# Machine Learning Automatically Detects COVID-19 using Chest CTs in a Large Multicenter Cohort

**Authors:**

Bogdan Georgescu, PhD (8), Shikha Chaganti, PhD (8), Gorka Bastarrika Aleman, MD (1), Eduardo Jose Mortani Barbosa Jr., MD (6), Jordi Broncano Cabrero, MD (2), Guillaume Chabin (9), Thomas Flohr, apl. Prof., PhD (7), Philippe Grenier, Prof., MD (3), Sasa Grbic, PhD (8), Nakul Gupta, MD (4), François Mellot, MD (3), Savvas Nicolaou, MD (11), Thomas Re, MD (8), Pina Sanelli, MD (5), Alexander W. Sauter, MD (10), Youngjin Yoo, PhD (8), Valentin Ziebandt (7), Dorin Comaniciu, PhD (8)

(1) Clínica Universidad de Navarra, Navarra, Spain
(2) Health Time, Jaén, Spain
(3) Hôpital Foch, Suresnes, France
(4) Houston Methodist, Houston, USA
(5) Donald and Barbara Zucker School of Medicine, Feinstein Institutes for Medical Research, Northwell Health, Manhasset, NY, USA
(6) Penn Medicine, Philadelphia, PA, USA
(7) Siemens Healthineers, Forchheim, Germany
(8) Siemens Healthineers, Princeton, NJ, USA
(9) Siemens Healthineers, Paris, France
(10) University Hospital Basel, Clinic of Radiology & Nuclear medicine, Basel, Switzerland
(11) Vancouver General Hospital, Vancouver, Canada

## Summary

Automated classification and triaging of chest computed tomography (CT) scans to distinguish COVID-19 from other lung diseases or normal scans can augment radiologist diagnostic accuracy and efficiency. This manuscript evaluates the performance of multiple machine learning methods (metric-based and deep learning) to classify chest CTs in suspected COVID-19 patients.

## Key Points

- Unsupervised clustering reveals the distribution of key tomographic features such as the percent of airspace opacity, percent of ground-glass opacities, percent of high opacity or consolidation, percent of peripheral and basilar opacities across COVID-19 and control groups.

- Test data comprises of COVID-19 positive and COVID-19 negative groups, the latter including a balanced distribution of non-COVID-19 pneumonia, interstitial lung disease (ILD), and normal chest CT scans. Classification accuracies for COVID-19, pneumonia, ILD, and normal CT scans are respectively 86%, 67%, 80%, and 100%, reiterating that there can be an overlap of imaging patterns associated with COVID-19 and non-COVID-19 pneumonia, as previously reported in the clinical literature.

- The deep learning (DL) -based classification method shows an overall performance of 0.90 AUC with a sensitivity of 86% and specificity of 81%. Machine learning methods applied to quantitative chest CT metrics can therefore be valuable to improve diagnostic accuracy in suspected COVID-19.

## List of Abbreviations

**COVID-19** Coronavirus Disease 2019
**SARS-CoV-2** Severe Acute Respiratory Syndrome Coronavirus 2

**RT-PCR** Reverse Transcript Polymerase Chain Reaction

**CT** Computed Tomography

**GGO** Ground Glass Opacity

**DL** Deep Learning

**ILD** Interstitial Lung Disease

**GBT** Gradient Boosted Trees

**LR** Logistic Regression

**RF** Random Forest

**PO** Percent of Opacity

**PHO** Percent of High Opacity

**CO-RADS** COVID-19 Reporting and Data System

**PACS** Picture Archiving and Communication System

**ROC** Receiver Operating Characteristic

**AUC** Area Under the Curve

**DICOM** Digital Imaging and Communications in Medicine

# ABSTRACT

## Purpose

To investigate if AI-based classifiers can distinguish COVID-19 from other pulmonary diseases and normal groups, using chest CT images. To study the interpretability of discriminative features for COVID-19 detection.

## Materials and Methods

Our database consists of 2096 CT exams that include CTs from 1150 COVID-19 patients. Training was performed on 1000 COVID-19, 131 ILD, 113 other pneumonias, 559 normal CTs, and testing on 100 COVID-19, 30 ILD, 30 other pneumonias, and 34 normal CTs. A metric-based approach for classification of COVID-19 used interpretable features, relying on logistic regression and random forests. A deep learning-based classifier differentiated COVID-19 based on 3D features extracted directly from CT intensities and from the probability distribution of airspace opacities.

## Results

Most discriminative features of COVID-19 are percentage of airspace opacity, ground glass opacities, consolidations, and peripheral and basal opacities, which coincide with the typical characterization of COVID-19 in the literature. Unsupervised hierarchical clustering compares the distribution of these features across COVID-19 and control cohorts. The metrics-based classifier achieved AUC, sensitivity, and specificity of respectively 0.85, 0.81, and 0.77. The DL-based classifier achieved AUC, sensitivity, and specificity of respectively 0.90, 0.86, and 0.81. Most of ambiguity comes from non-COVID-19 pneumonia with manifestations that overlap with COVID-19, as well as COVID-19 cases in early stages.

## Conclusion

A new method discriminates COVID-19 from other types of pneumonia, ILD, and normal, using quantitative patterns from chest CT. Our models balance interpretability of results and classification performance, and therefore may be useful to expedite and improve diagnosis of COVID-19.

# INTRODUCTION

Coronavirus disease 2019 or COVID-19 has caused a global pandemic associated with an immense human toll and health care burden across the world (1). COVID-19 can manifest as pneumonia, which may lead to acute hypoxemic respiratory failure, which is the main reason for hospitalization and mortality. A consensus statement provided by the Fleischner Society indicates the use of lung imaging for triage of patients with moderate to severe clinical symptoms, especially in resource-constrained environments (2). The most typical pulmonary CT imaging features related to COVID-19 are multi-focal (often bilateral and peripheral predominant) airspace opacities, comprised by ground glass opacities and/or consolidation, which may be associated with interlobular and intralobular septal thickening ("crazy-paving") (3). A study comparing the differences between COVID-19 and other types of viral pneumonia demonstrated that distinguishing features more typical of COVID-19 are predominance of ground glass opacities, peripheral distribution, and vascular thickening (4). A consensus statement on the reporting of COVID-19 by Radiological Society of North America (RSNA) indicates the typical appearance of COVID-19 as peripheral and bilateral distribution of ground glass opacities with or without consolidation or a crazy paving pattern, and possibly with the 'reverse halo' sign (5). Confirmatory diagnosis of COVID-19 requires identification of the virus on nasopharyngeal swabs via RT-PCR (reverse transcription - polymerase chain reaction), a test that is highly specific (>99%) but with sensitivity ranging from 50-80% (6,7). Given the imperfect sensitivity of RT-PCR and potential resource constraints, the role of chest CT imaging for diagnosis of COVID-19 is still under investigation.

