



OPEN ACCESS

EDITED BY

Jinming Duan,
University of Birmingham, United Kingdom

REVIEWED BY

Yuexiang Li,
Tencent Jarvis Lab, China
Fei He,
Coventry University, United Kingdom

*CORRESPONDENCE

Yalin Zheng
✉ Yalin.Zheng@liverpool.ac.uk

RECEIVED 30 November 2022

ACCEPTED 08 August 2023

PUBLISHED 23 August 2023

CITATION

Bridge J, Meng Y, Zhu W, Fitzmaurice T, McCann C, Addison C, Wang M, Merritt C, Franks S, Mackey M, Messenger S, Sun R, Zhao Y and Zheng Y (2023) Development and external validation of a mixed-effects deep learning model to diagnose COVID-19 from CT imaging.

Front. Med. 10:1113030.

doi: 10.3389/fmed.2023.1113030

COPYRIGHT

© 2023 Bridge, Meng, Zhu, Fitzmaurice, McCann, Addison, Wang, Merritt, Franks, Mackey, Messenger, Sun, Zhao and Zheng. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Development and external validation of a mixed-effects deep learning model to diagnose COVID-19 from CT imaging

Joshua Bridge¹, Yanda Meng¹, Wenyue Zhu¹, Thomas Fitzmaurice^{1,2}, Caroline McCann³, Cliff Addison⁴, Manhui Wang⁴, Cristin Merritt⁵, Stu Franks⁵, Maria Mackey⁶, Steve Messenger⁶, Renrong Sun⁷, Yitian Zhao⁸ and Yalin Zheng^{1*}

¹Institute of Life Course and Medical Sciences, University of Liverpool, Liverpool, United Kingdom, ²Department of Respiratory Medicine, Liverpool Heart and Chest Hospital NHS Foundation Trust, Liverpool, United Kingdom, ³Department of Radiology, Liverpool Heart and Chest Hospital NHS Foundation Trust, Liverpool, United Kingdom, ⁴Advanced Research Computing, University of Liverpool, Liverpool, United Kingdom, ⁵Alces Flight Limited, Bicester, United Kingdom, ⁶Amazon Web Services, London, United Kingdom, ⁷Department of Radiology, Hubei Provincial Hospital of Integrated Chinese and Western Medicine, Hubei University of Chinese Medicine, Wuhan, China, ⁸Cixi Institute of Biomedical Engineering, Ningbo Institute of Materials Technology and Engineering, Chinese Academy of Sciences, Ningbo, China

Background: The automatic analysis of medical images has the potential improve diagnostic accuracy while reducing the strain on clinicians. Current methods analyzing 3D-like imaging data, such as computerized tomography imaging, often treat each image slice as individual slices. This may not be able to appropriately model the relationship between slices.

Methods: Our proposed method utilizes a mixed-effects model within the deep learning framework to model the relationship between slices. We externally validated this method on a data set taken from a different country and compared our results against other proposed methods. We evaluated the discrimination, calibration, and clinical usefulness of our model using a range of measures. Finally, we carried out a sensitivity analysis to demonstrate our methods robustness to noise and missing data.

Results: In the external geographic validation set our model showed excellent performance with an AUROC of 0.930 (95%CI: 0.914, 0.947), with a sensitivity and specificity, PPV, and NPV of 0.778 (0.720, 0.828), 0.882 (0.853, 0.908), 0.744 (0.686, 0.797), and 0.900 (0.872, 0.924) at the 0.5 probability cut-off point. Our model also maintained good calibration in the external validation dataset, while other methods showed poor calibration.

Conclusion: Deep learning can reduce stress on healthcare systems by automatically screening CT imaging for COVID-19. Our method showed improved generalizability in external validation compared to previous published methods. However, deep learning models must be robustly assessed using various performance measures and externally validated in each setting. In addition, best practice guidelines for developing and reporting predictive models are vital for the safe adoption of such models.

KEYWORDS

CT, COVID-19, deep learning, diagnosis, imaging

1. Background

Coronavirus disease 2019 (COVID-19) is an infectious respiratory disease caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). Virus clinical presentation ranges from mild cold-like symptoms to severe viral pneumonia, which can be fatal (1). While some countries have achieved relative control through lockdowns, future outbreaks and new strains are expected to continue, with many experts believing the virus is here to stay (2). Detection and isolation is the most effective way to prevent further spread of the virus. Even with effective vaccines becoming widely available, with the threat of continued waves and new potentially vaccine-resistant variants, it is vital to further develop diagnostic tools for COVID-19. These tools will likely also apply to future outbreaks of other similar diseases as well as common diseases such as pneumonia.

The diagnosis of COVID-19 is usually determined by Reverse Transcription Polymerase Chain Reaction (RT-PCR), but this is far from being a gold standard. A negative test does not necessarily indicate a negative diagnosis, with one recent review finding that RT-PCR has a real-world sensitivity of around 70% and a specificity of 95% (3). Furthermore, an individual patient data systematic review (4) found that RT-PCR often fails to detect COVID-19, and early sampling is key to reducing false negatives. Therefore, these tests are often more helpful to rule in COVID-19 rather than ruling out. If a patient presents with symptoms of COVID-19, but an RT-PCR test is negative, then further tests are often required (1). Consecutive negative tests with at least a one-day gap are recommended; however, this still does not guarantee that the patient is negative for COVID-19 (5). Computed tomography (CT) can play a significant role in diagnosing COVID-19 (6). Given the excessive number of COVID-19 cases worldwide and the strain on resources expected, automated diagnosis might reduce the burden on reporting radiologists.

CT images are made up of many slices, creating a three dimensional (3D)-like structure. Previous approaches, such as those used by Li et al. (7) and Bai et al. (8) treat the image as separate slices and use a pooling layer to concatenate the slices. An alternative approach assumes the slices form a 3D structure and use a 3D CNN, such as that proposed in CoviNet (9). A fundamental limitation of these methods is the need for the same number of slices as their inputs, but the number of slices often varies between different CT volumes. Instead, we propose using a novel mixed-effects layer to consider the relationship between slices in each scan. Mixed-effects models are commonly used in traditional statistics (10, 11), but we believe this is the first time that mixed-effects models have been utilized in such a way. It has been observed that some lobes of the lung are more often affected by COVID-19 than others (12, 13) with lower lobe distribution being a prominent feature of COVID-19 (14), the fixed-effects take this into account by considering where each slice is located within the scan.

Deep learning has shown great potential in the automatic classification of disease, often achieving expert-level performance. Such models could screen and monitor COVID-19 by automatically analyzing routinely collected CT images. As observed by Wynants et al. (15) and Roberts et al. (16) many models are already developed to diagnose COVID-19, which often obtain excellent discriminative performance; however, very few of these models, if any, are suitable for clinical use, mainly due to a lack of robust analysis and reporting. These models often suffer from common pitfalls, making them

unsuitable for broader adoption. Roberts et al. (16) identified three common areas in which models often fail these are: (1) a lack of adequately documented methods for reproducibility, (2) failure to follow established guidelines and best practices for the development of deep learning models, and (3) an absence of external validation displaying the model's applicability to a broader range of data outside of the study sample. Failure to address these pitfalls leads to profoundly flawed and biased models, making them unsuitable for deployment.

In this work, we aim to address the problems associated with previous models by following guidelines for the reporting (17, 18) and development (19) of prediction models to ensure that we have rigorous documentation allowing the methods developed here to be replicated. In addition, we will make code and the trained model publicly available at: github.com/JTBridge/ME-COVID19 to promote reproducible research and facilitate adoption. Finally, we use a second dataset from a country other than the development dataset to externally validate the model and report a range of performance measures evaluating the model's discrimination, calibration, and clinical usefulness.

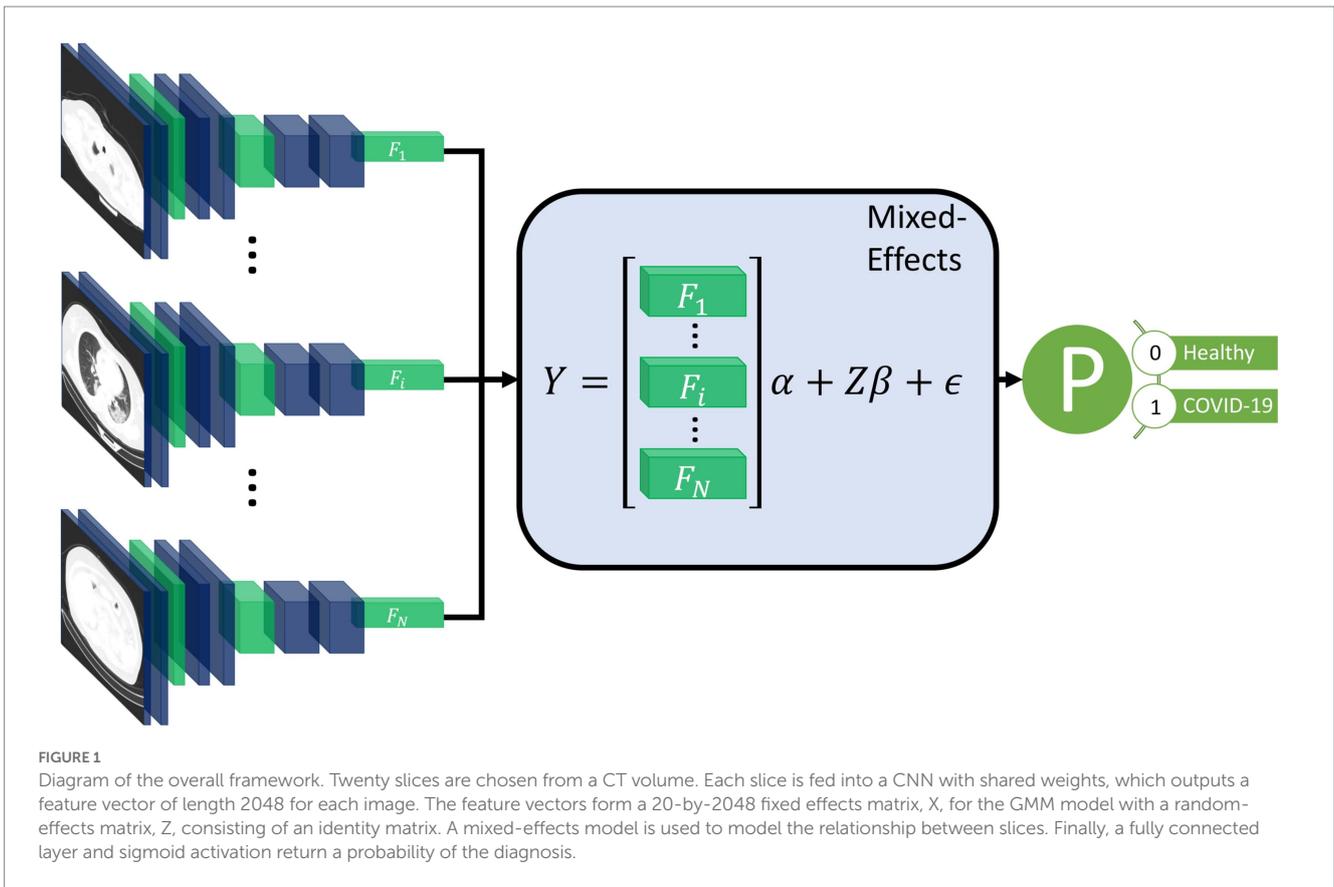
Hence, our main aim is to develop a mixed-effects deep learning model to accurately classify images as healthy or COVID-19, following best practice guidelines. Our secondary aim is to show how deep learning predictive algorithms can satisfy current best practice guidelines to create reproducible and less biased models.

