

Comparing a deep learning model's diagnostic performance to that of radiologists to detect Covid -19 features on chest radiographs

Sabitha Krishnamoorthy, Sudhakar Ramakrishnan¹, Lanson Brijesh Colaco², Akshay Dias³,
Indu K Gopi⁴, Gautham A G Gowda⁵, Aishwarya KC⁵, Veena Ramanan⁶, Manju Chandran⁷

Department of Internal Medicine, Saroja Multispecialty Hospital, Thrissur, Kerala, India, ¹Department of Computer Science Alumni, West Virginia University, WV, USA, ²K.V.G Medical College, Sullia, Rajiv Gandhi University of Health Sciences, Bangalore, ³Department of General Medicine, Father Muller Medical College Hospital, Mangalore, Karnataka, ⁴Jubilee Centre of Medical Research, Jubilee Mission Medical College and Research Institute, Thrissur, Kerala ⁵Department of Radiodiagnosis, K.V.G Medical College and Hospital, Sullia, Karnataka, ⁶Department of Radiodiagnosis, Travancore Scans, Thiruvananthapuram, Kerala, India, ⁷Osteoporosis and Bone Metabolism Unit, Department of Endocrinology, Division of Internal Medicine, Singapore General Hospital, Singapore

Correspondences: Dr. Sabitha Krishnamoorthy, MD FACP ABIM, Department of Internal Medicine, Saroja Multispecialty Hospital, Thrissur, Kerala, India. E-mail: drsabithakrishna@gmail.com

Prof. Manju Chandran, Senior Consultant and Director, Osteoporosis and Bone Metabolism Unit Department of Endocrinology, Division of Internal Medicine, Singapore General Hospital, Singapore. E-mail: Manju.chandran@singhealth.com.sg

Abstract

Background: Whether the sensitivity of Deep Learning (DL) models to screen chest radiographs (CXR) for CoVID-19 can approximate that of radiologists, so that they can be adopted and used if real-time review of CXRs by radiologists is not possible, has not been explored before. **Objective:** To evaluate the diagnostic performance of a doctor-trained DL model (Svita_DL8) to screen for COVID-19 on CXR, and to compare the performance of the DL model with that of expert radiologists. **Materials and Methods:** We used a pre-trained convolutional neural network to develop a publicly available online DL model to evaluate CXR examinations saved in .jpeg or .png format. The initial model was subsequently curated and trained by an internist and a radiologist using 1062 chest radiographs to classify a submitted CXR as either normal, COVID-19, or a non-COVID-19 abnormal. For validation, we collected a separate set of 430 CXR examinations from numerous publicly available datasets from 10 different countries, case presentations, and two hospital repositories. These examinations were assessed for COVID-19 by the DL model and by two independent radiologists. Diagnostic performance was compared between the model and the radiologists and the correlation coefficient calculated. **Results:** For detecting COVID-19 on CXR, our DL model demonstrated sensitivity of 91.5%, specificity of 55.3%, PPV 60.9%, NPV 77.9%, accuracy 70.1%, and AUC 0.73 (95% CI: 0.86, 0.95). There was a significant correlation ($r = 0.617$, $P = 0.000$) between the results of the DL model and the radiologists' interpretations. The sensitivity of the radiologists is 96% and their overall diagnostic accuracy is 90% in this study. **Conclusions:** The DL model demonstrated high sensitivity for detecting COVID-19 on CXR. **Clinical Impact:** The doctor trained DL

This is an open access journal, and articles are distributed under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 License, which allows others to remix, tweak, and build upon the work non-commercially, as long as appropriate credit is given and the new creations are licensed under the identical terms.

For reprints contact: WKHLRPMedknow_reprints@wolterskluwer.com

Cite this article as: Krishnamoorthy S, Ramakrishnan S, Colaco LB, Dias A, Gopi IK, Gowda GA, *et al.* Comparing a deep learning model's diagnostic performance to that of radiologists to detect Covid -19 features on chest radiographs. Indian J Radiol Imaging 2021;31:S53-60.

Received: 27-Nov-2020

Revised: 09-Dec-2020

Accepted: 17-Dec-2020

Published: 23-Jan-2021

Access this article online

Quick Response Code:



Website:

www.ijri.org

DOI:

10.4103/ijri.IJRI_914_20

tool Svita_DL8 can be used in resource-constrained settings to quickly triage patients with suspected COVID-19 for further in-depth review and testing.