Recently, several groups have shown that COVID-19 can be distinguished from other types of lung disease on CT with variable accuracy. Mei et al. showed that chest CT scans in patients who were positive for COVID-19 by RT-PCR testing could be distinguished from chest CT scans in patients that

tested negative with an AUC of 0.92 using machine learning and deep learning (8). While this classification is potentially valuable, it is limited by lack of details on the types and distribution of findings on negative cases. It is important to be able to distinguish COVID-19 related pulmonary disease not just from healthy subjects, but also from other types of lung diseases that are not related to COVID-19, including other infections, malignancy, ILD and COPD. This is especially important as COVID-19 can manifest similarly to other respiratory infections such as influenza, which can lead to confusion in triage and diagnosis. Bai et al. showed that an artificial intelligence system can assist radiologists to distinguish between COVID-19 and other types of pneumonia by improving their diagnostic sensitivity to 88% and specificity to 90% (9). The two cohorts compared in this study are from two different countries, therefore the generalizability of their model is limited. Similarly, some of the studies that show promising results in classification do not provide a detailed description of imaging cohorts in terms of acquisition protocols or countries from which the data is acquired (10,11). This information is important since different institutions will have varied CT acquisition protocols and different clinical indications for CT usage, which can lead to distinct patient populations.

In this manuscript, we compute CT derived quantitative imaging metrics corresponding to the typical clinical presentation of COVID-19 and evaluate the discriminative power of these metrics for the diagnosis of COVID-19. We perform unsupervised clustering of interpretable features to visualize how COVID-19 patients differ from controls. We compare the performance of metrics-based classifiers to a deep learning-based model. Our large training and test datasets are comprised of chest CTs obtained in COVID-19 confirmed patients and negative controls from North America and Europe, making this one of the first large studies to demonstrate differences in COVID –19 and non-COVID-19 imaging cohorts outside of China.

# MATERIALS AND METHODS

## Patient Selection and Imaging Data

The data used in this work has been acquired from 16 different centers in North America and Europe after anonymization and ethical review at the respective institutions. Our dataset consists of chest CT scans of 1150 patients who were positive for COVID-19, and 946 of chest CT scans of patients without COVID-19, including patients with pneumonia (n=159), interstitial lung disease (ILD) (n=177), and without any pathology on chest CT (n=610). All CT scans in the COVID-19 cohort from North America have been confirmed by an RT-PCR test. The COVID-19 cohort from Europe has been either confirmed by an RT-PCR test or diagnosed based on clinical symptoms, epidemiological exposure and radiological assessment. The pneumonia cohort consists of cases of patients with non-COVID-19 viral pneumonias, organizing pneumonia or aspiration pneumonia. The ILD cohort consists of patients with various types of ILD exhibiting ground glass opacities, reticulation, honeycombing and consolidation to different degrees. The dataset was divided into training, validation and test sets (see Table 1). Model training and selection was performed based on training and validation sets. The final performance of selected models is reported on the test dataset. Refer to Table S1 in the supplemental material for detailed breakdown of demographic and scanning information for each cohort. Note that some of the information is unavailable due to anonymization protocols of some centers.

## Metrics of Airspace Disease Severity

We computed several metrics of severity based on abnormalities known to be associated with COVID-19, as well as lung and lobar segmentation. We used a previously developed Deep Image-to-Image Network that was trained on a large cohort of healthy and abnormal cases for segmentation of lungs

and lobes (12). Next, we used a DenseUnet to identify the abnormalities related to COVID-19 such as GGO and consolidations (12). Based on these segmentations, we computed thirty severity metrics to summarize the distribution, location and extent of airspace disease in the two lungs. The complete list of metrics and their detailed description is provided in the supplementary section.

## Metric-based Analysis

### Unsupervised Feature Selection and Clustering

Mutual information was used to select the metrics of severity that are most discriminative between COVID-19 and non-COVID-19 abnormalities. The $k$ best features were incrementally selected based on an internal validation split. Based on the selected metrics, an unsupervised hierarchical cluster analysis was performed to identify clusters of images that have similar features. The pairwise Euclidean distance between two metrics was used to compute a distance matrix and the average linkage method is used for hierarchical clustering (13). The resulting clustering was visualized as a heatmap. The Python Seaborn package was used for this visualization (14).

### Supervised COVID-19 Classification

Two metrics-based classifiers were trained based on the thirty computed metrics. First, we trained a Random Forest classifier, M1, using $k$ selected features based on mutual information. Subsequently, we trained a second classifier that uses logistic regression (LR), after a feature transformation based on gradient boosted trees (GBT) (15). For training GBT, we used 2000 estimators with max depth 3 and 3 features for each split. The boosting fraction 0.8 was used for fitting the individual trees. The LR classifier, M2, was trained with L2 regularization (C=0.2). The class weights were adjusted to class frequencies for the class imbalance between COVID-19 and non-COVID-19 cases.

## Supervised Deep-learning based COVID-19 classification

A deep-learning-based 3D neural network model, M3, was trained to separate the positive class (COVID-19) vs negative class (non-COVID-19). As input, we considered a two-channel 3D tensor, with the first channel containing directly the CT Hounsfield units masked by the lung region segmentation and the second channel containing the probability map of a previously proposed opacity classifier (12). The 3D network uses anisotropic 3D kernels to balance resolution and speed and consists of deep dense blocks that gradually aggregate features down to a binary output. The network was trained end-to-end as a classification system using binary cross entropy and uses probabilistic sampling of the training data to adjust for the imbalance in the training dataset labels. A separate validation dataset was used for final model selection before the performance was measured on the testing set. The input 3D tensor size is fixed (2x128x384x384) corresponding to the lung segmentation from the CT data rescaled to a 3x1x1mm resolution. The first two blocks are anisotropic and consist of convolution (kernels 1x3x3) – batch normalization – LeakyReLU and Max-pooling (kernels 1x2x2, stride 1x2x2). The subsequent five blocks are isotropic with convolution (kernels 3x3x3) – batch normalization – LeakyReLU and Max-pooling (kernels 2x2x2, stride 2x2x2) followed by a final linear classifier with the input 144-dimensional. Figure 1 shows an overview of our 3D DL classifier.