2. Methods

Our proposed method consists of a feature extractor and a two-stage generalized linear mixed-effects model (GLMM) (20), with all parameters estimated within the deep learning framework using backpropagation. First, a series of CT slices forming a CT volume is input to the model. In our work, we use 20 slices. Next, a convolutional neural network (CNN) extracts relevant features from the model and creates a feature vector for each CT slice. Then, a mixed-effects layer concatenates the feature vectors into a single vector. Finally, a fully connected layer followed by a sigmoid activation gives a probability of COVID-19 for the whole volume. The mixed effects and fully connected layer with sigmoid activation are analogous to a linear GLMM in traditional statistics. The overall framework is shown in Figure 1.

2.1. Feature extractor

For the feature extractor, we use a CNN. In this work, we chose InceptionV3 (21) as it is relatively efficient and commonly used. InceptionV3 outputs a feature vector of length 2048. To reduce the time needed to reach convergence, we pretrained the CNN on ImageNet (22). A CNN is used for each slice, with shared weights between CNNs; this reduces the memory footprint of the model. Following the CNN, we used a global average pooling layer to reduce each image to a feature vector for each slice. We then added a dropout of 0.6 to improve generalizability to unseen images. We form the feature vectors into a matrix of shape $20 \times 2,048$. Although we used InceptionV3 (21) here, other networks would also work and may provide better performance on other similar tasks. We then need to concatenate these feature vectors into a single feature vector for the



whole volume; normally, pooling is used, in our work we propose using a mixed-effects models.

2.2. Mixed-effects network

We propose utilizing a novel mixed-effects layer to model the relationship between slices. Mixed-effects models are a statistical model consisting of a fixed-effects part and a random-effects part. The fixed-effects part models the relationship within the CT slice; the random effects can model the spatial correlation between CT slices within the same image (11). For volumetric data, the number of slices may differ significantly due to various factors such as imaging protocol and the size of the patient. Some CT volumes in the dataset may have fewer images than the model is designed to use, which leads to missing data. The number of slices depends upon many factors including the scanning protocol and the size of the patient. Mixed-effects models can deal with missing data provided the data are missing at random (23). It is likely that the data here is missing at random, although not completely at random. The mixed-effects model is given by

$$Y_i = X_i\alpha + Z_i\beta + e_i$$

where Y_i, X_i, Z_i, e_i are vectors of outcomes, fixed effects design matrix of shape $slices \times features$, random effects design matrix of shape $slices \times slices$, and vector of error unknown random errors of the i th patient of shape $slices$, respectively, and α, β are fixed and random effects parameters, both of length $features$ and $slices$, respectively. In our work, we have 20 slices and 2048 features and use

the identity matrix for the random effects design matrix. The values in the random effects design matrix can be changed to reflect a non-uniform distance between slices. We assume that the random effects β are normally distributed with mean 0 and variance G

$$\beta \sim N(0, G).$$

We also assume independence between the random effects and the error term.

The fixed effects design matrix, X , is made up of the feature vectors output from the feature extraction network. For the random effects design matrix, Z , we use an identity matrix with the same size as the number of slices; in our experiments, this is 20. The design matrix is then given by

$$Z_{20 \times 20} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}.$$

This matrix easily generalizes to any number of slices. If the distance between slices is not uniform, the values can be altered accordingly. We assumed no particular correlation matrix. We included the fixed and random intercept in the model. All parameters for the mixed-effects layer were initialized using the Gaussian distribution with mean 0 and standard deviation 0.05.

A type of mixed-effects modeling has previously been combined with deep learning for gaze estimation (24). However, their mixed-effects method is very different from our proposed method; they used the same design matrix for fixed and random effects. In addition, they also estimated random-effects parameters with an expectation-maximization algorithm, which was separate from the fixed effects estimation, which used deep learning. In our work, we utilize a spatial design matrix to model the spatial relationship between slices and estimate parameters within the deep learning framework using backpropagation without the need for multiple stages.

2.3. Loss function

As the parameters in the model are all estimated using backpropagation, we must ensure that the assumption of normally distributed random effects parameters with mean zero is valid. We achieve this by introducing a loss function for the random effects parameters, which enforces a mean, skewness, and excess kurtosis of 0. We measure skewness using the adjusted Fisher-Pearson standardized moment coefficient

$$Skew(\beta) = \frac{\sqrt{n(n-1)}}{n-2} \frac{E[(\beta - \bar{\beta})^3]}{(E[(\beta - \bar{\beta})^2])^{3/2}}$$

and the excess kurtosis using

$$Kurt(\beta) - 3 = \frac{1}{n^2} \sum_{i=1}^n \left(\frac{E[(\beta - \bar{\beta})^4]}{(E[(\beta - \bar{\beta})^2])^2} - 3 \right),$$

where n is the length of β , $\bar{\beta}$ is the mean of β and $E[\cdot]$ is the expectation function. The Gaussian distribution has a kurtosis of 3; therefore, the excess kurtosis is given by the kurtosis minus 3. This formula for this random-effects parameters loss function which we aim to minimize, is then given by

$$L_{random} = |E(\beta)| + |Skew(\beta)| + |Kurt(\beta) - 3|.$$

For the classification, we use the Brier Score (25) as the loss function, which is given by

$$L_{Brier} = \frac{1}{N} \sum_{i=1}^N (p_i - o_i)^2$$

where N is the total number of samples, p_i is the predicted probability of sample i and o_i is the observed outcome of sample i . The Brier score is the same as the mean squared error of the predicted probability.

We chose to use the Brier Score over the more commonly used binary cross-entropy because it can be decomposed into two

components: refinement and calibration. Calibration is often overlooked in deep learning models but is vital to assess the safety of any prediction model. The refinement component combines the model's resolution and uncertainty and measures the model's discrimination. The calibration component can be used as a measure of the model calibration. Therefore, the Brier Score can be used to optimize both the discrimination and calibration of the model. The overall loss function is given by

$$L = L_{Brier} + L_{random}.$$

A scaling factor could be introduced to weight one part of the loss function as more important than the other; however, we give both parts of the loss function equal weighting in our work.

We also transformed the labels as suggested by Platt (26) to reduce overfitting. The negative and positive labels become

$$o_- = \frac{1}{N_- + 2}$$

and

$$o_+ = \frac{N_+ + 1}{N_+ + 2}$$

respectively, where N_- and N_+ are the total number of negative and positive cases in the training set. This is similar to label smoothing as commonly used in deep learning, but the new targets are chosen by applying Bayes' Rule to the out-of-sample data to prevent overfitting.

2.4. Classification layer

The output of the mixed-effects layer is a single vector, which is the same length as the number of slices used. For example, in our work, we had a vector of length 20. Furthermore, we used a fully connected layer with sigmoid activation to obtain a probability of the scan showing COVID-19; the sigmoid activation is analogous to the logistic link function in traditional statistics. Finally, we added an L1 regularization term of 0.1 and an L2 regularization term of 0.01 to the kernel to reduce overfitting.

2.5. Model performance

Many deep learning models focus on assessing discriminative performance only, using measures such as the area under the receiver operating characteristic curve (AUROC), sensitivity, and specificity. To better understand the model performance and impact, we report performance measures in three broad areas: discrimination, calibration, and clinical usefulness (27). Discrimination assesses how well a model can discriminate between healthy and COVID-19 positive patients. Models with excellent discriminative performance can still produce unreliable results, with vastly overestimated probabilities regardless of the true diagnosis (28). Model calibration is often overlooked and rarely reported in deep learning, if at all; however, poorly calibrated models can be misleading and lead to dangerous clinical decisions (28). Calibration can be assessed using

four levels, with each level indicating better calibration than the last (29). The fourth and most stringent level (strong calibration) requires the correct model to be known, which in turn requires predictors to be non-continuous, and an infinite amount of data to be used and is therefore considered utopic. We consider the third level (moderate calibration) using calibration curves. Moderate calibration will ensure that the model is at least not clinically harmful. Finally, measures of clinical usefulness assess the clinical consequences of the decision and acknowledge that a false positive may be more or less severe than a false negative.

Firstly, the discriminative performance is assessed using AUROC using the pROC package in R (30), with confidence intervals constructed using DeLong's (31) method. For sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV), we use the epiR (32) package in R (30); with 95% confidence intervals constructed using Jeffrey's prior (33). We report performance at a range of probability thresholds to demonstrate how the thresholds can be adjusted to reduce false positives or false negatives depending on the setting (34). Secondly, we assess model calibration using calibration curves created using the CalibrationCurves package (29), which is based on the rms (35) package. Finally, we assess the clinical usefulness of the model using decision curve analysis (36). Net benefits are given at various thresholds, and models which reach zero net benefit at higher thresholds are considered more clinically useful. Two brief sensitivity analyses are performed, one assessing the model's ability to deal with missing data and the other assessing its ability to deal with noise. To improve the model's interpretability and reduce the black-box nature, we produce saliency maps (37) that show which areas of the image are helpful to the model in the prediction. We also check the assumption of normally distributed random-effects parameters.

2.6. Comparison models

To assess the added benefit of using our mixed-effects method, we compare against networks that use alternative methods. Both COVNet (7) and a method proposed by Bai et al. (8) propose deep learning models that consider the slices separately before concatenating the features using max pooling. COVNet uses a ResNet50 (38) CNN to extract features and pooling layers to concatenate the features before a fully connected classification layer. The model proposed by Bai et al. uses EfficientNetB4 (39) to extract features followed by a series of full-connected layers with batch normalization and dropout; average pooling is then used to concatenate the feature vectors before classification. While max pooling is simple and computationally efficient, it cannot deal with pose variance and does not model the relationship between slices.

An alternative method to pooling is treating the scans as 3D, such as in CoviNet (40). CoviNet takes the whole scan and uses a 16-layer 3D CNN followed by pooling and fully connected layers. We implemented these models as described in their respective papers.

In all comparison experiments, we kept hyperparameters, such as learning rate, learning rate decay, and data augmentation, the same to ensure the comparisons were fair. For COVNet (7) and the model proposed by Bai et al. (8) we pretrained the CNNs on ImageNet as

they also did; however, no pretrained models were available for CoviNet. For the loss function, we also used the Brier score (25).

2.7. Computing

Models were developed using an Amazon Web Services p3.xlarge node with four Tesla V100 16GiB GPUs and 244GiB available memory. Model inference was performed on a local Linux machine running Ubuntu 18.04, with a Titan X 12GiB GPU and 32GiB available memory. Model development and inference were performed using Tensorflow 2.4 (41, 42), and R 4.0.5 (30) was used to produce evaluation metrics (43, 44) and graphs (35, 45). We used mixed precision to reduce the computational cost, which uses 16-bit floating-point precision in all layers, except for the mixed-effects and classification layers, where 32-bit floating-point precision is used.

