



BMJ Open Assessment of the effect of a comprehensive chest radiograph deep learning model on radiologist reports and patient outcomes: a real-world observational study

Catherine M Jones,^{1,2} Luke Danaher,² Michael R Milne ^{1,2} Cyril Tang,¹ Jarrel Seah ^{1,3} Luke Oakden-Rayner,⁴ Andrew Johnson,¹ Quinlan D Buchlak,^{1,5} Nazanin Esmaili^{5,6}

To cite: Jones CM, Danaher L, Milne MR, *et al.* Assessment of the effect of a comprehensive chest radiograph deep learning model on radiologist reports and patient outcomes: a real-world observational study. *BMJ Open* 2021;**11**:e052902. doi:10.1136/bmjopen-2021-052902

► Prepublication history and additional supplemental material for this paper are available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2021-052902>).

Patients and public were not involved in the design, conduct, or reporting of this study.

Received 04 May 2021

Accepted 29 November 2021



© Author(s) (or their employer(s)) 2021. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

Correspondence to

Mr Michael R Milne;
michael.milne@annalise.ai

ABSTRACT

Objectives Artificial intelligence (AI) algorithms have been developed to detect imaging features on chest X-ray (CXR) with a comprehensive AI model capable of detecting 124 CXR findings being recently developed. The aim of this study was to evaluate the real-world usefulness of the model as a diagnostic assistance device for radiologists.

Design This prospective real-world multicentre study involved a group of radiologists using the model in their daily reporting workflow to report consecutive CXRs and recording their feedback on level of agreement with the model findings and whether this significantly affected their reporting.

Setting The study took place at radiology clinics and hospitals within a large radiology network in Australia between November and December 2020.

Participants Eleven consultant diagnostic radiologists of varying levels of experience participated in this study.

Primary and secondary outcome measures Proportion of CXR cases where use of the AI model led to significant material changes to the radiologist report, to patient management, or to imaging recommendations. Additionally, level of agreement between radiologists and the model findings, and radiologist attitudes towards the model were assessed.

Results Of 2972 cases reviewed with the model, 92 cases (3.1%) had significant report changes, 43 cases (1.4%) had changed patient management and 29 cases (1.0%) had further imaging recommendations. In terms of agreement with the model, 2569 cases showed complete agreement (86.5%). 390 (13%) cases had one or more findings rejected by the radiologist. There were 16 findings across 13 cases (0.5%) deemed to be missed by the model. Nine out of 10 radiologists felt their accuracy was improved with the model and were more positive towards AI poststudy.

Conclusions Use of an AI model in a real-world reporting environment significantly improved radiologist reporting and showed good agreement with radiologists, highlighting the potential for AI diagnostic support to improve clinical practice.

Strengths and limitations of this study

- This study substantially adds to the limited literature on real-world evaluation of comprehensive chest X-ray artificial intelligence models in radiology workflow.
- This was a multicentre study conducted across a mix of public hospitals, private hospitals and community clinic settings.
- Due to the design of the study, diagnostic accuracy of the decision support system was not a measurable outcome.
- Results of this study are self-reported and may therefore be prone to bias.
- Determination of the significance of report changes due to the model's recommendations was made at the discretion of each radiologist on a case-by-case basis.

INTRODUCTION

Radiology is a data-rich medical specialty and is well placed to embrace artificial intelligence (AI).¹ This is especially true in high volume imaging tasks such as chest X-ray (CXR) imaging. The rapid application of X-ray technology to diagnosing chest diseases at the end of the 19th century led to the CXR becoming a first-line diagnostic imaging tool² and it remains an essential component of the diagnostic pathway for chest disease. Due to advancements in digital image acquisition, low ionising radiation dose and low cost, the chest radiograph is more easily accessible worldwide than any other imaging modality.³