## RESULTS

Seven features were selected by computing mutual information between the feature and the class in the training dataset of 999 COVID-19 cases and 801 controls (pneumonia, ILD and healthy). Note that one case of COVID-19 was excluded from training due to field of view issues, one pneumonia control was excluded since the z-axis resolution was less than 10 mm and another pneumonia control was excluded due to incorrect DICOM parameters and artifact issues. The features are:

1) Percent of Ground Glass Opacities

2) Percent of High Opacity (PHO2) (corresponding to consolidation)

3) Percent of Opacity (PO) (corresponding to consolidation and ground-glass opacities)

4) Percent of Opacities in the Periphery (see appendix for definition)

5) Percent of Opacities in the Rind (see appendix for definition)

6) Percent of Opacities in the Right Lower Lobe

7) Percent of Opacities in the Left Lower Lobe

These features correspond to reported typical COVID-19 characteristics in the clinical literature, i.e., multifocal ground glass opacities and consolidation with basilar and peripheral distribution of the disease (3) according to the RSNA consensus statement (5) and the CO-RADS classification system, which has been proposed to classify the likelihood of COVID-19 based on the presence and extent of these findings on a scale of 1-6 (16). Figure 2. demonstrates the hierarchical clustering of these metrics, along with the ground-truth diagnosis cohort membership (COVID-19, pneumonia, ILD and healthy) shown on the band on the left of heat map. The metric values are standardized and rescaled to a value between 0 and 1. In Figure 2(a), the clustering is performed on the entire training set of 1800 subjects. The probability of belonging to the COVID-19 class increases towards the bottom of the heat map, which corresponds to higher values of the metrics, i.e., more opacities (both GGO and consolidation), and more peripheral and basilar distribution. The middle of the heatmap shows the ambiguous region, where there is an overlap of features from different disease cohorts. Figure 2(b) shows the same clustering in the test dataset for each of the disease cohorts. While there is a cluster of COVID-19 subjects that have characteristic features, there are also many which do not show all characteristics. Moreover, some cases of pneumonia and ILD overlap with the typical features of COVID-19

The seven selected features were used to train a random forest classifier (M1). The performance of this classifier on a test dataset has an AUC of 0.80 (95% CI: [0.73, 0.86]) as shown in Figure 3. The figure

shows bootstrapped ROC and AUC values, along with their 95% confidence intervals, which were computed on 1000 samples with replacement. The sensitivity and specificity of this model are 0.74 and 0.73, respectively. The performance is improved by training a second classifier on all thirty metrics using a logistic regression model (M2). The metrics are first transformed to a higher-dimensional space using feature embedding with gradient boosted trees. This model produces an AUC of 0.85 (95% CI: [0.80, 0.90])  with a sensitivity of 0.81 and a specificity of 0.77. While the performance improves, some of the interpretability is lost since the features are transformed to a higher dimension.

Finally, our deep learning-based classifier (M3) has the best performance with an AUC of 0.90 (95% CI: [0.85, 0.94]), improving the sensitivity and specificity of the system to 0.86 and 0.81 respectively. The improvement is mostly due to a reduction of the false positives from the ILD and non-COVID 19 pneumonia categories. The optimal operating point for all the models was chosen as the point with the shortest distance from the top left corner on the ROC computed on the whole test dataset, without bootstrapping (17). The corresponding confusion matrices for the three models are shown in Table 2. Figure 4 and Figure 5 illustrate examples of correctly labeled samples by the metrics-based classifier and the DL-based classifier. Figure 4 shows typical CT images from COVID-19 patients and Figure 5 shows negative examples from ILD and non-COVID-19 pneumonia patients. Overlaid in red are the areas identified by the opacity classifier. Figure 6 illustrates examples of cases incorrectly labeled by both classifiers and Figure 7 shows cases that are incorrectly labeled by the metric-based classifier but correctly labeled by the DL classifier that uses additional texture features extracted directly from the images.

## DISCUSSION

In this research, we evaluated the ability of machine learning algorithms to distinguish between chest CTs in patients positive for COVID-19 and a control cohort comprising of chest CTs obtained to evaluate

other pneumonias, ILD and normal cases. We performed an analysis based on clinically interpretable severity metrics computed from automated segmentation of abnormal regions in a chest CT scan, as well as a black-box approach using a deep learning system. Unsupervised clustering on selected severity metrics shows that while there are dominant characteristics that can be observed in COVID-19 such as the presence of ground glass opacities as well as peripheral and basal distribution, these characteristics are not observed in all cases of COVID-19. On the other hand, some subjects with ILD and pneumonia can exhibit similar characteristics. We found that the performance of the system can be improved by mapping these metrics into a higher dimensional space prior to training a classifier, as shown by model M2 in Figure 2. The best classification accuracy is achieved by the deep learning system, which is essentially a high-dimensional, non-linear model.

The deep learning method achieves a reduced false positive and false negative rate relative to the metrics-based classifier suggesting that there might be other latent radiological representations of COVID-19 that distinguish it from interstitial lung diseases or other types of pneumonia. It would be interesting to investigate how to incorporate the common imaging features into our 3D DL classifier as prior information. The proposed AI-based method has been trained and tested on a database of 2096 CT datasets with 1150 COVID-19 patients and 946 datasets coming from other categories. We also show how our method compares to the one published by Li et al (10) and found that our method achieves a higher AUC as well as sensitivity. Further details are provided in the supplementary section.

One limitation of this study is that our training set is biased toward COVID-19 and healthy controls. This bias could have influenced the specificity for discriminating against other types of lung pathology. Another limitation is that the validation set size is relatively small, which might not capture the entire data distribution of clinical use cases for proper model selection. Among the strengths of this study are the diversity of training and testing CT scans used, which were acquired from a variety of manufacturers, institutions, and regions as shown in Table S1, ensuring that our results are robust and likely

generalizable to different environments. We included not only healthy subjects but also various types of lung pathology from ILD and pneumonia to the COVID-19 negative control group.

The system described in this paper provides clinical value in several aspects. It can be used for rapid triage of positive cases, particularly in resource constrained environments where radiologic expertise may not be immediately available, whereas RT-PCR results may take up to several hours. This system could help radiologist to prioritize interpreting CTs in patients with COVID-19 by screening out lower probability cases. In addition to rapidity and efficiency concerns, the output of our deep learning classifier is easily reproducible and replicable, mitigating inter-reader variability in manually read radiology studies. While RT-PCR will remain the reference standard for confirmatory diagnosis of COVID-19, machine learning methods applied to quantitative CT can perform with high diagnostic accuracy, increasing the value of imaging in diagnosis and management of this disease.

Furthermore, the algorithms described in this paper could potentially be integrated in a surveillance effort for COVID-19, even in unsuspected patients. All chest CT scans for pulmonary and non-pulmonary pathology (i.e. coronary artery exams, chest trauma evaluation) would be automatically assessed for evidence of COVID-19 lung disease as well as for non-COVID-19 pneumonia and referring clinicians could be alerted, allowing more rapid institution of isolation protocols. Finally, it could potentially be applied retrospectively to large numbers of chest CT exams from institutional PACS systems worldwide to uncover the origin and trace the diffuse of SARS-CoV-2 in communities prior to the implementation of widespread testing efforts.