We used the Adam optimizer (46) with an initial learning rate of $1e-4$; if the internal validation loss did not improve for three epochs, we reduced the learning rate to 20%. In addition, we assumed convergence and stopped training if the loss did not improve for 10 epochs to reduce the time spent training and the energy used.

2.8. Data

There is currently no established method for estimating the sample size estimate in deep learning. We propose treating the final fully connected classification layer as the model and treating previous layers as feature extraction. We can then use the number of parameters in the final layer to estimate the required sample size. Using the "pmsampsize" package (47) in R, we estimate the required minimum sample size in the development set. We use a conservative expected C-statistic of 0.8, with 21 parameters and an estimated disease prevalence of 80% based on datasets used in other studies. This gives a minimum required sample size of 923 patients in the training set. For model validation, around 200 patients with the disease and 200 patients without the disease are estimated to be needed to assess calibration (29).

All data used here is retrospectively collected and contains hospital patients with CT scans performed during the COVID-19 pandemic. The diagnosis was determined by examining radiological features of the CT scan for signs of COVID-19, such as ground-glass opacities. For model development, we use the MosMed dataset (48), which consists of a total of 1,110 CT scans displaying either healthy or COVID-19 infected lungs. The scans were performed in Moscow hospitals between March 1, 2020, and April 25, 2020. We split the dataset into two sets for training and internal validation on the patient level. The training set is used to train the model, and the internal validation set is used to select the best model based on the loss at each epoch; this helps prevent overfitting on the training set. In addition, we obtained images from a publicly available dataset published by Zhang et al. (49) consisting of CT images from a consortium of Chinese hospitals.

Overall, this allows us to perform external geographical validation in another country and to better evaluate the developed model. In addition, we will be able to assess how well a deep learning model generalizes to other populations. A summary of all the datasets used

is shown in Table 1. We have 923 patients in the training set and at least 200 patients in each class for the external validation set.

2.9. Patient and public involvement

Patients or the public were not involved in the design, conduct, reporting, or dissemination of our research.

2.10. Data pre-processing and augmentation

The MosMed dataset was converted from Dicom image format into PNG, normalized to have a mean of 120 and a variance of 95. Images were ordered from the top of the lungs to the bottom. During training, we applied random online data augmentation to the images. This alters the image slightly and gives the effect of increasing the training dataset size, although this is not as good as expanding the training dataset with more samples. First, we adjusted the brightness and contrast between 80 and 120%. We then rotated the image plus or minus 5 degrees and cropped the image up to 20% on each side. Finally, we flipped the image horizontally and vertically with a probability of 50% each. All random values were

chosen using the uniform distribution except for the flips, which were chosen using a random bit. Example images are shown in Figure 2A.

The dataset taken from Zhang et al. (49) required a large amount of sorting to be made suitable for use. Some of the scans were pre-segmented and only showed the lung areas, while others showed the whole CT scan. We removed any pre-segmented images. Identifying information on some images had to be cropped to reduce bias in the algorithm. In addition, many of the scans were duplicates but were not labeled as such, and many scans were incomplete, only showing a few lung slices or not showing any lung tissue at all. We only used complete scans with one scan per patient. Finally, some scans needed to be ordered top to bottom. Using the bilinear sampling algorithm, all images were resized to 256 by 256 pixels, and image values were divided by 255 to normalize between 0 and 1. Example images are shown in Figure 2B.

The MosMed dataset has a median of 41 slices, a minimum of 31 slices and a maximum of 72 slices. The Zhang et al. dataset has much greater variability in scan size with a median of 61 slices, a minimum of 19 slices, and a maximum of 415 slices. We present histograms showing the number of slices per scan in Figure 3. We require a fixed number of slices as input, and we chose to use 20 slices. For all scans, we included the first and last images. If scans had more than 20 slices, we sampled uniformly to select 20. Only one scan in the Zhang et al. dataset had less than 20 slices; a blank slice replaced this slice; the mixed-effects model can account for missing data.

While removing slices may waste some information available to us, using the full 415 slices that some images have would be impractical due to the large memory footprint. An alternative to removing slices would be to reduce the resolution of each slice; however, this again would waste information. Choosing to use 20 slices of each CT image is a compromise between the amount of information used and the practicality of processing the CT scans.

TABLE 1 Summary of the datasets used.

| Dataset | Location | Use | Healthy/COVID19 |
|-------------------|----------------|---------------------|-----------------|
| MosMed training | Moscow, Russia | Training | 169/856 |
| MosMed validation | Moscow, Russia | Internal validation | 85/285 |
| Zhang et al. (48) | China | External validation | 243/553 |

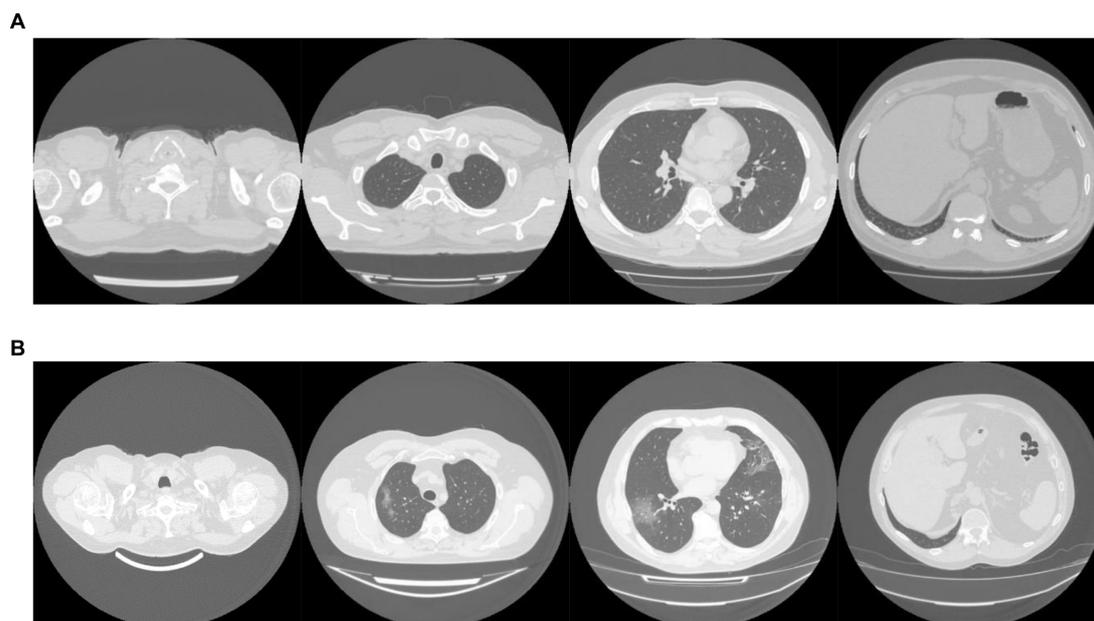
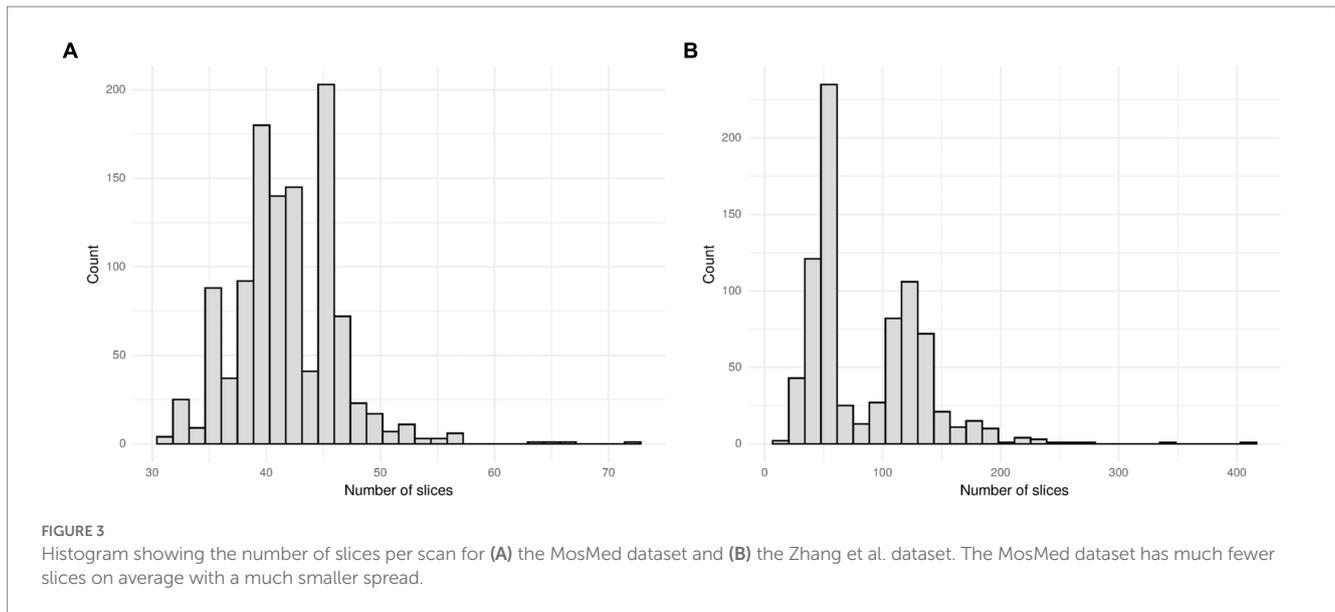


FIGURE 2 Example images showing (A) healthy and (B) COVID-19 lungs taken from the Mosmed dataset.



3. Results

On the internal validation dataset, the proposed model attained an AUROC of 0.936 (95%CI: 0.910, 0.961). Using a probability threshold of 0.5, the sensitivity, specificity, NPV, and PPV were 0.753 (0.647, 0.840), 0.909 (0.869, 0.940), 0.711 (0.606, 0.802), and 0.925 (0.888, 0.953), respectively. The model proposed by Bai et al. (8) attained an AUROC of 0.731 (0.674, 0.80). However, despite attaining a reasonable AUC value, the model was badly calibrated, and the predicted probabilities of COVID-19 were all clustered around 0.42, meaning that the sensitivity, specificity, PPV, and NPV are meaningless. We tried to retrain the model and rechecked the code implementation; however, we could not obtain more meaningful results. Covinet (9) attained an AUROC of 0.810 (0.748, 0.853). Using a probability threshold of 0.5, the sensitivity, specificity, NPV, and PPV were 0.824 (0.726, 0.898), 0.596 (0.537, 0.654), 0.378 (0.308, 0.452), and 0.919 (0.870, 0.954), respectively. COVNet (7) attained an AUROC of 0.935 (0.912, 0.959). Using a probability threshold of 0.5, the sensitivity, specificity, NPV, and PPV were 1.0 (0.958, 1.0), 0.796 (0.745, 0.842), 0.594 (0.509, 0.676), and 1.0 (0.984, 1.0), respectively. Full results for a range of probability thresholds are shown in Table 2, with ROC curves shown in Figure 4.

Calibration curves in Figure 5 show reasonable calibration for the mixed-effects model, although the model may still benefit from some recalibration. The other models do not have good calibration and likely provide harmful predictions. The decision curve in Figure 6 shows that the proposed model is of great clinical benefit compared to the treat all and treat-none approach.