Key words: Artificial intelligence; COVID 19; CXR; deep learning

Introduction

It is imperative that rapid and sensitive methods to detect infection by SARS-CoV-2 be developed to meet the demands and challenges brought about by the ongoing COVID-19 pandemic. Although COVID-19 infection is typically confirmed by reverse-transcription polymerase chain reaction (RT-PCR), other non-invasive imaging methods can be employed to supplement the diagnosis.^[1] CT has demonstrated good, albeit variable, sensitivity (60 to 98%) in detecting the characteristic lung manifestations of COVID-19.^[1-3] However, CT is impractical as a screening modality given the substantial radiation exposure associated with it, lack of easy availability, and higher cost.

An additional challenge in the use of CT for initial diagnostic workup of COVID-19 is the need to establish and follow protocols for maintaining corridors for the transport of patients with suspected COVID-19 to the CT scanning room. Chest radiography (CXR) has also been used for the initial diagnosis of COVID-19 and is simpler to implement for this purpose. Its portability is an additional advantage. Moreover, the findings of COVID-19 on CXR mirror those on CT.^[4,5] Though variability may exist, the findings typically include hazy areas of increased opacities, peripheral air space and diffuse lung opacities, and bilateral lower lobe consolidations.^[4] Guidelines by the American College of Radiology support the use of portable CXR units to prevent cross-infection during the COVID-19 pandemic.^[6]

CXR reporting may experience substantial delays due to radiologist staffing issues that may be exacerbated during the pandemic. There may also be a shortage of both general and thoracic radiologists in resource-limited geographic regions. This may lead to misreporting and misdiagnosis by less trained providers.^[7] Artificial intelligence (AI) using deep learning (DL) has been applied to various imaging analyses, substantiated by the ImageNet Large Scale Visual Recognition Competition^[8] (hereafter referred to as ImageNet). DL has previously been used for detection of pulmonary pathologies such as tuberculosis^[9] and more recently has been applied to successfully identify features of COVID-19 on CT.^[10]

The few existing studies of the application of AI to identify characteristic findings of COVID-19 on CXR focused on the technical aspects of DL, were performed in experimental settings, or used single-center data sets.^[11-15] Physicians with limited technical expertise cannot be expected to readily adopt such models despite the strong need for rapid and real

time CXR interpretation in resource-limited settings.^[5,16] For widespread implementation, a DL model must not only be practical and accessible across clinical settings, but also be able to detect COVID-19 features on CXR with reasonable sensitivity not far removed from that of expert radiologists. A central role of physicians in the AI model's curation techniques can improve such models' generalizability and to facilitate use in a wide variety of clinical contexts. We thus aimed to develop and evaluate the diagnostic performance of a physician-trained DL model to screen for characteristic features of COVID-19 on CXR, comparing the performance of the DL model with that of expert radiologists.

Materials and Methods

This retrospective study was approved by the Institutional Review and the Ethics Committees (Reference Number 11/20/IRC/JMMC&RI). The requirement for informed consent was waived by the Institutional Review Committee. We compiled the CXR examinations used for this study from publicly available datasets of confirmed cases of COVID-19 from numerous countries including Australia, India, Iran, Israel, Italy, Jordan, Pakistan, Qatar, Spain, and the United Kingdom. We also acquired with permission, CXR from case presentations on COVID-19 from three states in India (Kerala, Tamil Nadu, and Karnataka), as well as from patients without COVID-19, from two partnering hospital repositories in India. These latter examinations were performed in 2017 and 2018, before the first known reported case of COVID-19 in December 2019.^[17]

Development of the DL model

1. We curated the DL model using a 42-layer deep convolutional neural network (CNN) based on the Inception-v3 network architecture. We employed transfer learning, pre-training the network initially with ImageNet data. The model was deployed as a publicly available online tool (<https://svita.in/>).
2. We transferred the knowledge of the ImageNet dataset to another dataset designated as the COVID_COLLECT_TRAINING_SET (CCTS). This served as the primary dataset for curation. CCTS contained 94535 random non-radiograph images obtained from a publicly available collection, listed in Table 1 and 1684 random radiographs that included 1062 CXR examinations [Table 1]. An internist with 15 years of experience in CXR interpretation (SK) and a general radiologist with 15 years of experience (VR) curated the CCTS. These two individuals analyzed all CXR examinations in the primary dataset and assigned them to

