with nvidia-smi command.*

**Figure 2** shows the script using 3437MiB of the available 4036MiB memory from the GPU. Which corresponds to a GPU memory utilization of 85.16%.

# Discussion

We noticed that the model training using a CPU took 15 times longer than when using the GPU. This is likely due to GPU's being better at floating point and matrix calculations. While the use of one GPU already increased training performance, using multiple GPU's can certainly be more beneficial. Given the available memory of the single GPU used in this project a batch size of 32 was used. By increasing the number of GPU's, the available GPU memory can be increased significantly and consequently a greater batch size can be used.

A dataset containing the images of three diseases were used (effusion, cardiomegaly and pneumothorax) resulting in scores described in **Table 2**. The scores resulting from the CPU and GPU trained models are very similar. The GPU did perform better when it comes to the running time (8h v.s. 0.5h). However, the difference in recall score (0.40 v.s. 0.41), that is the proportion of actual positives that are identified correctly, is minimal and may be a consequence of chance. The use of a smaller dataset may contribute to overall lower scores. In addition, the fact that at least 10% of the images are incorrectly annotated further reduces these final scores. To improve on the model accuracy the incorrectly annotated images could be identified and the provided bounding boxes could be used.

# Data Availability

The code used in this project can be found at [https://github.com/MichelledeGroot/img_analysis](https://github.com/MichelledeGroot/img_analysis). Furthermore, the chest x-ray images used to train and validate our model were made publicly available by the National Institutes of Health (NIH) on their website; [https://nihcc.app.box.com/v/ChestXray-NIHCC](https://nihcc.app.box.com/v/ChestXray-NIHCC).


# Acknowledgements

Parts of Greg Chase's code ([https://github.com/gregwchase/nih-chest-xray](https://github.com/gregwchase/nih-chest-xray)) were used in this project.