It is important to remember that the model was selected using this internal testing set to avoid overfitting on the training set; therefore, these results are biased, and the external validation results are more representative of the true model performance.

On the external geographical validation dataset, the proposed model attained an AUROC of 0.930 (0.914, 0.947). With a probability threshold of 0.5, the sensitivity, specificity, NPV, and PPV were 0.778 (0.720, 0.828), 0.882 (0.853, 0.908), 0.744 (0.686, 0.797), and 0.90 (0.872, 0.924), respectively. The model proposed by Bai et al. (8) again

attained a reasonable AUROC of 0.805 (0.774, 0.836); however, the sensitivity, specificity, NPV, and PPV were meaningless. Covinet (9) attained an AUROC of 0.651 (0.610, 0.691). Using a probability threshold of 0.5, the sensitivity, specificity, NPV, and PPV were 0.008 (0.001, 0.029), 0.991 (0.979, 0.997), 0.286 (0.037, 0.710), and 0.695 (0.661, 0.727), respectively. COVNet (7) attained an AUROC of 0.808 (0.775, 0.841). With a cut-off point of 0.5, the sensitivity, specificity, NPV, and PPV were 0.387 (0.325, 0.451), 0.940 (0.917, 0.959), 0.740 (0.655, 0.814), and 0.777 (0.744, 0.808), respectively. Full results are shown in Table 3.

Similar to the internal validation, Figure 7 shows reasonable calibration for the mixed-effects model, although some recalibration may improve performance. Again, the comparison models could give harmful predictions as they are poorly calibrated. The decision curve in Figure 8 shows that the model is of great clinical benefit compared to the treat all and treat-none approach.

Although our proposed method and the Covnet model showed comparable performance on the internal validation set, the Covnet model could not generalize to the external geographical validation set, and calibration showed that the Covnet model would provide harmful risk estimates. This highlights the need for robust external validation in each intended setting. Nevertheless, the results show that the proposed method better generalizes to external geographical datasets and provides less harmful predictions when compared to the four previously proposed methods based on the calibration curves.

3.1. Saliency maps

It is vital to understand how the algorithm makes decisions and to check that it identifies the correct features within the image. Saliency maps can be used as a visual check to see what features the algorithm is learning. For example, the saliency maps in Figure 9 show that the model correctly identifies the diseased areas of the scans. We used 100 samples with a smoothing noise of 0.05 to create these saliency maps.

TABLE 2 Area under the receiver operating characteristic curve (AUROC), sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) on the internal validation dataset.

| Model | AUROC | Threshold | Sensitivity | Specificity | PPV | NPV |
|-------------------------|-------------------------|-----------|----------------------|----------------------|----------------------|----------------------|
| Bai et al. | 0.731 (0.674, 0.80) | 0.3 | 0 0.0 (0.0, 0.042) | 1.0 (0.987, 1.0) | NA | 0.77 (0.724, 0.812) |
| | | 0.4 | 0.012 (0, 0.064) | 0.996 (0.981, 1.0) | 0.50 (0.013, 0.987) | 0.772 (0.725, 0.814) |
| | | 0.5 | 1.0 (0.958, 1.0) | 0.0 (0.0, 0.013) | 0.230 (0.188, 0.276) | NA |
| | | 0.6 | 1.0 (0.958, 1.0) | 0.0 (0.0, 0.013) | 0.230 (0.188, 0.276) | NA |
| | | 0.7 | 1.0 (0.958, 1.0) | 0.0 (0.0, 0.013) | 0.230 (0.188, 0.276) | NA |
| CoviNet | 0.801 (0.748, 0.853) | 0.3 | 0.459 (0.350, 0.570) | 0.898 (0.857, 0.931) | 0.574 (0.448, 0.693) | 0.848 (0.802, 0.886) |
| | | 0.4 | 0.706 (0.597, 0.80) | 0.761 (0.708, 0.810) | 0.469 (0.380, 0.559) | 0.897 (0.851, 0.932) |
| | | 0.5 | 0.824 (0.726, 0.898) | 0.596 (0.537, 0.654) | 0.378 (0.308, 0.452) | 0.919 (0.870, 0.954) |
| | | 0.6 | 0.918 (0.838, 0.966) | 0.446 (0.387, 0.505) | 0.331 (0.271, 0.394) | 0.948 (0.895, 0.979) |
| | | 0.7 | 0.965 (0.90, 0.993) | 0.246 (0.197, 0.30) | 0.276 (0.226, 0.331) | 0.959 (0.885, 0.991) |
| CovNet | 0.935 (0.912, 0.959) | 0.3 | 0.941 (0.868, 0.981) | 0.839 (0.791, 0.879) | 0.635 (0.544, 0.719) | 0.98 (0.953, 0.993) |
| | | 0.4 | 0.965 (0.90, 0.993) | 0.825 (0.775, 0.867) | 0.621 (0.533, 0.704) | 0.987 (0.964, 0.997) |
| | | 0.5 | 1.0 (0.958, 1.0) | 0.796 (0.745, 0.842) | 0.594 (0.509, 0.676) | 1.0 (0.984, 1.0) |
| | | 0.6 | 1.0 (0.958, 1.0) | 0.779 (0.726, 0.826) | 0.574 (0.490, 0.655) | 1.0 (0.984, 1.0) |
| | | 0.7 | 1.0 (0.958, 1.0) | 0.761 (0.708, 0.810) | 0.556 (0.473, 0.636) | 1.0 (0.984, 1.0) |
| Mixed-effects (ours) | 0.936 (0.910, 0.961) | 0.3 | 0.588 (0.476, 0.694) | 0.961 (0.932, 0.981) | 0.820 (0.70, 0.906) | 0.887 (0.846, 0.920) |
| | | 0.4 | 0.659 (0.548, 0.758) | 0.933 (0.898, 0.959) | 0.747 (0.633, 0.840) | 0.902 (0.862, 0.933) |
| | | 0.5 | 0.753 (0.647, 0.840) | 0.909 (0.869, 0.940) | 0.711 (0.606, 0.802) | 0.925 (0.888, 0.953) |
| | | 0.6 | 0.812 (0.712, 0.888) | 0.884 (0.841, 0.919) | 0.676 (0.577, 0.766) | 0.940 (0.905, 0.960) |
| | | 0.7 | 0.906 (0.823, 0.958) | 0.832 (0.783, 0.873) | 0.616 (0.525, 0.702) | 0.967 (0.937, 0.986) |

Point estimates and 95% confidence intervals were calculated using De Long’s method for AUROC and Jeffrey’s interval for sensitivity, specificity, PPV, and NPV. Results are shown at a range of probability thresholds.

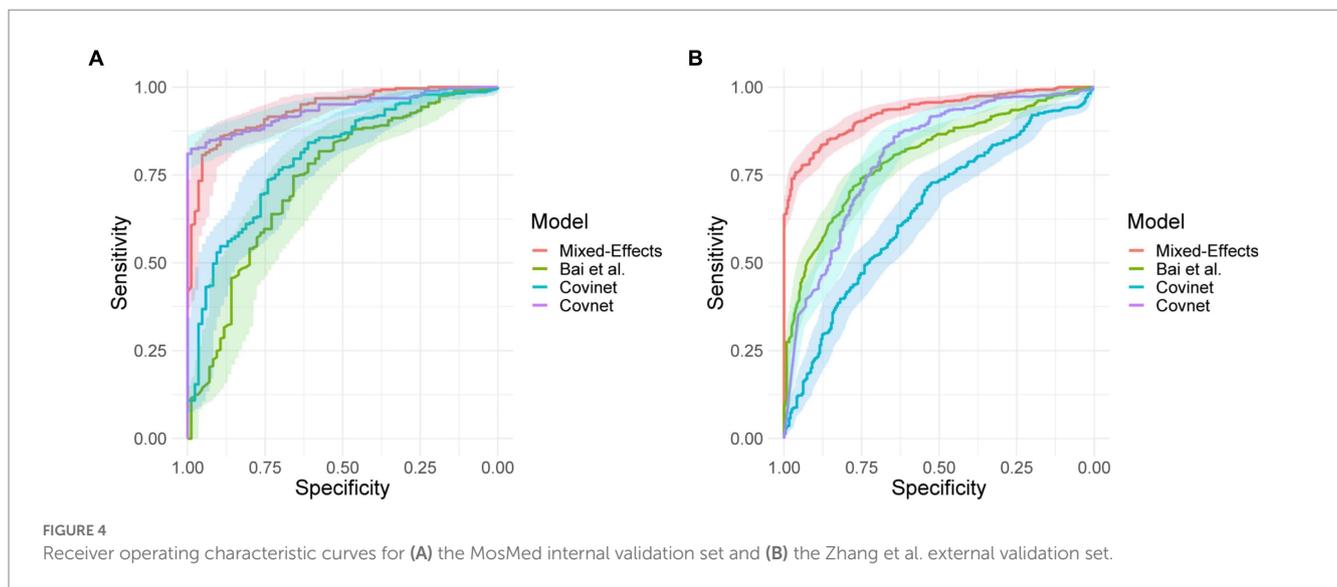
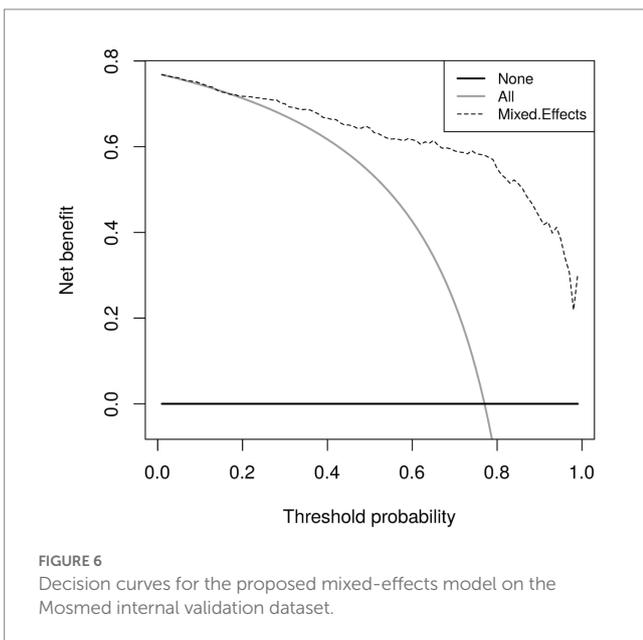
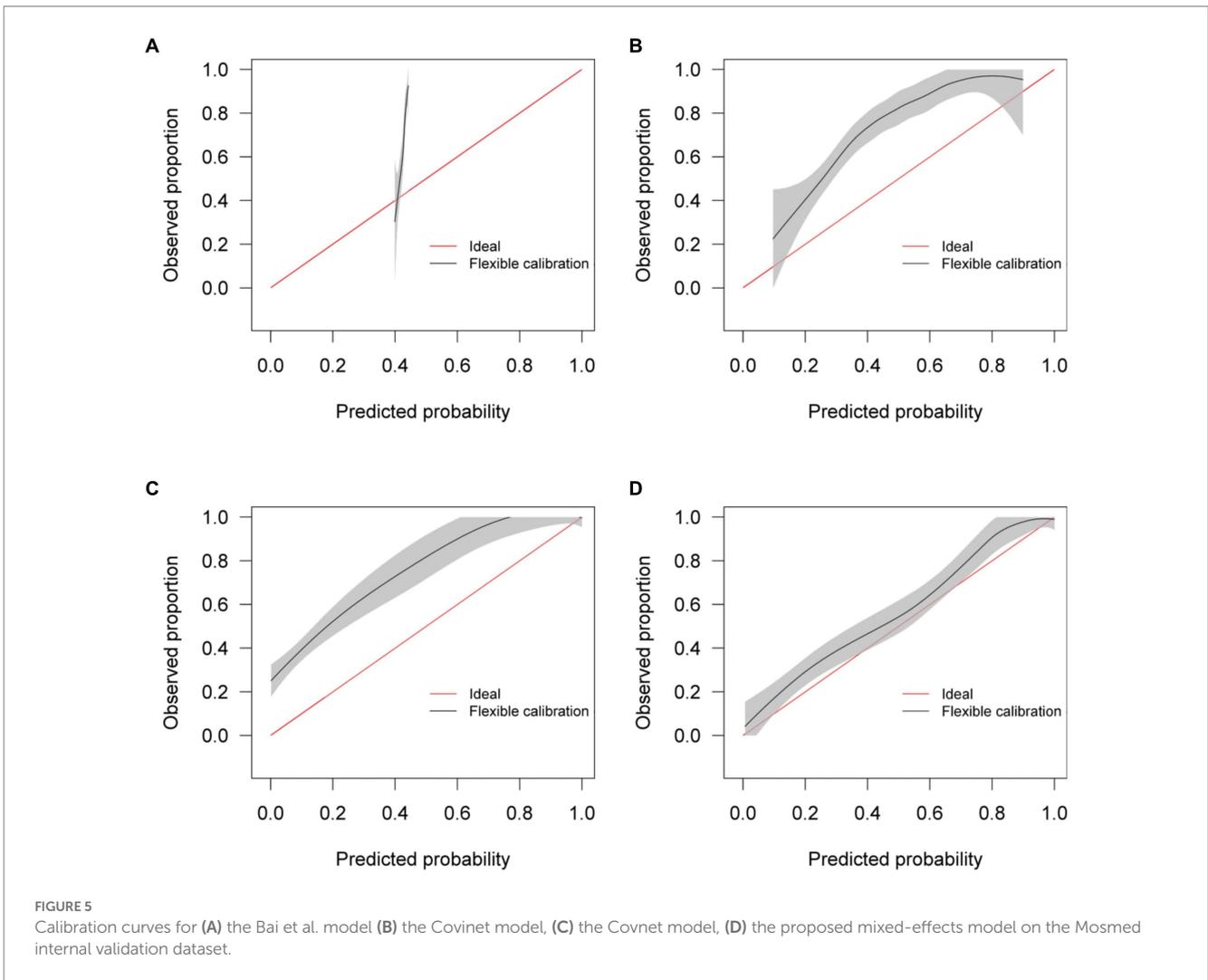