Table 1: COVID_COLLECT_TRAINING_SET (CCTS)

| | USE | SOURCES |
|---------|---------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Stage 1 | | |
| DS1 | Random (non-radiograph) images for training | tiny-imagenet.herokuapp.com (<i>n</i> =9435) |
| DS2 | Random radiographs for training | Repository of all types of radiographs belonging to investigators and some images from NIH_imagenet (<i>n</i> =1684) |
| Stage 2 | | |
| DS3 | Non-CXR radiographs for training | Repository of Non-CXRs belonging to investigators (<i>n</i> =391) |
| DS4 | "Normal" CXRs for training the DL model | NIH_imagenet (<i>n</i> =77) montgomery_dataset (<i>n</i> =60) thrissurnormal_dataset (<i>n</i> =113) Normal_KVG (<i>n</i> =15) Normal_NIH (<i>n</i> =15) Normal_SNH pelvis-osteoporotic 18 (<i>n</i> =17) |
| DS5 | "Abnormal" Mix includes COVID and other chest abnormalities for training the DL model | NIH_imagenet (<i>n</i> =360) COVID-19 (<i>n</i> =170) Covidspain_dataset (<i>n</i> =5) Covid19 cases_cxr (<i>n</i> =31) lee8023-covid-chestxray-dataset (<i>n</i> =97) Montgomery_dataset (<i>n</i> =18) SNH (<i>n</i> =84) |
| Stage 3 | | |
| DS6 | "COVID CXR" for training the DL model | Covidspain_dataset (<i>n</i> =5) lee8023-covid-chestxray-dataset (<i>n</i> =97) NIH_imagenet (<i>n</i> =121) Covid19 cases_cxr (<i>n</i> =31) COVID-19 (<i>n</i> =170) SNH_dataset (<i>n</i> =42) |
| DS7 | "Non-COVID Abnormal CXR" for training the DL model | NIH_imagenet (<i>n</i> =226) Montgomery_dataset (<i>n</i> =17) SNH_dataset (<i>n</i> =39) |

one of three labels: (1) normal CXR, (2) COVID-19 pattern, and (3) abnormal CXR unlikely to represent COVID-19. The images were randomly split into a test set with 10% of the data, a validation set with another 10% of the data, and a training set with the remaining 80%.

- The CXRs labeled as COVID-19 exhibited multifocal peripheral and basal lung opacities.^[5] The classical consolidation in COVID-19 is bilateral with lower lobe predominance. As the disease progresses, a diffuse distribution of lung opacities develops. CXR examinations labeled as abnormal and unlikely to represent COVID-19, exhibited central and peribronchovascular apical distribution of opacities, hilar lymphadenopathy, cavitation, calcifications, mass lesions, and pleural effusions. CXRs labeled as normal were those without any intrapulmonary opacities, bone lesions, pleural thickening associated with coarse calcifications, or pleural effusions.
- The DL model was trained to generate an output classifying the radiographs into the three labels. We introduced an expert system classifier that piped the input image through an AI pipeline decision tree. The

first stage was a fraud detection system that allowed only a valid chest radiograph to pass to subsequent steps. If the submitted image was deemed to be a valid CXR, the model evaluated if it was normal or abnormal, and if abnormal, if it was likely to represent COVID-19. At each step, the DL model generated a percentage probability of the submitted image satisfying the particular criterion. We set threshold values for the percentage probability at each step for the model to further pipe the image through the multiple stages of the decision tree [Figure S1].

- We used a different dataset for validation termed the COVID_COLLECT_REALWORLD_VALIDATION_SET (CCRVS) [Table 2]. CCRVS used three mutually exclusive groups of data from entirely different clinical resources: (1) CXR examinations from patients with confirmed COVID-19 (RT-PCR positive) [DS9, DS11] taken from case presentations by physicians from three states in India and for which the co-investigators (SK, AD, LC) validated the CXR examinations. (2) publicly available CXR images from patients with confirmed COVID-19, collected from open forums from various countries, as has been done previously^[14]; (3) CXR