In the future, we plan to deploy and validate the algorithm in a clinical setting and evaluate the clinical utility and diagnostic accuracy on prospective data, as well as to investigate the correlation of the proposed metrics with the clinical severity of COVID-19 and disease progression over time. COVID-19 severity can be further quantified by using features from contrast CT angiography such as detection and

measurement of acute pulmonary embolism which was reported to be associated with severe COVID-19 infections (18,19). In addition, a clinical decision models could be improved by training a classifier that incorporates other clinical data such as pulse oximetry, cell counts, liver enzymes, etc. in addition to imaging features.

# SUPPLEMENTAL INFORMATION

## Airspace Disease Severity Metrics

### Metric #1-6: Percentage of Opacity (%) or PO

The total percent volume of the lung parenchyma that is affected by the airspace disease. Computed for both lungs and for each lobe.

### Metric #7-12: Percentage of High Opacity (%) or PHO

The total percent volume of the lung parenchyma that is affected by severe disease i.e., high opacity regions including consolidation and vascular thickening. High opacity is defined as the airspace disease region with mean H.U. greater than -200. Computed for both lungs and for each lobe.

### Metric #13-18: Percentage of High Opacity (%) 2

The total percent volume of the lung parenchyma that is affected by denser airspace disease i.e., high opacity regions including consolidation. High opacity is defined as the airspace disease region with mean H.U. between -200 and 50. Computed for both lungs and for each lobe.

## Metric #19: Lung severity Score (LSS)

Sum of severity score for each of the five lobes. Based on PO for each lobe, severity score of a lobe is: **0** if lobe not affected, **1** if 1-25% is affected, **2** is 25-50% is affected, **3** is 50-75% is affected, **4** is 75-100% affected. (20)

## Metric #20: Lung High Opacity Score (LHOS)

Sum of severity score for each of the five lobes, for high opacity regions only. Based on PHO for each lobe, severity score of a lobe is: **0** if lobe not affected, **1** if 1-25% is affected, **2** is 25-50% is affected, **3** is 50-75% is affected, **4** is 75-100% affected.

## Metric #21: Lung High Opacity Score (LHOS)

Sum of severity score for each of the five lobes, for high opacity regions excluding vasculature (threshold 50 HU). Based on PHO for each lobe, severity score of a lobe is: **0** if lobe not affected, **1** if 1-25% is affected, **2** is 25-50% is affected, **3** is 50-75% is affected, **4** is 75-100% affected.

## Metric #22: Bilaterality

True if both right and left lungs are involved, false if only one of the two or none is involved.

## Metric #23: Number of Affected Lobes

Number of lobes affected by the disease.

## Metric #24: Number of Total Lesions

Number of affected regions in the lung.

### Metric #25: Number of Peripheral Lesions

Number of lesions that are in the periphery of the lung. Not including apex and mediastinal regions. See Fig S1(a). Any abnormality that intersects with the peripheral border is considered a peripheral lesion. (16)

### Metric #26: Number of Lesions in the Rind

Number of regions that are in the rind of the lung as defined in (17) (See Fig S1(b)). Any abnormality that intersects with the "rind" is considered a lesion in the rind.

### Metric #27: Number of Lesions in the Core

Number of regions that are in the core of the lung as defined in (17) (See Fig S1(b)). Any abnormality that does not intersect with the rind, is considered a core lesion.

### Metric #28: Percent of Peripheral Distribution

Given by the number of peripheral lesions divided by the number of total lesions.

### Metric #29: Percent of Peripheral Lesions

The total percent volume of the lung parenchyma that is affected by disease for peripheral lesions only.

### Metric #30: Percentage of Ground Glass Opacity

The total percent volume of the lung parenchyma that is affected by less dense airspace disease i.e., lesions which are characterized as GGO only. GGO is defined as the airspace disease region with mean H.U. less than -200.

## Comparison with literature

We compared the models in this work to those published by Li et al(10). They investigated a deep learning method to distinguish COVID-19 from community-acquired pneumonia and healthy subjects

using chest CT. Their proposed DL method is based on extracting 2D features on each CT slice followed by feature pooling across slices and a final linear classifier. While the results were promising, the distribution of data for train and test was not specified in detail in terms of scanning protocols and geography.

There are two main differences between the DL method proposed in this article and the one proposed by Li et al (10). First, our method is fundamentally based on 3D deep learning, which exploits better the 3D image context, and second, our method is using as input the location of lung regions affected by opacities, which focuses the classifier on the regions of interest.

We trained and tested on our dataset using the published code by Li et al (10) and achieved an AUC of 0.86 (95% CI: [0.81, 0.91]) as shown in Figure S2. The optimal operating point, which was selected as the point closest to the top left corner of the ROC computed on the whole test dataset, without bootstrapping, produced a sensitivity of 0.72 and specificity of 0.82. Our best model, M3, on the other hand achieved an AUC of 0.9 with a sensitivity of 0.86 and specificity of 0.81.  The confusion matrix is shown in Table S2.

## ACKNOWLEDGEMENTS

## REFERENCES

1.      JHU. Coronavirus COVID-19 Global Cases by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU) [Internet]. 2020. Available from:

https://gisanddata.maps.arcgis.com/apps/opsdashboard/index.html#/bda7594740fd402994234
67b48e9ecf6

2. Rubin GD, Ryerson CJ, Haramati LB, Sverzellati N, Kanne JP, Raoof S, et al. The Role of Chest Imaging in Patient Management during the COVID-19 Pandemic: A Multinational Consensus Statement from the Fleischner Society. Chest. 2020;

3. Bernheim A, Mei X, Huang M, Yang Y, Fayad ZA, Zhang N, et al. Chest CT findings in coronavirus disease-19 (COVID-19): relationship to duration of infection. Radiology. 2020;200463.

4. Bai HX, Hsieh B, Xiong Z, Halsey K, Choi JW, Tran TML, et al. Performance of radiologists in differentiating COVID-19 from viral pneumonia on chest CT. Radiology. 2020;200823.

5. Simpson S, Kay FU, Abbara S, Bhalla S, Chung JH, Chung M, et al. Radiological Society of North America Expert Consensus Statement on Reporting Chest CT Findings Related to COVID-19. Endorsed by the Society of Thoracic Radiology, the American College of Radiology, and RSNA. Radiol Cardiothorac Imaging. 2020;2(2):e200152.