The challenges of interpreting CXR, however, have not lessened over the last half-century. CXR images are two-dimensional representations of complex three-dimensional structures, relying on soft-tissue

contrast between structures of different densities. Multiple overlapping structures lead to reduced visibility of both normal and abnormal structures,⁴ with up to 40% of the lung parenchyma obscured by overlying ribs and the mediastinum.⁵ This can be further exacerbated by other factors including the degree of inspiration, other devices in the field of view, and patient positioning. In addition, there is a wide range of pathology in the chest which is visible to varying degrees on the CXR. These factors combine to make CXRs difficult to accurately interpret, with an error rate of 20%–50% for CXRs containing radiographic evidence of disease reported in the literature.⁶ Notably, lung cancer is one of the most common cancers worldwide and is the most common cause of cancer death,⁷ and CXR interpretation error accounts for 90% of cases where lung cancer is missed.⁸ Despite technological advancements in CXR over the past 50 years, this level of diagnostic error has remained constant.⁶

A rapidly developing field attempting to assist radiologists in radiological interpretation involves the application of machine learning, in particular deep neural networks.⁹ Deep neural networks learn patterns in large, complex datasets, enabling the detection of subtle features and outcome prediction.^{10–11} The potential of these algorithms has grown rapidly in the past decade thanks to the development of more useful neural network models, advancements in computational power and an increase in the volume and availability of digital imaging datasets.¹¹ Of note is the rise of convolutional neural networks (CNNs), a type of deep neural network that excels at image feature extraction and classification, and demonstrates strong performance in medical image analysis, leading to the rapid advancement of computer vision in medical imaging.^{12–13} CNNs have been used to develop models to successfully detect targeted clinical findings on CXR, including lung cancer,^{14–15} pneumonia,^{16–17} COVID-19,¹⁸ pneumothorax,^{19–22} pneumoconiosis,²³ cardiomegaly,²⁴ pulmonary hypertension²⁵ and tuberculosis.^{26–30} These studies highlight the effectiveness of applied machine learning in CXR interpretation, however, most of these deep learning systems are limited in scope to a single finding or a small set of findings, therefore lacking the broad utility that would make them useful in clinical practice.

Recently, our group developed a comprehensive deep learning CXR diagnostic assist device, which was designed to assist clinicians in CXR interpretation and improve diagnostic accuracy, validated for 124 clinically relevant findings seen on frontal and lateral chest radiographs.³¹ The primary objective of the current study was to evaluate the real-world usefulness of the model as a diagnostic assist device for radiologists in both hospital and community clinic settings. This involved examining the frequency at which the model's recommendations led to a 'significant impact on the report', defined as the inclusion of findings recommended by the model which altered the radiologists report in a meaningful way. The frequency of change in patient management and recommendations

for further imaging were also evaluated. Secondary endpoints included: (1) investigating agreement between radiologists and the findings detected by the model; and (2) assessing radiologist attitudes towards the tool and AI models in general.

METHODS

Model development and validation

A modified version of a commercially available AI tool for use as a diagnostic assist device displaying results within a viewer (CXR viewer; Annalise CXR V.1.2, Annalise-AI, Sydney, Australia) was evaluated.³² The AI tool deploys an underlying machine learning model, developed and validated by Seah *et al*,³¹ which consists of attribute and classification CNNs based on the EfficientNet architecture³³ and a segmentation CNN based on U-Net³⁴ with EfficientNet backbone. The model was trained on 821 681 de-identified CXR images from 284 649 patients originating from inpatient, outpatient and emergency settings across Australia, Europe and North America. Training dataset labelling involved independent triple labelling of all images by three radiologists selected from a wider pool of 120 consultant radiologists (none of whom were employed by the radiology network involved in this current study). The model was validated for 124 clinical findings in a multireader, multicase study.³¹ Thirty-four of these findings were deemed priority findings based on their clinical importance. The full list of 124 findings is available in online supplemental table 1. Ground truth labels for the validation study dataset were determined by a consensus of three independent radiologists drawn from a pool of seven fully credentialed subspecialty thoracic radiologists. The algorithm is publicly available at <https://cxrdemo.annalise.ai>. The AI model was used in line with pre-existing regulatory approval.³⁵

Technical integration

Prior to the start of the study, technical integration of the software into existing radiology practice systems and testing occurred over several weeks. First, an integration adapter was installed on the IT network of each radiology clinic and acted as a gateway between the internal IT infrastructure and the AI model. Auto-routing rules were established ensuring only CXR studies were forwarded to the integration adapter from the picture archiving and communication system (PACS). Following a successful testing period, the Annalise CXR viewer was installed and configured on workstations for the group of study radiologists.