Table 2: Validation study Dataset

| COVID_COLLECT_REALWORLD_VALIDATION_SET (CCRVs) | REAL WORLD IMAGES TEST | | | Total No. of images in the dataset | Images run through the DL model | Images EXCLUDED (from analysis) | Images used for analysis of DL model performance |
|--------------------------------------------------------|---------------------------------------------|-------|-----|--------------------------------------|---------------------------------|---------------------------------|--------------------------------------------------|
| Dataset compiled by investigators [DS8] | RANDOM Non-COVID CXRs [DS10] | ADFR | 60 | 60 | 18 | 42 | |
| | | ADFR2 | 79 | 79 | 12 | 67 | |
| | | ADFR3 | 32 | 32 | 4 | 28 | |
| | | NKVG | 16 | 16 | 0 | 16 | |
| | | NKVG2 | 8 | 8 | 1 | 7 | |
| | | NKVG3 | 22 | 22 | 0 | 22 | |
| | | N_SNH | 40 | 40 | 1 | 39 | |
| | CXRs of RT-PCR confirmed COVID cases [DS11] | N_NIH | 32 | 32 | 0 | 32 | |
| | | BK | 8 | 8 | 0 | 8 | |
| | | CBE | 9 | 9 | 3 | 6 | |
| | | DP | 5 | 5 | 0 | 5 | |
| | | MJRI | 2 | 2 | 0 | 2 | |
| | | OJ | 6 | 6 | 0 | 6 | |
| | | RP | 31 | 31 | 4 | 27 | |
| | | AP | 5 | 5 | 1 | 4 | |
| Covidspain_dataset [DS9] | CXRs of RT-PCR positive COVID cases | SPAIN | 129 | 129 | 10 | 119 | |
| DL Model VALIDATION STUDY on June 10, 2020 using CCRVS | | | | TOTAL IMAGES IN CCRVS | | 484 | |
| | | | | IMAGES EXCLUDED FROM CCRVS | | 54 | |
| | | | | IMAGES ANALYSED FOR STUDY FROM CCRVS | | 430 | |

images from patients without COVID-19 [DS10] from the two previously noted hospital repositories. The CXR examinations in the final group were either normal, or abnormal due to a non-COVID-19 etiology.

- The validation study had a comparative cross-sectional design and used 484 CXRs in anteroposterior or posteroanterior projection that were saved in .png or .jpeg format. The radiologists (GG, AC) visually evaluated these CXRs for overall quality and cleared them as interpretable for the study. We excluded lateral chest radiographs as well as any radiographs that either the radiologists or the DL model identified as poor-quality, for example if they were over- or underexposed or showed motion blur artifacts. These exclusions resulted in a final sample of 430 radiographs for validation.
- Multiple rounds of internal testing and optimization to maximize sensitivity were employed to determine the DL model's classification thresholds for each of the three outcomes (normal, COVID-19, and non-COVID-19 abnormality). For the purpose of this study, the threshold value for classifying a CXR as COVID-19 was set at $\geq 40\%$ probability. As an example, when we changed the thresholds for detecting COVID-19 from 40% to 60% in the internal trials, the sensitivity of the model decreased to 83%, while specificity increased to 60%.

Radiologist review

Two general radiologists with 9 and 12 years of experience (GG, AC) independently reviewed all CXR

examinations in CCRVS. These radiologists were blinded to the source of the CXR file, RT-PCR results, as well as other clinical data. The radiologists recorded for each CXR whether or not they suspected COVID-19. The radiologists also performed a subjective post-hoc assessment of imaging findings on CXRs for which the DL model provided a false positive or false negative interpretation.

Statistical analysis

Diagnostic performance of the DL model and of the two radiologists in detecting COVID-19 was assessed using the RT-PCR results as reference standard. True positives, true negatives, false positives, and false negatives were identified. The sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) were computed. The ROC curve was plotted, and the area under the curve (AUC) was calculated. The DL model's overall diagnostic accuracy was calculated. Analysis was performed using statistical software (IBM SPSS (Statistical Package of Social Sciences) version 25).

Results

Performance of the DL model

The DL model achieved sensitivity of 91.5%, specificity of 55.3%. PPV of 60.9%, and NPV of 77.9%. The overall diagnostic accuracy was 70.2%. The AUC was 0.73 (95% CI: 0.69, 0.78). Table 3 provides further classification of the model's performance.