6. Kanne JP, Little BP, Chung JH, Elicker BM, Ketai LH. Essentials for radiologists on COVID-19: an update—radiology scientific expert panel. Radiology. 2020;200527.

7. Fang Y, Zhang H, Xie J, Lin M, Ying L, Pang P, et al. Sensitivity of chest CT for COVID-19: comparison to RT-PCR. Radiology. 2020;200432.

8. Mei X, Lee H-C, Diao K, Huang M, Lin B, Liu C, et al. Artificial intelligence-enabled rapid diagnosis of COVID-19 patients. medRxiv. 2020;

9. Bai HX, Wang R, Xiong Z, Hsieh B, Chang K, Halsey K, et al. AI Augmentation of Radiologist Performance in Distinguishing COVID-19 from Pneumonia of Other Etiology on Chest CT. Radiology. 2020;201491.

10. Li L, Qin L, Xu Z, Yin Y, Wang X, Kong B, et al. Artificial intelligence distinguishes COVID-19 from community acquired pneumonia on chest CT. Radiology. 2020;200905.

11. Singh D, Kumar V, Kaur M. Classification of COVID-19 patients from chest CT images using multi-objective differential evolution–based convolutional neural networks. Eur J Clin Microbiol Infect Dis. 2020;1–11.

12. Chaganti S, Balachandran A, Chabin G, Cohen S, Flohr T, Georgescu B, et al. Quantification of tomographic patterns associated with COVID-19 from chest CT. arXiv Prepr arXiv200401279. 2020;

13. Müllner D. Modern hierarchical, agglomerative clustering algorithms. arXiv Prepr arXiv11092378. 2011;

14. Waskom M. seaborn: statistical data visualization. Python 2.7 and 3.5.

15. Friedman JH. Greedy function approximation: a gradient boosting machine. Ann Stat. 2001;1189–232.

16. Prokop M, van Everdingen W, van Rees Vellinga T, Quarles van Ufford J, Stöger L, Beenen L, et al. CO-RADS–A categorical CT assessment scheme for patients with suspected COVID-19: definition and evaluation. Radiology. 2020;201473.

17. Carter J V, Pan J, Rai SN, Galandiuk S. ROC-ing along: Evaluation and interpretation of receiver operating characteristic curves. Surgery. 2016;159(6):1638–45.

18. Grillet F, Behr J, Calame P, Aubry S, Delabrousse E. Acute pulmonary embolism associated with COVID-19 pneumonia detected by pulmonary CT angiography. Radiology. 2020;201544.

19. Lang M, Som A, Mendoza DP, Flores EJ, Reid N, Carey D, et al. Hypoxaemia related to COVID-19: vascular and perfusion abnormalities on dual-energy CT. Lancet Infect Dis. 2020;

20. Adam Bernheim. Chest CT findings in COVID-19. Radiology. 2020;19.

# TABLES

Table 1. Data-split table.

| 2 classes | 4 categories | Train | Validation | Test |
|---|---|---|---|---|
| Positive | COVID-19 | 1000 | 50 | 100 |
| Negative | Pneumonia | 113 | 16 | 30 |
| | ILD | 131 | 16 | 30 |
| | Without pathology | 559 | 17 | 34 |

Table 2. Metrics based classifier confusion matrices. The models were evaluated with 100 COVID-19 positive, 30 ILD, 30 pneumonia and 34 healthy scans. The operating point was chosen as the closest point to the top left corner on the ROC computed over the test dataset (without bootstrapping).

| | | Ground Truth | | | |
|---|---|---|---|---|---|
| | | Positive | Negative | | |
| | | COVID-19 | ILD | Pneumonia | Healthy |
| | Positive | **74** | 13 | 12 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| **Predicted (M1)** | Negative | 26 | **17** | **18** | **34** |
| **Predicted (M2)** | Positive | **81** | 12 | 10 | 0 |
| | Negative | 19 | **18** | **20** | **34** |
| **Predicted (M3)** | Positive | **86** | 6 | 12 | 0 |
| | Negative | 14 | **24** | **18** | **34** |

# FIGURES

Figure 1. Overview of the deep learning based COVID-19 classifier. Preprocessing consists of lung segmentation and opacities probability distribution computation (12) followed by a 3D deep neural network trained to distinguish between the COVID-19 class and nonCOVID-19 class.