Study participants

Eleven consultant radiologists working for a large Australian radiology network were invited to participate in the study through their local radiologist network. This group included general diagnostic radiologists who had completed specialist radiology training and passed all diagnostic radiology college examinations required for

consultant accreditation in Australia. All radiologists reported the minimum of 2000 chest radiographs per year (either within the radiology network or through other institutions) suggested to maintain competency.³⁶ No subspecialist chest radiologists were included.

The group included radiologists with a range of experience levels: five radiologists had 0–5 years post-training experience, three radiologists had 6–10 years of experience, and three radiologists had more than 10 years of experience. Radiologists were situated across four states in Australia and worked in public hospitals, private hospitals and community clinic settings. Both on site and remote reporting was included, in line with regular workflow. Prior to study commencement, each radiologist attended a training seminar and a one-on-one training session to fully understand the CXR viewer and its features. In addition, the participating radiologists were able to familiarise themselves with the viewer prior to commencement of data collection.

CXR case selection

In this multicentre real-world prospective study, all consecutive chest radiographs reported by the radiologists originating from inpatient, outpatient and emergency settings were included for a period covering nearly 6 weeks. The CXR cases were reported with the assistance of the AI tool in real-world clinical practice, using high-resolution diagnostic radiology monitors within the radiologists' normal reporting environment. As per usual workflow across a large radiology network spanning a geographically large area with many regional and remote clinics, both on-site and remote reporting of CXR cases was undertaken. A total of 106 sites contributed cases with case numbers varying from one case up to a maximum of 271 cases at the busiest site.

At least one frontal chest radiograph was required for analysis by the model, and cases that did not include at least one were excluded. Chest radiographs from patients aged younger than 16 years were excluded. Data from all sources was de-identified for analysis.

AI-assisted reporting

For each CXR case, radiologists produced their clinical report with access to clinical information, the referral and available patient history, in line with the normal workflow. The AI model analyses the CXR image(s) for each case but does not incorporate clinical inputs (such as previous imaging, referral information or patient demographic data) into the analysis. Model output was displayed to the radiologist in a user interface, linked to the image in the PACS, automatically launching when a CXR case was opened (figure 1).

A modified version of the commercially available AI software was employed for this study, which incorporated changes into the user interface to allow radiologists to provide feedback on model recommendations. No changes were made to the underlying model. An example of the modified model user interface is presented in