Table 3: Statistical Classification of DL models Performance

| Svita_DL8 June 10, 2020 test run results | DL model misdiagnoses a normal CXR as COVID | | DL model misdiagnoses a COVID CXR as Normal | | DL model correctly classifies as COVID | | DL model correctly classifies as Non-COVID | | When the DL model predicts COVID positive, how often is it correct? TP/(TP + FP) | | Recall How many of the COVID cases did DL model identify? TP/(TP + FN) | | The proportion of CXRs of Actual non-COVID patients who are identified as Non-COVID by the DL model TM/(TN + FP) | | F1 Overall measure of a model's accuracy: $2 * (P * R) / (P + R)$ | | Diagnostic Accuracy (TP + FN) / (TP + FN + FP + TN) | | Reasons are mentioned under exclusion criteria in Discussion |
|------------------------------------------------|---------------------------------------------|---------------------|---------------------------------------------|--------------------|----------------------------------------|------------------|--------------------------------------------|----------|----------------------------------------------------------------------------------|------------------------|------------------------------------------------------------------------|--|------------------------------------------------------------------------------------------------------------------|--|-------------------------------------------------------------------|--|-----------------------------------------------------|--|--------------------------------------------------------------|
| | False Positive (FP) | False Negative (FN) | True Positive (TP) | True Negative (TN) | Precision (PPV) | Sensitivity (SN) | Specificity (SP) | F1 Score | Diagnostic Accuracy | Excluded from analysis | | | | | | | | | |
| ADFR | 23 | 0 | 0 | 19 | 0 | 0 | 0.452381 | 0 | 0.452381 | 18 | | | | | | | | | |
| ADFR2 | 49 | 0 | 0 | 18 | 0 | 0 | 0.268657 | 0 | 0.268657 | 12 | | | | | | | | | |
| ADFR3 | 20 | 0 | 0 | 8 | 0 | 0 | 0.285714 | 0 | 0.285714 | 4 | | | | | | | | | |
| COVID_BK | 0 | 1 | 7 | 0 | 1 | 0.875 | 0 | 0.933333 | 0.875 | 0 | | | | | | | | | |
| COVID_CBE | 0 | 0 | 6 | 0 | 1 | 1 | 0 | 1 | 1 | 3 | | | | | | | | | |
| COVID_DP | 0 | 0 | 5 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | | | | | | | | | |
| COVID_MJRI | 0 | 0 | 2 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | | | | | | | | | |
| COVID_OJ | 0 | 0 | 6 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | | | | | | | | | |
| COVID_RP | 0 | 1 | 9 | 0 | 1 | 0.9 | 0 | 0.947368 | 0.9 | 4 | | | | | | | | | |
| COVID_RP15-31 | 0 | 1 | 16 | 0 | 1 | 0.941176 | 0 | 0.969697 | 0.941176 | 0 | | | | | | | | | |
| COVID_SPAIN | 0 | 12 | 107 | 0 | 1 | 0.899160 | 0 | 0.946903 | 0.899160 | 10 | | | | | | | | | |
| COVID_AP | 0 | 0 | 4 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | | | | | | | | | |
| NKVG | 3 | 0 | 0 | 13 | 0 | 0 | 0.8125 | 0 | 0.8125 | 0 | | | | | | | | | |
| NKVG2 | 0 | 0 | 0 | 7 | 0 | 0 | 1 | 0 | 1 | 1 | | | | | | | | | |
| NKVG3 | 15 | 0 | 0 | 7 | 0 | 0 | 0.318182 | 0 | 0.318182 | 0 | | | | | | | | | |
| N_NIH | 3 | 0 | 0 | 29 | 0 | 0 | 0.90625 | 0 | 0.90625 | 0 | | | | | | | | | |
| N_SNH | 0 | 0 | 0 | 39 | 0 | 0 | 1 | 0 | 1 | 1 | | | | | | | | | |
| OVERALL | 113 | 15 | 162 | 140 | 0.589091 | 0.915254 | 0.553360 | 0.716814 | 0.702326 | 54 | | | | | | | | | |

Comparison of the DL model and radiologists

There was a statistically significant correlation ($r = 0.617, P = 0.000$) between the interpretation by the DL model and the interpretation of the CXRs by the radiologists in classifying the CXR as a suspected COVID case or not [Figure 1]. The first radiologist's sensitivity was 96%, specificity was 85.7%, PPV 82.5%, NPV 96.8%, AUC was 0.91 (95% CI: 0.88, 0.94).

The second radiologist's sensitivity was 95.5%, specificity was 85.65%, PPV was 84.47%, NPV was 96.21% and AUC 0.88 (95% CI: 0.84, 0.91). The overall diagnostic accuracy of the radiologists was 90% and 90.2%, respectively.

Figure 2 shows the distribution of findings on CXR examinations with a false-negative interpretation by the DL model for COVID-19, based on post hoc assessment by the two radiologists. The most common findings were an ARDS like picture and effusion, pneumonia, or dense infiltrates. Figure 3 shows the distribution of findings on CXR examinations with a false-positive interpretation by the DL model for COVID-19, based on post hoc assessment by the two radiologists. The most common findings were an ARDS like picture and a normal radiograph.