FIGURE 4 Receiver operating characteristic curves for (A) the MosMed internal validation set and (B) the Zhang et al. external validation set.

3.2. Sensitivity analysis

Mixed-effects models are capable of accounting for missing data. However, only one image had less than 20 slices; hence, we could not adequately assess if our model can indeed maintain good performance with missing data. Here, we rerun the analysis using the same dataset, using the same model and weights; however,

we reduce the number of slices available as testing data inputs to simulate missing data. Blank images replace these slices. We uniformly sampled the slices choosing between 10 and 19 slices; this equates to between 5 and 50% missing data for the model. We ran inference at each level of missingness and briefly show the AUROC to determine at which point the predictive performance is significantly reduced.



The plot of AUROCs at different levels of missingness is shown in Figure 10, along with 95% confidence intervals. We can see that at 20% missingness, there is a statistically significant decrease in

predictive performance. Although, even at 50% missingness, the model still performs relatively well, with an AUROC of 0.890 (95% CI: 0.868, 0.912). It should be noted that this does not mean that there is no reduction in performance at 5–15% missingness, only that the reduction was not statistically significant at the 95% confidence level.

Deep learning models can be susceptible to adversarial attacks (50), where minor artifacts or noise on an image can cause the image to be misclassified, even when the image does not look significantly different to a human observer. Here, we perform a brief sensitivity analysis by adding a small Gaussian noise to the image. We tested the model performance on the external dataset, with each image having a random Gaussian noise added. Experiments were conducted with standard deviations of 0 up to 0.005 in increments of 0.001 added to the normalized image. We did not add Gaussian noise in the data augmentation so that the model is not explicitly trained to deal with this kind of attack.

When using a variance of 0, the images are unchanged, and the results are the same as the standard results above. We present results on the Zhang et al. (49) dataset. Example images for each level of variance are shown in Figure 11, and a graph showing the reduction in AUROC is shown in Figure 12.

TABLE 3 Area under the receiver operating characteristic curve (AUROC), sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) on the external validation dataset.

| Model | AUROC | Threshold | Sensitivity | Specificity | PPV | NPV |
|-------------------------|-------------------------|-----------|----------------------|----------------------|----------------------|----------------------|
| Bai et al. | 0.805 (0.774, 0.836) | 0.3 | 0.0 (0.0, 0.015) | 1.0 (0.993, 1.0) | NA | 0.695 (0.661, 0.727) |
| | | 0.4 | 0.0 (0.0, 0.015) | 1.0 (0.993, 1.0) | NA | 0.695 (0.661, 0.727) |
| | | 0.5 | 1.0 (0.985, 1.0) | 1.0 (0.0, 0.007) | 0.305 (0.273, 0.339) | NA |
| | | 0.6 | 1.0 (0.985, 1.0) | 1.0 (0.0, 0.007) | 0.305 (0.273, 0.339) | NA |
| | | 0.7 | 1.0 (0.985, 1.0) | 1.0 (0.0, 0.007) | 0.305 (0.273, 0.339) | NA |
| CoviNet | 0.651 (0.610, 0.691) | 0.3 | 0.0 (0.0, 0.015) | 1.0 (0.993, 1.0) | NA | 0.695 (0.661, 0.727) |
| | | 0.4 | 0.0 (0.0, 0.015) | 1.0 (0.993, 1.0) | NA | 0.695 (0.661, 0.727) |
| | | 0.5 | 0.008 (0.001, 0.029) | 0.991 (0.979, 0.997) | 0.286 (0.037, 0.710) | 0.695 (0.661, 0.727) |
| | | 0.6 | 0.160 (0.117, 0.213) | 0.929 (0.905, 0.949) | 0.50 (0.385, 0.615) | 0.716 (0.681, 0.749) |
| | | 0.7 | 0.551 (0.487, 0.615) | 0.694 (0.654, 0.733) | 0.442 (0.385, 0.50) | 0.779 (0.740, 0.815) |
| CovNet | 0.808 (0.775, 0.841) | 0.3 | 0.305 (0.247, 0.367) | 0.969 (0.951, 0.982) | 0.813 (0.718, 0.887) | 0.760 (0.727, 0.791) |
| | | 0.4 | 0.354 (0.294, 0.418) | 0.955 (0.934, 0.971) | 0.775 (0.686, 0.849) | 0.771 (0.737, 0.802) |
| | | 0.5 | 0.387 (0.325, 0.451) | 0.940 (0.917, 0.959) | 0.740 (0.655, 0.814) | 0.777 (0.744, 0.808) |
| | | 0.6 | 0.432 (0.369, 0.497) | 0.937 (0.913, 0.956) | 0.750 (0.670, 0.819) | 0.790 (0.756, 0.820) |
| | | 0.7 | 0.473 (0.409, 0.538) | 0.931 (0.907, 0.951) | 0.752 (0.675, 0.818) | 0.801 (0.768, 0.831) |
| Mixed-effects (ours) | 0.930 (0.914, 0.947) | 0.3 | 0.675 (0.612, 0.733) | 0.935 (0.911, 0.954) | 0.820 (0.760, 0.871) | 0.867 (0.838, 0.894) |
| | | 0.4 | 0.741 (0.681, 0.795) | 0.904 (0.877, 0.927) | 0.773 (0.713, 0.825) | 0.888 (0.859, 0.913) |
| | | 0.5 | 0.778 (0.720, 0.828) | 0.882 (0.853, 0.908) | 0.744 (0.686, 0.797) | 0.90 (0.872, 0.924) |
| | | 0.6 | 0.827 (0.774, 0.873) | 0.859 (0.827, 0.887) | 0.720 (0.664, 0.772) | 0.919 (0.892, 0.941) |
| | | 0.7 | 0.885 (0.838, 0.922) | 0.828 (0.794, 0.859) | 0.694 (0.639, 0.744) | 0.942 (0.918, 0.961) |

Point estimates and 95% confidence intervals were calculated using De Long’s method for AUROC and Jeffrey’s interval for sensitivity, specificity, PPV, and NPV. Results are shown at a range of probability thresholds.

3.3. Fixed-effects only

To show that the mixed-effects method improves prediction over the fixed-effects method alone, we removed the random-effects part of the model to leave the fixed effects only. This was the only change to the model and allowed us to see the added benefit of the mixed-effects part. The full results are shown in Tables 4, 5. This experiment shows much worse performance when the random effects are removed from the model.