Figure 2. Heat Map of Hierarchical Clustering. This illustrates the unsupervised hierarchical clustering of the seven metrics along with cohort membership (COVID-19, pneumonia, ILD and healthy) from the entire training set of 1800 cases. The metric values are standardized and rescaled to a value between 0 and 1.
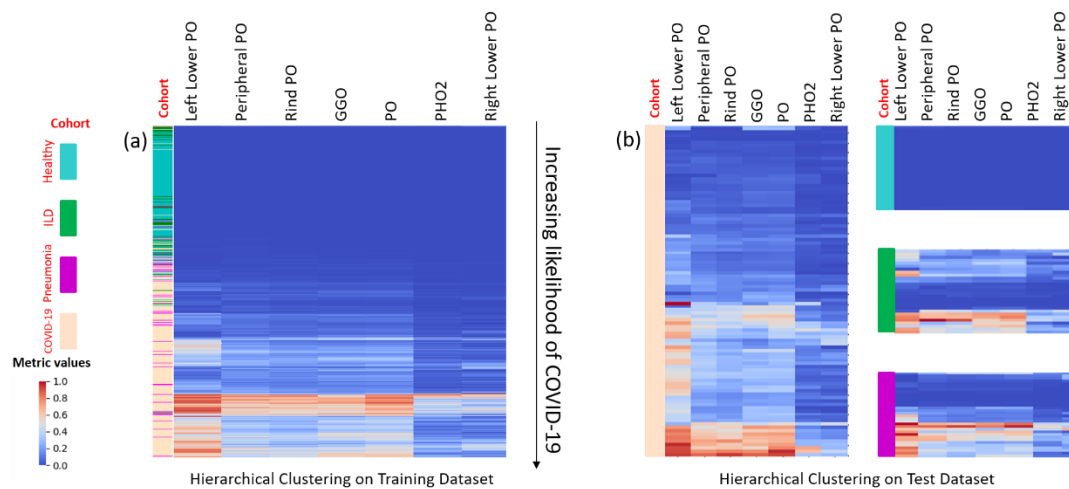
Figure 3. Bootstrapped ROCs for discriminating COVID-19 from ILD, pneumonia and healthy control by the models proposed in this study. The models were evaluated with 100 COVID-19 positive, 30 ILD, 30 pneumonia and 34 healthy scans. The 95% confidence intervals (shown as a band) are computed by bootstrapping over 1000 samples with replacement from the predicted scores.
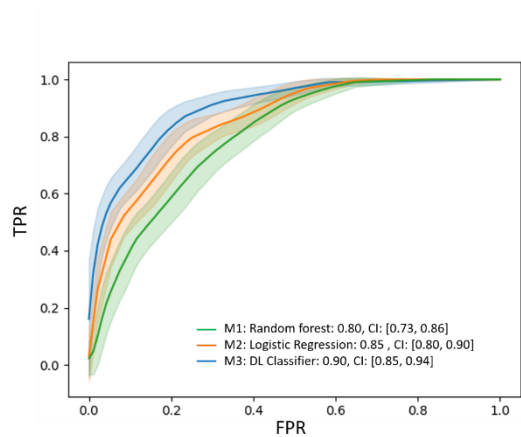
Figure 4. Examples of correctly classified COVID-19 positive samples from both methods. Red marks abnormalities associated with COVID-19.
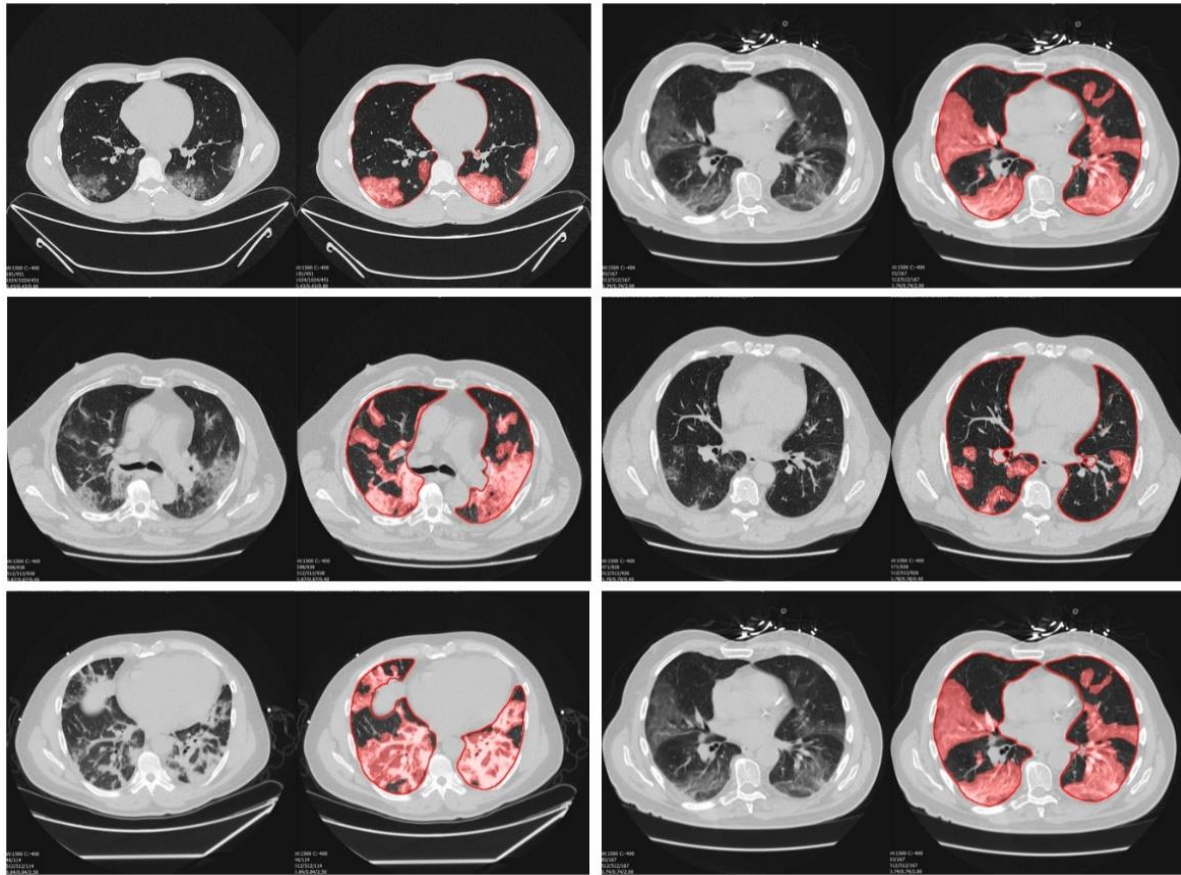