Discussion

In this study, we evaluated the performance of a DL model that was trained to detect COVID-19 on CXR using multiple real-world datasets curated by an internist and radiologist. Model validation was performed using CXR data sets from multiple hospitals and outpatient clinics from 10 countries and 3 different states in India. The model had high sensitivity for COVID-19 on CXR in both the training and validation data sets, with sensitivities of 92% and 91.5%, respectively.

Of prior studies^[11-15] that used DL models to detect COVID-19 on CXR, one used a class decomposition approach employing a deep CNN architecture (Detrac ResNet). This had obtained a sensitivity of 97.9%, though it was conducted experimentally on a smaller dataset.^[13] The other study that used a pre-trained CNN (ResNet 50) included only abnormal CXR examinations and showed an overall accuracy of 89.2%.^[14] The techniques presented in these two studies require machine learning expertise and are difficult to incorporate into clinical settings. Another study that explored the diagnostic performance of AI for COVID-19 in a clinical setting used only a small sample of abnormal CXR examinations from a single center.^[15]

Our DL model and study have several strengths. We not only studied the DL's model's performance characteristics, but also conducted independent analysis of 2 radiologists'

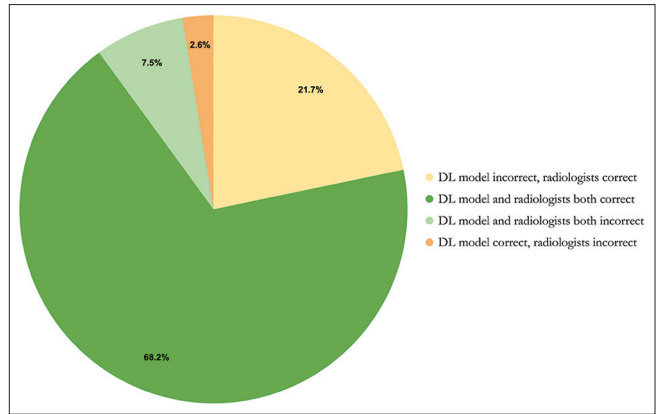


Figure 1: Comparison of deep learning model to radiologist interpretation

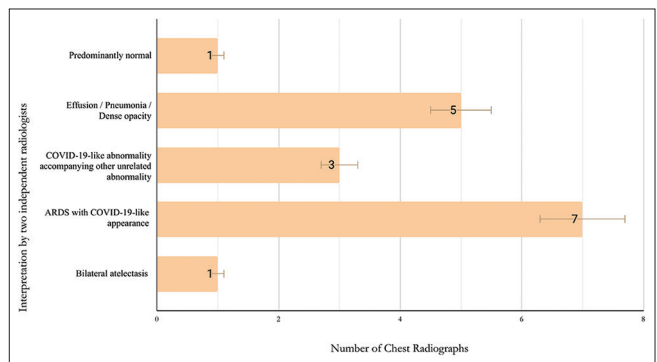


Figure 2: Distribution of false-negative CXR interpretations for COVID-19 by the deep learning model

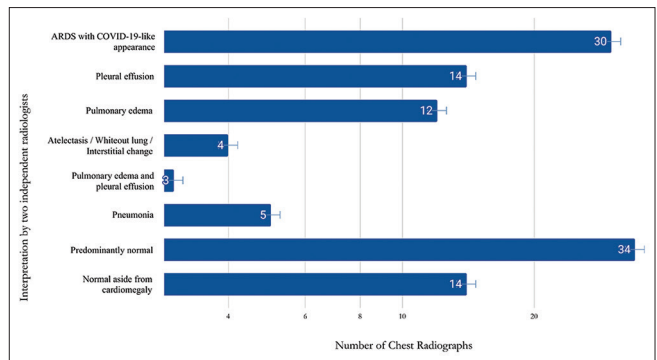


Figure 3: Distribution of false-positive CXR interpretations for COVID-19 by the deep learning model

interpretations and compared their performance characteristics with that of the model. We tested the model over a large spectrum of CXRs and demonstrated high sensitivity. This high sensitivity makes the model a potential screening tool to guide patient triage. We used a distinctive curation technique in introducing a step by step classification with an expert system overlay between the steps. The fraud detection system built into the application eliminates “junk” images from entering the model, thereby enabling the model's use in real-world clinical settings. The introduction of variable percentage thresholds at different

steps of the expert system helped to change and assess the model's performance before determining the final cut-off [Figure S1]. The expert system thresholds allow users of the model to tailor its sensitivity depending upon the pretest probability for COVID-19 in the setting in which the model is used (e.g., outpatient, inpatient, or ICU). For instance, the threshold can be changed to a higher level in an ICU setting.