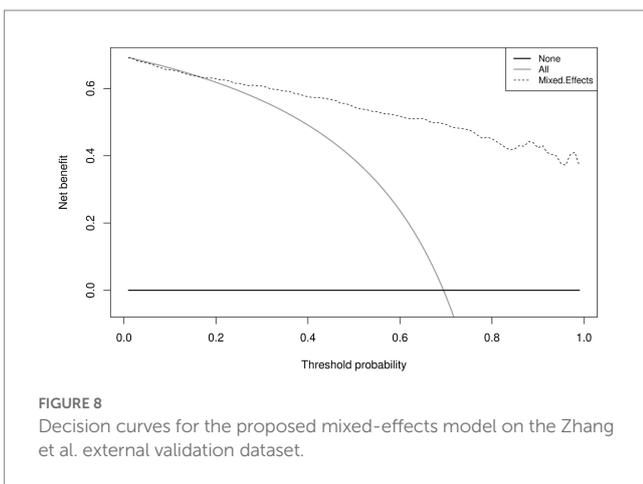
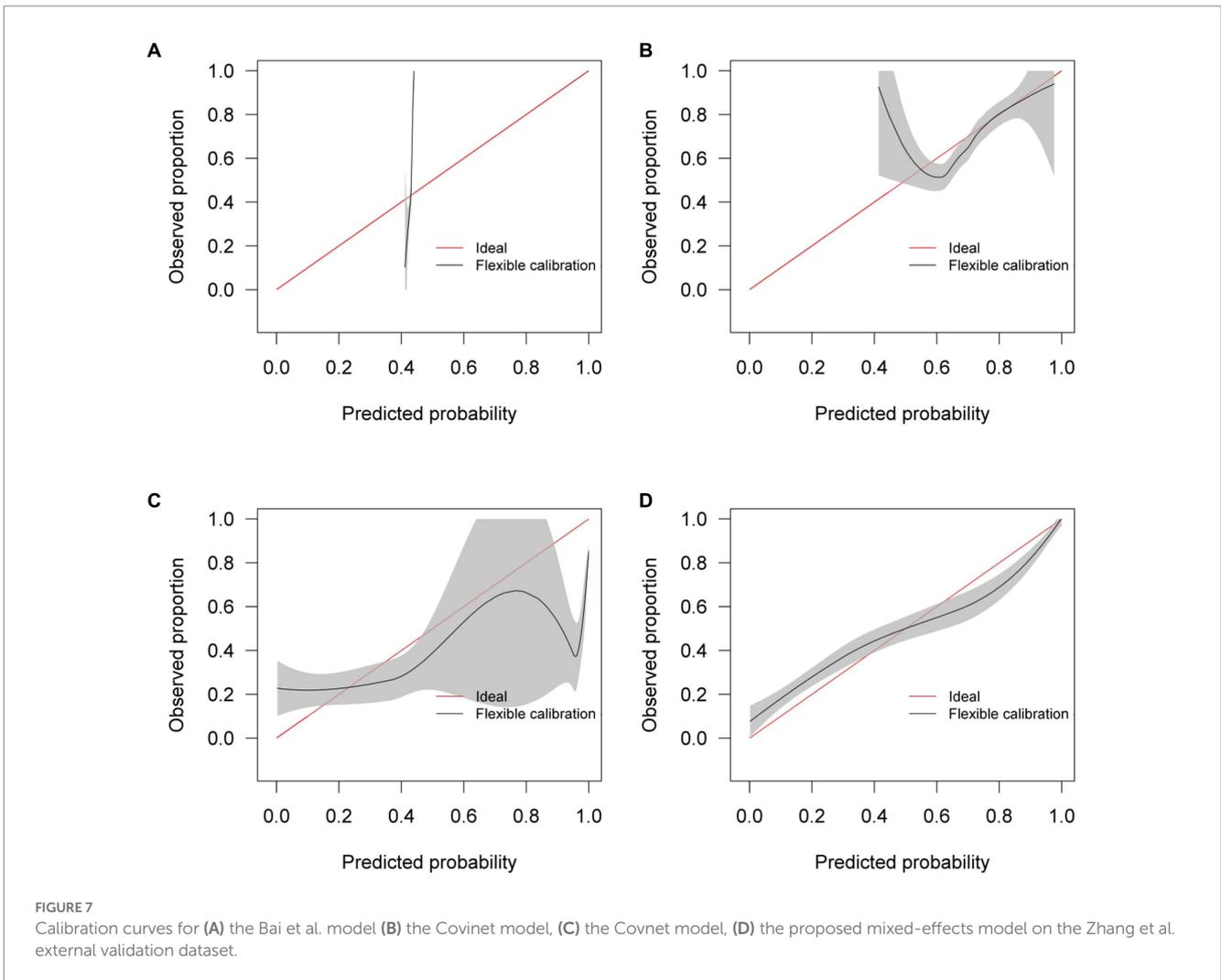
4. Discussion

Artificial intelligence is set to revolutionize healthcare, allowing large amounts of data to be processed and analyzed automatically, reducing pressure on stretched healthcare services. These tools can aid clinicians in monitoring and managing both common conditions and outbreaks of novel diseases. However, these tools must be assessed adequately, and best practice guidelines for reporting and development must be followed closely to increase reproducibility and reduce bias. We have developed a deep learning model to classify CT scans as healthy or COVID-19 using a novel mixed-effects model. Following best practice guidelines, we have externally validated the model. In addition, we robustly externally geographically validated the developed model in several performance areas, which are not routinely reported. For example, discriminative performance

measures show that the model can discriminate between healthy and COVID-19 CT scans well, calibration shows that the model is not clinically harmful. Finally, the clinical usefulness measures show that the model may be useful in a clinical setting. From the results presented here, it would seem that our deep learning model outperforms the RT-PCR tests as shown in the review by Watson et al. (3); however, those results are conservative estimates and were conducted under real-world clinical settings. A prospective study is required to determine if this is the case.

Compared to previously proposed models, our model showed similar discriminative performance to one existing method; however, our method generalized better to an external geographical validation set and showed improved calibration performance. Interestingly, in both internal and external validation, the sensitivity and NPV are similar in all models. However, the specificity and PPV are statistically significantly improved for the mixed-effects model in the external validation dataset. The performance of the proposed model in the external validation set is similar to that reported by PCR testing (3). However, a direct comparison should not be made as PCR testing on this exact dataset is unavailable.

There are several limitations of the study that should be highlighted and improved in future work. Firstly, we have only performed external geographical validation in a single dataset. Further external validation, both geographical and temporal, is needed on many datasets to determine if the model is correct in each intended setting. Although we performed a brief sensitivity analysis



here, more extensive work on adversarial attacks is needed. Future studies could consider following the method proposed by Goodfellow et al. (50) to improve robustness against adversarial examples. Patient demographic data were not available for this study, but future studies could incorporate this data into the model to

improve results. Finally, rules of thumb for assessing sample size calculations in the validation set can lead to imprecise results (51). Simulating data is a better alternative; however, it is difficult to anticipate the distribution of the model's linear predictor. Therefore, we were required to revert to the rule of thumb using a minimum of 200 samples in each group (29).

Initial experiments used the Zhang et al. (49) dataset for training; this showed promising results on the internal validation set; however, external validation showed random results. In addition, saliency maps showed that the model was not using the features of COVID-19 to make the diagnosis and was instead using the area around the image. We concluded that the images for each class were slightly different, perhaps due to different imaging protocols, and the algorithm was learning the image format rather than the disease. We then used the MosMed dataset for training and the Zhang et al. (49) dataset for external validation. This highlights the need for good quality training data and external validation and visualization.

Future studies should validate models and follow reporting guidelines such as TRIPOD (17) or the upcoming QUADAD-AI (52) and TRIPOD-AI (53) to bring about clinically useful and deployable models. Further research could look deeper into the areas of images identified by the algorithm as shown on the

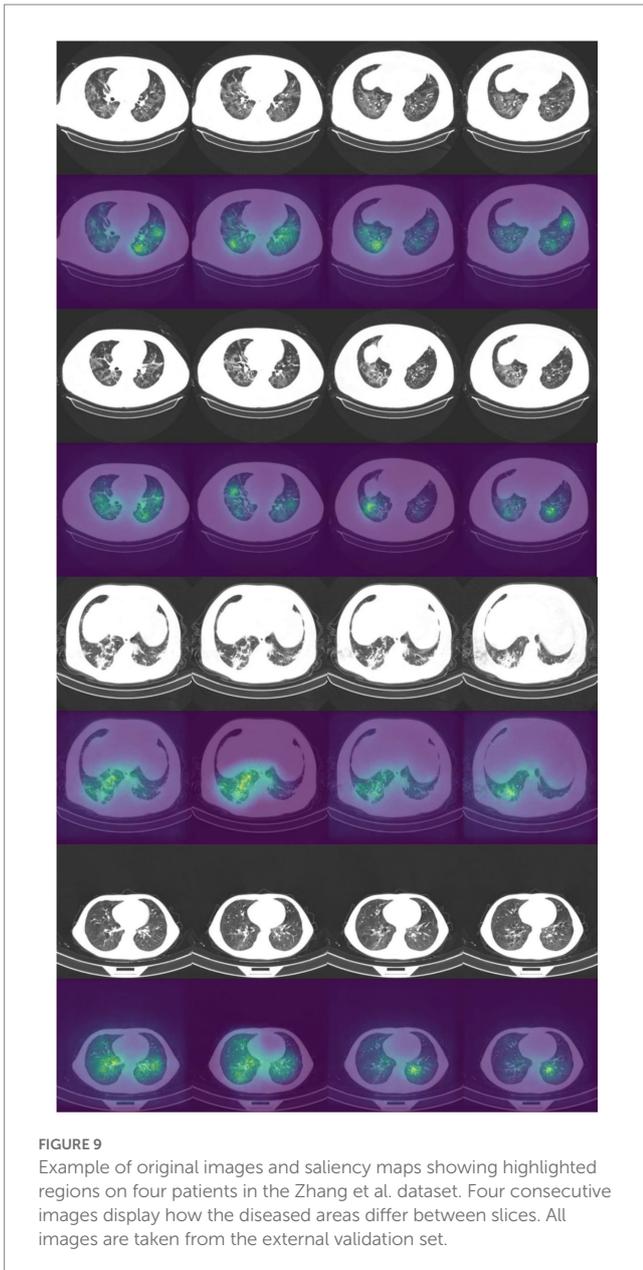


FIGURE 9
Example of original images and saliency maps showing highlighted regions on four patients in the Zhang et al. dataset. Four consecutive images display how the diseased areas differ between slices. All images are taken from the external validation set.

saliency maps; this could potentially identify new features of COVID-19 which have gone unnoticed. Before any model can be fully deployed, clinical trials are needed to study the full impact of using such algorithms to diagnose COVID-19 and the exact situations in which such a model may be used. In-clinic prospective studies comparing the performance deep learning models with RT-PCR and lateral flow tests should be carried out to determine how deep learning compares; this will show whether deep learning could be used as an automated alternative to RT-PCR testing.

This study indicates that deep learning could be suitable for screening and monitoring of COVID-19 in a clinical setting; however, validation in the intended setting is vital, and models should not be adopted without this. It has been observed that the quality of reporting of deep learning prediction models is usually very poor;

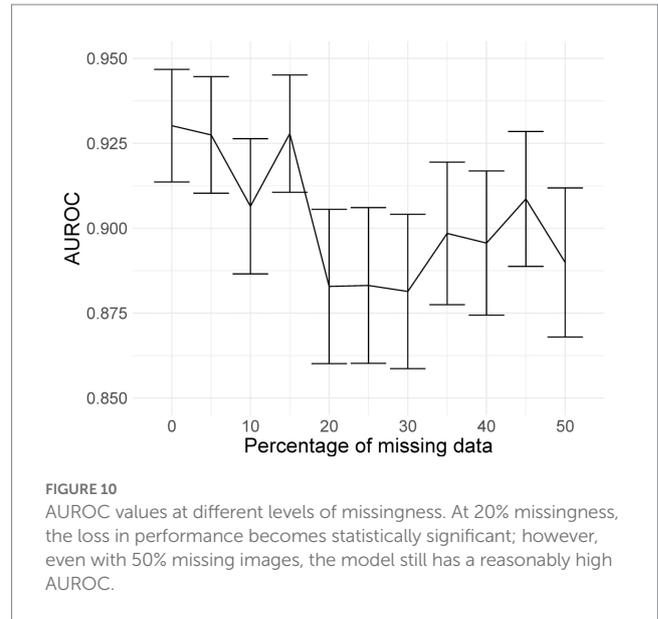


FIGURE 10
AUROC values at different levels of missingness. At 20% missingness, the loss in performance becomes statistically significant; however, even with 50% missing images, the model still has a reasonably high AUROC.

however, with a bit of extra work and by following best practice guidelines, this problem can be overcome. This study highlights the importance of robust analysis and reporting of models with external validation.

Data availability statement

The source code and datasets presented in this study can be found in an online repository. The accession link to the source code is: <https://github.com/JTBridge/ME-COVID19>. Publicly available datasets were analyzed in this study. These data can be found at: the CNCB and MosMed repositories, available from doi: 10.17816/DD46826 and doi: 10.1016/j.cell.2020.04.045.