Figure 5. Examples of correctly classified negative samples from both methods – top row ILD, bottom row Pneumonia. Red marks abnormalities.
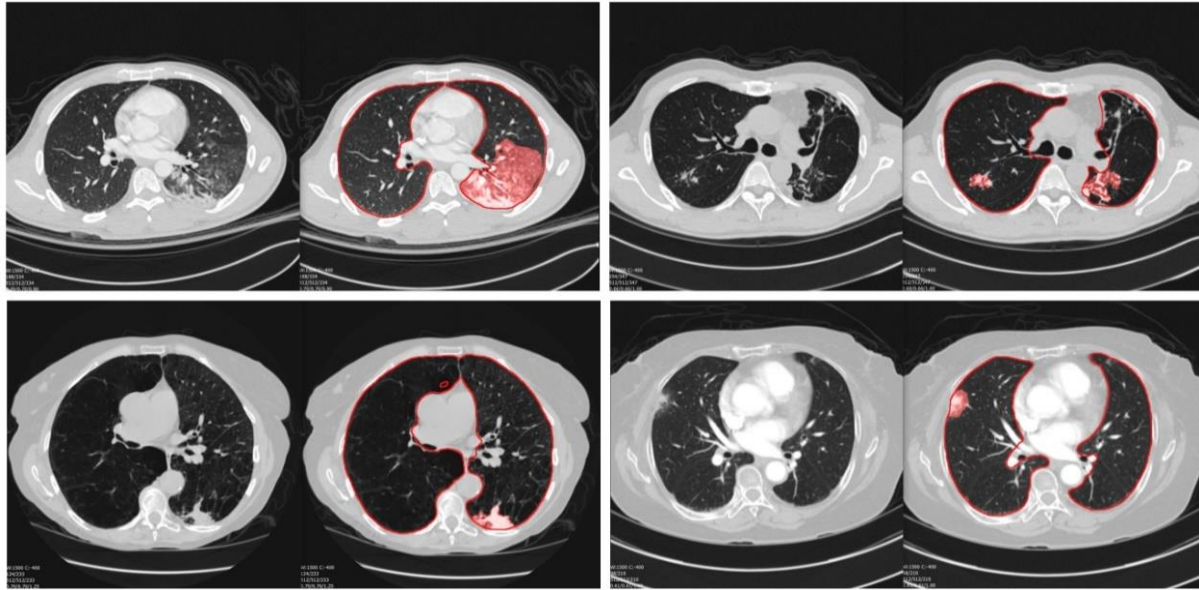
Figure 6. Examples of incorrectly classified samples by both methods: top-row COVID-19 (FN), middle-row ILD (FP), bottom-row Pneumonia (FP). Red marks abnormalities.
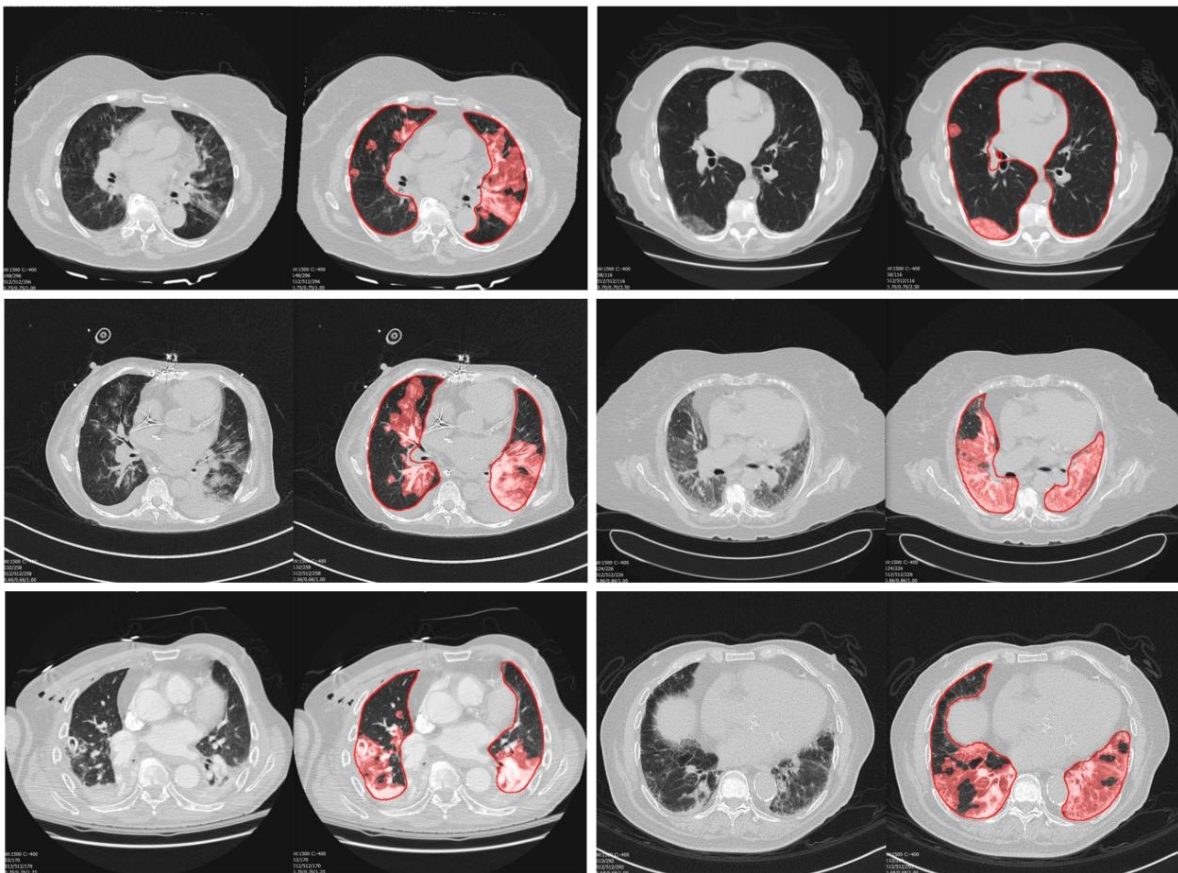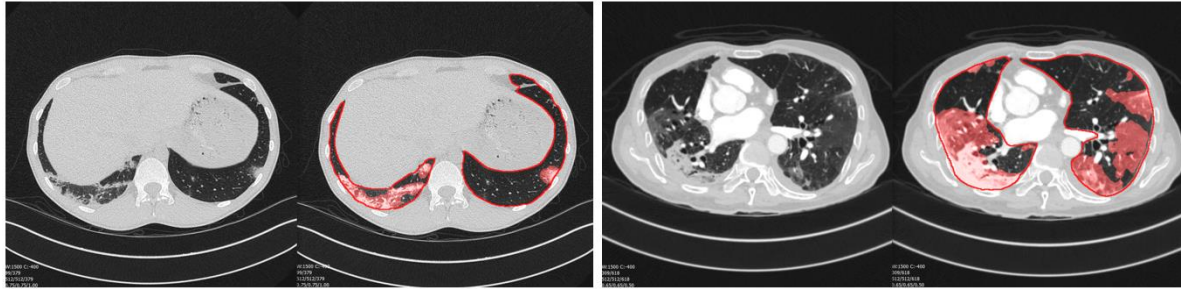
Figure 7. Examples of samples correctly classified by the DL classifier but incorrectly classified by the metric-based classifier. Red marks abnormalities.



# SUPPLEMENTARY TABLES

Table S1

| | COVID-19 | | |
|---|---|---|---|
| **# Data Sets** | Train: 1000 | Validation: 50 | Test: 100 |
| **Data Origin** | EU: 810; North America: 190 | EU:40; North America: 10 | EU:46; North America: 54 |
| **Sex** | F:232, M:345, Unknown:423 | F:15, M:18, Unknown:17 | F:27, M:57, Unknown:16 |
| **Age** | Median:61 yrs, IQR: 49-73 | Median:61 yrs, IQR: 50-74.5 | Median:49 yrs, IQR:39-66 |
| **Manufacturer** | Siemens: 777; GE: 135; NMD: 22; Philips: 10; Toshiba: 20; Unknown: 36 | Siemens: 33; GE: 9; NMD: 2; Philips: 3; Toshiba: None; Unknown: 3 | Siemens: 21; GE: 35; NMD: 11; Philips: 11; Toshiba: 16; Unknown: 17 |
| **Slice Thickness [mm]** | <1.5: 831; [1.5,3]: 124; >3: 6; unknown:39 | <1.5: 39; [1.5,3]: 11; >3: None | <1.5: 41; [1.5, 3]: 48; >3: 11 |
| **Reconstruction kernel** | Hard: 842 Soft: 97; Unknown: 61 | Hard: 39; Soft:8; Unknown: 3 | Hard:25; Soft:57; Unknown: 18 |
| | ILD | | |
| **# Data Sets** | Train: 131 | Validation: 16 | Test: 30 |