The DL model was deployed as an online tool (<https://svita.in/>). Physicians can readily access the model and upload images of CXR examinations. Furthermore, our dataset was collected randomly from various sources, minimizing selection bias. An additional unique feature of the model is that we trained the model using CXR examinations saved in .jpeg or .png format. This facilitates the model being applied in clinical scenarios in which resource constraints exist. For example, the model can assess an image of a CXR taken using a camera of a smart phone on a viewing box. A unique feature of our model is that it has been configured in such a way that it can enable "fan-in" of other pathologies into the AI pipeline without having to reconstruct the system.

A limitation is the model's low specificity (55.3%). The false positives may be explained by the fact that the model is not able to differentiate subtle changes on CXRs unlike human radiologists. Therefore, any minimal increase in opacities or haziness due to other etiologies or artifacts were also labelled as COVID-19 by the model. ARDS is a sequela of multiple diseases, and it may be difficult to determine the etiology of ARDS solely by CXR.^[18] CXR examinations

that show ARDS, as well as pulmonary edema, due to causes other than COVID-19 are difficult to differentiate from those that show ARDS resulting from COVID-19.^[19] The CXR examinations from patients without COVID-19 were procured from a tertiary care center with many critically ill patients with other lung pathologies that caused an ARDS-like picture [Figure 4]. These were likely misinterpreted by the model as COVID-19.^[19] However, in general, it has to be noted that the reported specificity of radiologists in diagnosing COVID-19 from CXRs is only about 69%.^[5] While the model had low specificity, the model is not intended to serve as a standalone tool but rather as a screening tool given its excellent sensitivity. Screened patients considered likely to have COVID-19 on CXR can be triaged to a separate area where an RT-PCR swab test can be performed.

Our model aimed to detect COVID-19 characteristics on CXR and not disease extent or severity. Future research and development could supplement the model by introducing an object detection layer, and re-designing, and re-training it to categorize the disease extent.

Conclusion

Our experienced radiologists' sensitivity and diagnostic performance in detecting COVID-19 characteristics on CXRs is higher than previously reported. The specificity of the DL model is lower than that of experienced radiologists' to rule out COVID-19 on CXR. This makes it quite clear that it cannot replace expert radiologists. However, its relatively high sensitivity makes it useful as a rapid and

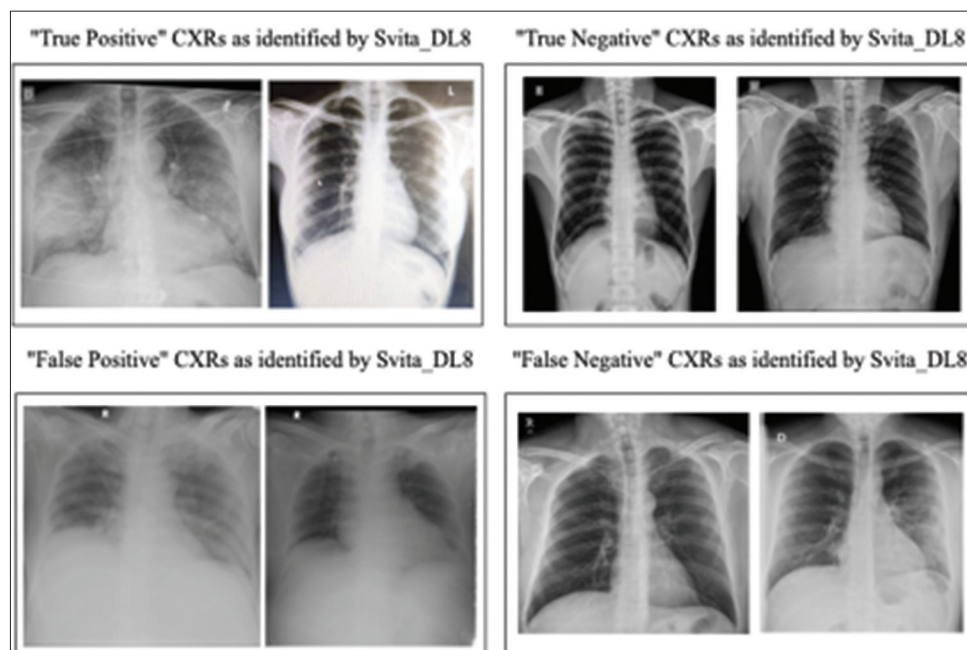


Figure 4: Examples of the interpretation of CXR by the deep learning model

real-time screening solution for COVID-19, particularly in resource-constrained environments. Such AI models can help automate some aspects of initial screening, for example, to filter in the CXRs images that need to be flagged for earlier expert radiology assessment.