Ethics statement

Ethical approval was not required for the study involving humans in accordance with the local legislation and institutional requirements. Written informed consent to participate in this study was not required from the participants or the participants' legal guardians/next of kin in accordance with the national legislation and the institutional requirements.

Author contributions

JB, YM, and YaZ: conception. JB, YM, YiZ, and YaZ: methodology. CA, MW, CMe, SF, MM, SM, and YaZ: administration. JB, YM, WZ, TF, CMc, RS, YiZ, and YaZ: investigation. JB, YM, CA, MW, CMe, SF, MM, SM, RS, and YiZ: data curation. JB: analysis and validation. JB, YM, WZ, and YaZ: writing of first draft. All authors made substantial contributions to the reviewing and editing of the manuscript.

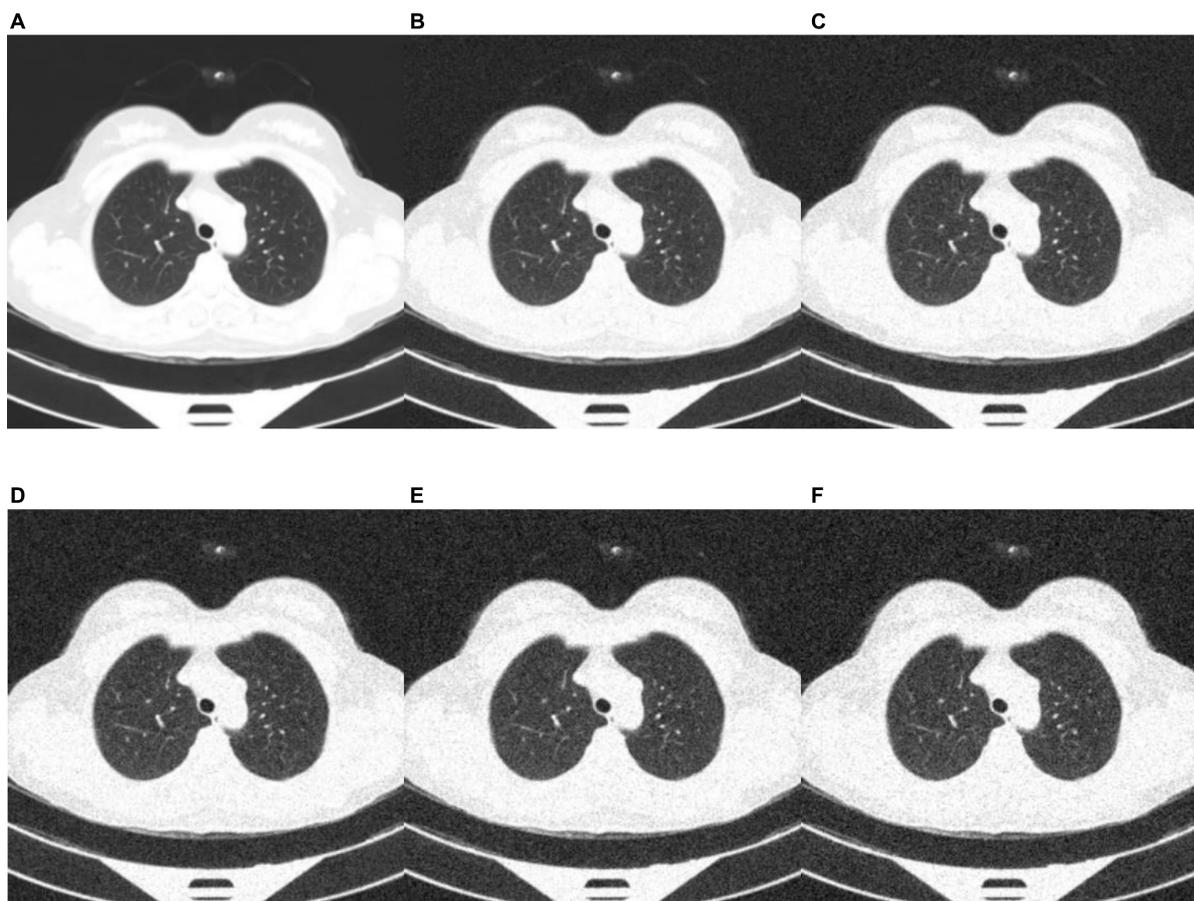


FIGURE 11
 Example images showing the effect of increasing the amount of noise in the image input. (A) no noise; (B) deviation = 0.001; (C) deviation = 0.002; (D) deviation = 0.003; (E) deviation = 0.004; (F) deviation = 0.005.

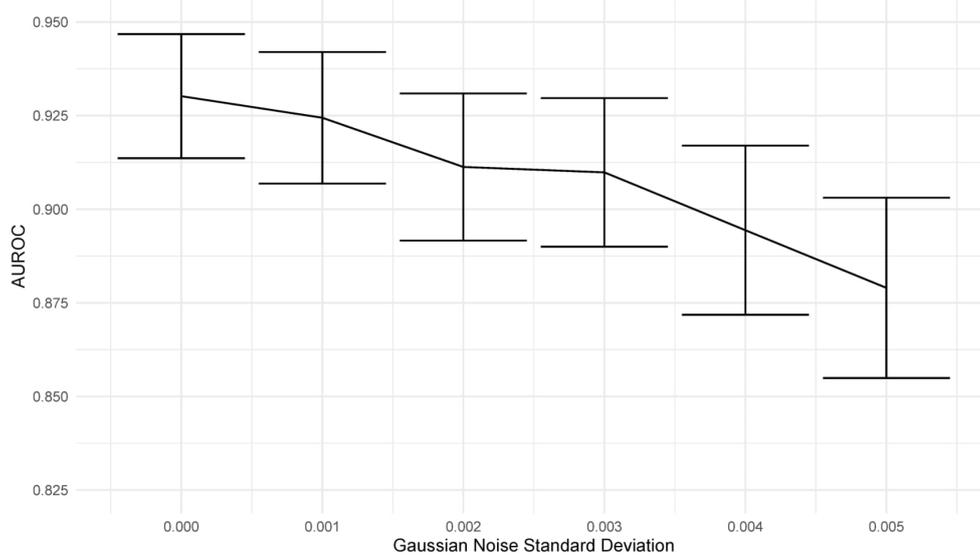


FIGURE 12
 Graph showing the drop in AUROC as the amount of noise in the image input increases. The AUROC falls steadily with increased noise in the image.

TABLE 4 Area under the receiver operating characteristic curve (AUROC), sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) on the internal validation dataset for our proposed model and the fixed effects model.

| Model | AUROC | Threshold | Sensitivity | Specificity | PPV | NPV |
|-----------------------------------|-------------------------|-----------|----------------------|----------------------|----------------------|----------------------|
| Fixed effects | 0.494 (0.427, 0.561) | 0.3 | 0.859 (0.766, 0.925) | 0.165 (0.124, 0.213) | 0.235 (0.189, 0.286) | 0.797 (0.672, 0.890) |
| | | 0.4 | 0.953 (0.884, 0.987) | 0.046 (0.025, 0.077) | 0.229 (0.187, 0.277) | 0.765 (0.501, 0.932) |
| | | 0.5 | 0.988 (0.936, 1.0) | 0.014 (0.004, 0.036) | 0.231 (0.188, 0.277) | 0.80 (0.284, 0.995) |
| | | 0.6 | 1.0 (0.958, 1.0) | 0.0 (0.0, 1.0) | 0.230 (0.188, 0.276) | NA (NA, NA) |
| | | 0.7 | 1.0 (0.958, 1.0) | 0.0 (0.0, 1.0) | 0.230 (0.188, 0.276) | NA (NA, NA) |
| Mixed-effects (fixed + random) | 0.936 (0.910, 0.961) | 0.3 | 0.588 (0.476, 0.694) | 0.961 (0.932, 0.981) | 0.820 (0.70, 0.906) | 0.887 (0.846, 0.920) |
| | | 0.4 | 0.659 (0.548, 0.758) | 0.933 (0.898, 0.959) | 0.747 (0.633, 0.840) | 0.902 (0.862, 0.933) |
| | | 0.5 | 0.753 (0.647, 0.840) | 0.909 (0.869, 0.940) | 0.711 (0.606, 0.802) | 0.925 (0.888, 0.953) |
| | | 0.6 | 0.812 (0.712, 0.888) | 0.884 (0.841, 0.919) | 0.676 (0.577, 0.766) | 0.940 (0.905, 0.960) |
| | | 0.7 | 0.906 (0.823, 0.958) | 0.832 (0.783, 0.873) | 0.616 (0.525, 0.702) | 0.967 (0.937, 0.986) |

Point estimates and 95% confidence intervals were calculated using De Long's method for AUROC and Jeffrey's interval for sensitivity, specificity, PPV, and NPV. Results are shown at a range of probability thresholds.

TABLE 5 Area under the receiver operating characteristic curve (AUROC), sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) on the external validation dataset for our proposed model and the fixed effects model.

| Model | AUROC | Threshold | Sensitivity | Specificity | PPV | NPV |
|-----------------------------------|-------------------------|-----------|----------------------|----------------------|----------------------|----------------------|
| Fixed effects | 0.630 (0.590, 0.670) | 0.3 | 0.794 (0.738, 0.843) | 0.374 (0.334, 0.416) | 0.358 (0.317, 0.40) | 0.805 (0.752, 0.852) |
| | | 0.4 | 0.971 (0.942, 0.988) | 0.159 (0.130, 0.192) | 0.337 (0.302, 0.373) | 0.926 (0.854, 0.970) |
| | | 0.5 | 1.0 (0.984, 1.0) | 0.063 (0.044, 0.087) | 0.319 (0.286, 0.354) | 1.0 (0.90, 1.0) |
| | | 0.6 | 1.0 (0.985, 1.0) | 0.018 (0.277, 0.343) | 0.309 (0.277, 0.343) | 1.0 (0.692, 1.0) |
| | | 0.7 | 1.0 (0.985, 1.0) | 0.004 (0.0, 0.013) | 0.306 (0.274, 0.339) | 1.0 (0.158, 1.0) |
| Mixed-effects (fixed + random) | 0.930 (0.914, 0.947) | 0.3 | 0.675 (0.612, 0.733) | 0.935 (0.911, 0.954) | 0.820 (0.760, 0.871) | 0.867 (0.838, 0.894) |
| | | 0.4 | 0.741 (0.681, 0.795) | 0.904 (0.877, 0.927) | 0.773 (0.713, 0.825) | 0.888 (0.859, 0.913) |
| | | 0.5 | 0.778 (0.720, 0.828) | 0.882 (0.853, 0.908) | 0.744 (0.686, 0.797) | 0.90 (0.872, 0.924) |
| | | 0.6 | 0.827 (0.774, 0.873) | 0.859 (0.827, 0.887) | 0.720 (0.664, 0.772) | 0.919 (0.892, 0.941) |
| | | 0.7 | 0.885 (0.838, 0.922) | 0.828 (0.794, 0.859) | 0.694 (0.639, 0.744) | 0.942 (0.918, 0.961) |

Point estimates and 95% confidence intervals were calculated using De Long's method for AUROC and Jeffrey's interval for sensitivity, specificity, PPV, and NPV. Results are shown at a range of probability thresholds.