| Data Origin | North America: 131 | North America: 16 | EU:5; North America:25 |
|---|---|---|---|
| **Sex** | F:54, M:54, Unknown:23 | F:9, M:4, Unknown:3 | F:13, M:14 Unknown:3 |
| **Age** | Median:59 yrs, IQR: 56.5-73 | Median:58 yrs, IQR: 50-69.25 | Median:66 yrs IQR:53.5-75.75 |
| **Manufacturer** | Siemens: 43; GE: 55; Philips: 7; Toshiba: 5; Unknown: 21 | Siemens: 16; | Siemens: 24; GE: 1; Unknown: 5 |
| **Slice Thickness [mm]** | <1.5: 16; [1.5,3]: 54; >3: 13; Unknown: 48 | <1.5: 14; [1.5,3]: 1; >3: 1 | <1.5: 28; [1.5, 3]: 1; >3: 1 |
| **Reconstruction kernel** | Hard: 15; Soft: 73; Unknown: 43 | Hard: 0; Soft:2; Unknown: 14 | Hard: 5; Soft:0; Unknown: 25 |
| | **Pneumonia** | | |
| **# Data Sets** | Train: 113 | Validation: 16 | Test: 30 |
| **Data Origin** | EU:31; North America: 82 | North America: 16 | North America: 30 |
| **Sex** | F:53, M:53, unknown:7 | F:8, M:5, Unknown:3 | F:11, M:14, Unknown:5 |
| **Age** | Median:70 yrs, IQR: 49-73 | Median:60 yrs, IQR: 50-73.75 | Median:56 yrs IQR:41.5-66 |
| **Manufacturer** | Siemens: 70; GE: 7; Philips: 1; Toshiba: 2; Unknown:33 | Siemens: 14 , GE:1; Unknown:1 | Siemens: 25; GE: 5 |
| **Slice Thickness [mm]** | <1.5: 58; [1.5,3]: 16; >3: 39 | <1.5: 13; [1.5,3]: 1; >3: 2 | <1.5: 26; [1.5, 3]: 4; >3: 0 |
| **Reconstruction kernel** | Hard: 85; Soft: 12; Unknown:16 | Soft: 1; Unknown:15 | Hard: 1; Soft: 3; Unknown:26 |
| | **Healthy** | | |
| **# Data Sets** | Train: 559 | Validation: 17 | Test: 34 |
| **Data Origin** | North America: 559 | North America: 17 | North America: 34 |
| **Sex** | F:302, M:209, Unknown:48 | F:8, M:8, Unknown:1 | F:17, M:17 |
| **Age** | Median:57 yrs, IQR: 45-66 | Median:60 yrs, IQR: 56-64 | Median:61 yrs, IQR: 56.25-65.75 |
| **Manufacturer** | Siemens: 291; GE: 184; Philips: 16; Toshiba: 7; Unknown: 61 | Siemens: 7; GE: 6; Philips: 3; Toshiba: 1 | Siemens: 11; GE: 10; Philips: 9; Toshiba: 4 |
| **Slice Thickness [mm]** | <1.5: 515; [1.5, 3]: 43; >3: 1 | <1.5: 0; [1.5, 3]: 14; >3: 3 | <1.5: 6; [1.5, 3]: 23; >3: 5 |

| Reconstruction kernel | Hard: 149; Soft: 19; Unknown: 391 | Hard: 4 Soft: 12, Unknown: 1 | Hard: 12 Soft: 22 |
|---|---|---|---|

Table S2. Confusion matrix for the model from Li et al. The operating point was chosen as the closest point to the top left corner on the ROC computed over the test dataset (without bootstrapping).

| Li et al | | GT | | | |
|---|---|---|---|---|---|
| | | Positive | Negative | | |
| | | COVID-19 | ILD | Pneumonia | Healthy |
| Predicted | Positive | **72** | 6 | 10 | 1 |
| | Negative | 28 | **24** | **20** | **33** |

# SUPPLEMENTARY FIGURES

Figure S1 (a) Shows the periphery region. This definition excludes the apex and mediastinal border. (b) Shows the core and rind regions.
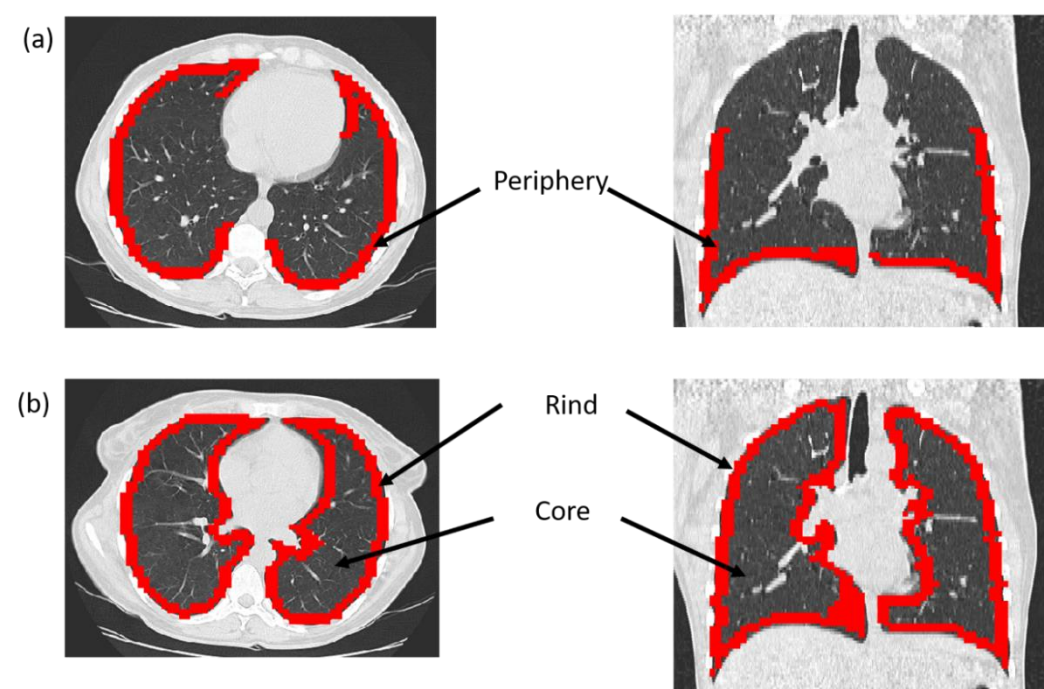
Figure S2. Bootstrapped ROCs for our 3D DL classifier and the model proposed by Li et al (10). For the model proposed by Li et al, we trained and tested on our dataset using the code provided by the authors. The 95% confidence intervals (shown as a band) are computed by bootstrapping over 1000 samples with replacement from the predicted scores.