Ethics committee

Acknowledgements

Mr. Eric Dee who provided copy editing support. Mr. Unnikrishnan of JMMC&RI provided assistance with the statistical analysis.

Financial support and sponsorship

Nil.

Conflicts of interest

There are no conflicts of interest.

References

- 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;296:E115-7.
- 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. *Radiology: J Thorac Imaging* 2020;2:e200152.
- Li Y, Xia L. Coronavirus disease 2019 (COVID-19): Role of chest CT in diagnosis and management. *AJR Am J Roentgenol* 2020;214:1280-6.
- Jacobi A, Chung M, Bernheim A, Eber C. Portable chest X-ray in coronavirus disease-19 (COVID-19): A pictorial review. *Clin Imaging* 2020;64:35-42.
- Wong HYF, Lam HYS, Fong AH-T, Leung ST, Chin TW-Y, Lo CSY, *et al.* Frequency and distribution of chest radiographic findings in COVID-19 positive patients. *Radiology* 2020;296:E72-8.
- ACR Recommendations for the use of Chest Radiography and Computed Tomography (CT) for Suspected COVID-19 Infection [Internet]. *Acr.org*. 2020. Available from: <https://www.acr.org/Advocacy-and-Economics/ACR-Position-Statements/Recommendations-for-Chest-Radiography-and-CT-for-Suspected-COVID19-Infection> [Last cited on 2020 May 07].
- Nolan TM, Oberklaid F, Boldt D. Radiological services in a hospital emergency department—An evaluation of service delivery and radiograph interpretation. *J Paediatr Child Health* 1984;20:109-12.
- Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, *et al.* ImageNet large scale visual recognition challenge. *Int J Comput Vis* 2015;115:211-52.
- Qin ZZ, Sander MS, Rai B, Titahong CN, Sudrungrot S, Laah SN, *et al.* Using artificial intelligence to read chest radiographs for tuberculosis detection: A multi-site evaluation of the diagnostic accuracy of three deep learning systems. *Sci Rep* 2019;9:15000.
- 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. Published online 2020 Mar 19. doi: 10.1148/radiol.2020200905
- Oh Y, Park S, Ye JC. Deep learning Covid-19 features on CXR using limited training data sets. *IEEE Trans Med Imaging* 2020;39:2688-700.
- Rajaraman S, Antani S. Weakly labeled data augmentation for deep learning: A study on COVID-19 detection in chest X-rays. *Diagnostics* 2020;10:358.
- Abbas A, Abdelsamea MM, Gaber MM. Classification of COVID-19 in chest X-ray images using DeTraC deep convolutional neural network. *Appl Intell* 2020. <https://doi.org/10.1007/s10489-020-01829-7>.
- Hall LO, Paul R, Goldgof DB, Goldgof GM. Finding covid-19 from chest x-rays using deep learning on a small dataset. *arXiv preprint arXiv: 2004.02060*. 2020. [arXiv: 2004.02060](https://arxiv.org/abs/2004.02060)
- Murphy K, Smits H, Knoops AJG, Korst MBJM, Samson T, Scholten ET, *et al.* COVID-19 on the chest radiograph: A multi-reader evaluation of an AI system. *Radiology* 2020;296:E166-72.
- Chau T-N, Lee P-O, Choi K-W, Lee C-M, Ma K-F, Tsang T-Y, *et al.* Value of initial chest radiographs for predicting clinical outcomes in patients with severe acute respiratory syndrome. *Am J Med* 2004;117:249-54.
- Timeline of WHO's response to COVID-19 [Internet]. *Who.int*. 2020. Available from: <https://www.who.int/news-room/detail/29-06-2020-covidtimeline> [Last cited on 2020 Jun 29].
- Force AD, Ranieri VM, Rubenfeld GD, Thompson BT, Ferguson ND, Caldwell E. Acute respiratory distress syndrome. *JAMA* 2012;307:2526-33.
- Gibson PG, Qin L, Puah S. COVID-19 acute respiratory distress syndrome (ARDS): Clinical features and differences from typical pre-COVID-19 ARDS. *Med J Aust* 2020;213:54-56.e1.