Funding

This study received funding from EPSRC studentship (No. 2110275), EPSRC Impact Acceleration Account (IAA) Awards, and Amazon Web Services. The funders were not involved in the study design, collection, analysis, interpretation of data, the writing of this article or the decision to submit it for publication.

Conflict of interest

CMe and SF were employed by Alces Flight Ltd. MM and SM were employed by Amazon Web Services.

The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmed.2023.1113030/full#supplementary-material>

References

1. *Coronavirus disease 2019 (COVID-19) - symptoms, diagnosis and treatment* | *BMJ Best Practice* BMJ Publishing Group (2020) Available at: <https://bestpractice.bmj.com/topics/en-gb/3000201>.
2. Torjesen I. Covid-19 will become endemic but with decreased potency over time, scientists believe. *BMJ*. (2021) 372:n494. doi: 10.1136/bmj.n494
3. Watson J, Whiting PF, Brush JE. Interpreting a covid-19 test result. *BMJ*. (2020) 369:m1808. doi: 10.1136/bmj.m1808
4. Mallett S, Allen AJ, Graziadio S, Taylor SA, Sakai NS, Green K, et al. At what times during infection is SARS-CoV-2 detectable and no longer detectable using RT-PCR-based tests? A systematic review of individual participant data. *BMC Med*. (2020) 18:346. doi: 10.1186/s12916-020-01810-8
5. Ruan Z-R, Gong P, Han W, Huang M-Q, Han M. A case of coronavirus disease 2019 with twice negative nucleic acid testing within 8 days. *Chin Med J*. (2020) 133:1487–8. doi: 10.1097/CM9.0000000000000788
6. Pontone G, Scafuri S, Mancini ME, Agalato C, Guglielmo M, Baggiano A, et al. Role of computed tomography in COVID-19. *J Cardiovasc Comput Tomogr*. (2020) 15:27–36. doi: 10.1016/j.jcct.2020.08.013
7. Li L, Qin L, Xu Z, Yin Y, Wang X, Kong B, et al. Using artificial intelligence to detect COVID-19 and community-acquired pneumonia based on pulmonary CT: evaluation of the diagnostic accuracy. *Radiology*. (2020) 296:E65–71. doi: 10.1148/radiol.20200905
8. Bai HX, Wang R, Xiong Z, Hsieh B, Chang K, Halsey K, et al. Artificial intelligence augmentation of radiologist performance in distinguishing COVID-19 from pneumonia of other origin at chest CT. *Radiology*. (2020) 296:E156–65. doi: 10.1148/radiol.2020201491
9. Mittal B, Oh J. CoviNet: Covid-19 diagnosis using machine learning analyses for computerized tomography images. *Thirteenth International Conference on Digital Image Processing (ICDIP 2021)*. SPIE (2021).
10. MacCormick IJC, Zheng Y, Czanner S, Zhao Y, Diggle PJ, Harding SP, et al. Spatial statistical modelling of capillary non-perfusion in the retina. *Sci Rep*. (2017) 7:16792. doi: 10.1038/s41598-017-16620-x
11. Zhu W, Ku JY, Zheng Y, Knox PC, Kolamunnage-Dona R, Czanner G. Spatial linear mixed effects modelling for OCT images: SLME model. *J Imaging*. (2020) 6:44. doi: 10.3390/jimaging6060044
12. Albtoush OM, Al-Shdefat RB, Al-Akaleh A. Chest CT scan features from 302 patients with COVID-19 in Jordan. *Eur J Radiol Open*. (2020) 7:100295. doi: 10.1016/j.ejro.2020.100295
13. Haseli S, Khalili N, Bakhshayeshkaram M, Sanei Taheri M, Moharramzad Y. Lobar distribution of COVID-19 pneumonia based on chest computed tomography findings; a retrospective study. *Arch Acad Emerg Med*. (2020) 8:e55-e. doi: 10.22037/aaem.v8i1.665
14. Xiang C, Lu J, Zhou J, Guan L, Yang C, Chai C. CT findings in a novel coronavirus disease (COVID-19) pneumonia at initial presentation. *Biomed Res Int*. (2020) 2020:1–10. doi: 10.1155/2020/5436025
15. Wynants L, Van Calster B, Collins GS, Riley RD, Heinze G, Schuit E, et al. Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal. *BMJ*. (2020) 369:m1328. doi: 10.1136/bmj.m1328
16. Roberts M, Driggs D, Thorpe M, Gilbey J, Yeung M, Ursprung S, et al. Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans. *Nat Mach Intell*. (2021) 3:199–217. doi: 10.1038/s42256-021-00307-0
17. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *BMJ*. (2015) 350:g7594. doi: 10.1136/bmj.g7594
18. Mongan J, Moy L, Kahn CE. Checklist for artificial intelligence in medical imaging (CLAIM): a guide for authors and reviewers. *Radiol Artif Intell*. (2020) 2:e200029. doi: 10.1148/ryai.2020200029
19. Wolff RF, Moons KGM, Riley RD, Whiting PF, Westwood M, Collins GS, et al. PROBAST: a tool to assess the risk of bias and applicability of prediction model studies. *Ann Intern Med*. (2019) 170:51–8. doi: 10.7326/M18-1376
20. Jiang J, Nguyen T. *Linear and generalized linear mixed models and their applications*. New York, NY: Springer (2007).
21. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z. Rethinking the inception architecture for computer vision. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE (2016);2818–2826.
22. Deng J, Dong W, Socher R, Li L, Kai L, Li F-F. ImageNet: a large-scale hierarchical image database. *2019 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE (2009); 248–255.
23. Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, et al. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*. (2009) 338:b2393. doi: 10.1136/bmj.b2393
24. Xiong Y, Kim HJ, Singh V. Mixed effects neural networks (MeNets) with applications to gaze estimation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE (2019);7735–7744.
25. Brier GW. Verification of forecasts expressed in terms of probability. *Mon Weather Rev*. (1950) 78:1–3. doi: 10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2
26. Platt J. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Adv Large Marg Classif*. (1999) 10:61–74.
27. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al. Assessing the performance of prediction models: a framework for traditional and novel measures. *Epidemiology*. (2010) 21:128–38. doi: 10.1097/EDE.0b013e3181c30fb2
28. Van Calster B, McLernon DJ, van Smeden M, Wynants L, Steyerberg EW, Bossuyt P, et al. Calibration: the Achilles heel of predictive analytics. *BMC Med*. (2019) 17:230. doi: 10.1186/s12916-019-1466-7
29. Van Calster B, Nieboer D, Vergouwe Y, De Cock B, Pencina MJ, Steyerberg EW. A calibration hierarchy for risk models was defined: from utopia to empirical data. *J Clin Epidemiol*. (2016) 74:167–76. doi: 10.1016/j.jclinepi.2015.12.005
30. R Core Team. *R: a language and environment for statistical computing*. R Foundation for Statistical Computing. (2021).
31. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*. (1988) 44:837–45. doi: 10.2307/2531595
32. Stevenson M, Sergeant E, Nunes T, Heuer C, Marshall J, Sanchez J, et al. *epiR: tools for the analysis of epidemiological data*. (2022). Available at: <https://CRAN.R-project.org/package=epiR>
33. Brown LD, Cai TT, DasGupta A. Interval estimation for a binomial proportion. *Stat Sci*. (2001) 16:101–17. doi: 10.1214/ss/1009213286
34. Wynants L, van Smeden M, McLernon DJ, Timmerman D, Steyerberg EW, Van Calster B, et al. Three myths about risk thresholds for prediction models. *BMC Med*. (2019) 17:192. doi: 10.1186/s12916-019-1425-3
35. Harrell FE Jr. *rms: regression modeling strategies*. (2021). Available at: <https://CRAN.R-project.org/package=rms>
36. Vickers AJ, Elkin EB. Decision curve analysis: a novel method for evaluating prediction models. *Med Decis Mak*. (2006) 26:565–74. doi: 10.1177/0272989X06295361
37. Smilkov D, Thorat N, Kim B, Viégas F, Wattenberg M. Smoothgrad: removing noise by adding noise. *arXiv*. (2017) arXiv:170603825. doi: 10.48550/arXiv.1706.03825
38. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE (2016);770–778.
39. Tan M, Le Q. Efficientnet: rethinking model scaling for convolutional neural networks. *International Conference on Machine Learning* (2019) PMLR.
40. Mittal B, Oh J. CoviNet: Covid-19 diagnosis using machine learning analyses for computerized tomography images. *SPIE Proceeding* (2021).
41. Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al. Tensorflow: large-scale machine learning on heterogeneous distributed systems. *SPIE* (2016) arXiv:160304467. doi: 10.48550/arXiv.1603.04467
42. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, et al., editors. *Tensorflow: a system for large-scale machine learning*. *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. USENIX Association. (2016).
43. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, et al. pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics*. (2011) 12:77. doi: 10.1186/1471-2105-12-77
44. Du Z, Hao Y. *reportROC: an easy way to report ROC analysis*. R package version 3.5. (2020). Available at: <https://CRAN.R-project.org/package=reportROC>
45. Wickham H. *ggplot2: elegant graphics for data analysis*. New York, NY: Springer-Verlag (2016).
46. Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. *arXiv*. (2014) 1412.6980. doi: 10.48550/arXiv.1412.6980
47. Riley RD, Snell KIE, Ensor J, Burke DL, Harrell FE Jr, Moons KGM, et al. Minimum sample size for developing a multivariable prediction model: PART II - binary and time-to-event outcomes. *Stat Med*. (2019) 38:1276–96. doi: 10.1002/sim.7992
48. Morozov SP, Andreychenko AE, Blokhin IA, Gelezhe PB, Gonchar AP, Nikolaev AE, et al. MosMedData: data set of 1110 chest CT scans performed during the COVID-19 epidemic. *Digit Diagn*. (2020) 1:49–59. doi: 10.17816/DD46826
49. Zhang K, Liu X, Shen J, Li Z, Sang Y, Wu X, et al. Clinically applicable AI system for accurate diagnosis, quantitative measurements, and prognosis of COVID-19 pneumonia using computed tomography. *Cells*. (2020) 181:1423–33.e11. doi: 10.1016/j.cell.2020.04.045
50. Goodfellow IJ, Shlens J, Szegedy C. Explaining and harnessing adversarial examples. *arXiv*. (2014) 1412.6572. doi: 10.48550/arXiv.1412.6572

51. Snell KIE, Archer L, Ensor J, Bonnett LJ, Debray TPA, Phillips B, et al. External validation of clinical prediction models: simulation-based sample size calculations were more reliable than rules-of-thumb. *J Clin Epidemiol.* (2021) 135:79–89. doi: 10.1016/j.jclinepi.2021.02.011

52. Sounderajah V, Ashrafian H, Rose S, Shah NH, Ghassemi M, Golub R, et al. A quality assessment tool for artificial intelligence-centered diagnostic test

accuracy studies: QUADAS-AI. *Nat Med.* (2021) 27:1663–5. doi: 10.1038/s41591-021-01517-0

53. Collins GS, Dhiman P, Andaur Navarro CL, Ma J, Hoof L, Reitsma JB, et al. Protocol for development of a reporting guideline (TRIPOD-AI) and risk of bias tool (PROBAST-AI) for diagnostic and prognostic prediction model studies based on artificial intelligence. *BMJ Open.* (2021) 11:e048008. doi: 10.1136/bmjopen-2020-048008