

# Interpreting Covid-19 Prediction Models using Information Bottleneck

## General Info

Project Title: Interpreting Covid-19 Prediction Models using Information Bottleneck

Contact Person: Ashkan Khakzar, Dr. Seong Tae Kim

Contact Email: [ashkan.khakzar@tum.de](mailto:ashkan.khakzar@tum.de), [seongtae.kim@tum.de](mailto:seongtae.kim@tum.de)

## Project Abstract

In order to establish trust in the clinical routine for the use of neural network models for screening Covid-19, it is vital to know how the models reach their predictions. In this project we investigate what parts of the input image are responsible for each network prediction using a state of the art feature-attribution method<sup>1</sup>. Validating that the model is focusing on regions where the pathologies are located reinforces the trust over these models. The feature-attribution method that we use here has the advantage that highlights input regions based on the amount of information contained in that area for the prediction of the model. In this project we use the NIH Chest X-ray<sup>2</sup> dataset for pretraining the network and later use the Covid dataset<sup>3</sup> for finetuning the model.

## Background

There are couple of papers which present methods for automatic detection of COVID-19 disease<sup>45</sup> or quantifying the COVID-19 severity score<sup>6</sup> from X-ray images.

A huge body of works in neural network interpretability are devoted to understanding what parts of the input are responsible for the predicted output. These methods are known as input feature attribution methods. A well-known method within this category is Class Activation Maps (CAM)<sup>7</sup>. Recently a new feature attribution method<sup>1</sup> grounded on information theory is proposed that can result in more detailed maps, and the highlighted areas are based on the information the model uses for the prediction.

## Technical Prerequisites

- Basic understanding of convolutional neural networks
- Very good skills in Python
- Good skills in PyTorch,

---

<sup>1</sup> Schulz, Karl, et al. "Restricting the flow: Information bottlenecks for attribution." arXiv preprint arXiv:2001.00396 (2020).

<sup>2</sup> Wang, Xiaosong, et al. "Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.

<sup>3</sup> <https://brixia.github.io/#get-the-data>

<sup>4</sup> Oh, Yujin, Sangjoon Park, and Jong Chul Ye. "Deep learning covid-19 features on cxr using limited training data sets." IEEE Transactions on Medical Imaging (2020).

<sup>5</sup> Ozturk, Tulin, et al. "Automated detection of COVID-19 cases using deep neural networks with X-ray images." Computers in Biology and Medicine (2020): 103792.

<sup>6</sup> Signoroni, Alberto, et al. "End-to-end learning for semiquantitative rating of COVID-19 severity on Chest X-rays." arXiv preprint arXiv:2006.04603 (2020).

<sup>7</sup> Zhou, Bolei, et al. "Learning deep features for discriminative localization." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.

**Benefits:**

- Working on a Covid19 related project!
- Possible publication (based on the results)

**Work-packages and Time-plan:**

	Description	#Students	From	To
<b>WP1</b>	Familiarize with the literature. (ChestX ray classification, and feature attribution) and the datasets	all	01.11	
<b>WP2</b>	Familiarize with PyTorch. Come up with a detailed time-plan (Gantt)	all		
<b>WP3</b>	Train a model on NIH ChestX-ray dataset	Group1		
<b>WP4</b>	Fine-tune model on Covid Dataset	Group1		
<b>WP5</b>	Setup Information Bottleneck method	Group2		
<b>M1</b>	Intermediate Presentation II	all	17.12.2020	
<b>WP6</b>	Run Information Bottleneck on NIH ChestXray and Covid datasets	all		
<b>WP7</b>	Share results with clinical partners	all		
<b>WP+</b>	(Optional) Run the framework on private CT dataset	all		
<b>WP8</b>	Documentation	all		
<b>M2</b>	Final Presentation	all	26.02.2021	