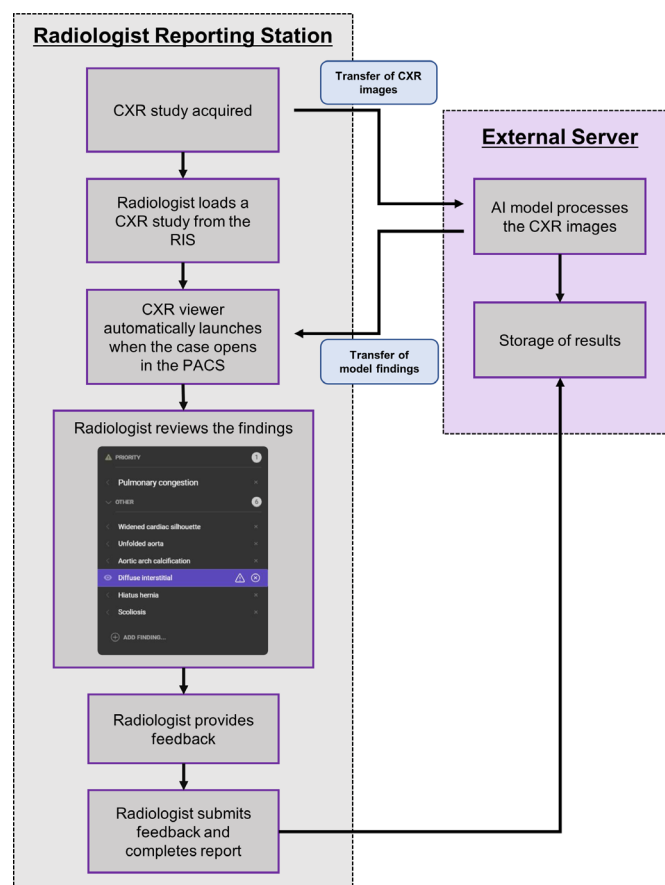


Figure 1 Flow diagram illustrating the AI-assisted reporting process described in this study. AI, artificial intelligence; CXR, chest X-ray; PACS, picture archiving and communication system; RIS, radiological information system.

figure 2. For each case, the model provided a list of suggested findings, listed as ‘priority’ or ‘other’, along with a confidence indicator. For a subset of findings, a region of interest localiser was overlaid on the image and the model indicated whether the finding was on the left or the right side, or both (see online supplemental table 1). The CXR viewer was configured to display its findings after the radiologists’ initial read of the case. For each case, radiologists were asked to review the CXR viewer’s findings and provide feedback within the viewer. The options presented to the radiologists in the viewer are listed in table 1.

The outcome measure of ‘significant impact on the report’ was the primary outcome measure. A significant change was described as the inclusion of findings recommended by the model, which altered the radiologists’ report in a meaningful way. As this varied by patient and clinical setting, it was left to the discretion of the radiologist. During the analysis of radiologist feedback, it was assumed that a change in patient management or further imaging recommendation would not occur without radiologists indicating a material change in the CXR report, and thus management and imaging questions were dependent on a significant change in the report. This was also patient-specific; for example, missing a pneumothorax in

Figure 2 Example of the modified user interface used by the participating radiologists in this study. The red box highlights the feedback options added to the interface for this study.

a ventilated patient with known pneumothorax would not have the same impact on patient management as a previously unknown pneumothorax in an outpatient. Free-text input describing missed findings or other relevant data were manually added after data collection was complete.

No formal adjudication of cases showing discrepancy between radiologist and model interpretation was performed. The study was not designed as a diagnostic accuracy validation. No review or ground truthing process was performed. Radiologists remained responsible for image interpretation and formulation of the report.

Poststudy survey

On completion of data collection, a poststudy survey was distributed to all participating radiologists to obtain feedback on the usefulness of the CXR viewer and how it affected their opinion of AI in radiology. A table of the survey questions is presented in online supplemental table 2.

Table 1 List of review options presented to the radiologist with each case

Review option	Description
Rejected clinical finding	A model-detected finding disputed by the radiologist
Missed clinical finding	A model-detected finding missed by the radiologist
Add additional findings	Finding(s) identified by the radiologist but not identified by the model
These findings significantly impacted my report	A yes/no binary question relating to the effect of the model output on the radiologist report
These findings may impact patient management	A yes/no binary question relating to the effect of the model output on patient management, as perceived by the reporting radiologist
These findings led to additional imaging recommendations	A binary yes/no question related to whether the radiologist recommended further imaging based on the model output

Statistics and data analysis

A 1% rate of significant changes in reports (the primary outcome measure) was deemed to be clinically significant prior to commencing the study. Based on estimations of the prevalence of missed critical findings on CXR, preliminary power calculations estimated that the number of cases required to detect at least a 1% rate of significant changes in reports was approximately 2000 cases in total, with alpha value 0.05 and desired power of 0.90. To account for any dropout in radiologists or cases, a target of 3000 cases was set for the study. Ten radiologists were recruited, with an 11th included for any unexpected participant drop out and to achieve this target in a reasonable time period.

A two-tailed binomial test was used to test the hypothesis that the rate of significant report change, patient management change or imaging recommendation change was at least 1%. To ensure that the sampling of CXRs reasonably approximated a random snapshot of the true population, radiologists in various states, with varying experience levels, as well as in different conditions of practice (community clinic vs hospital based) were selected. Additionally, the study was conducted prospectively which further aligned the structure of the sampled data with the expected structure of the population, justifying the choice of analysing the sample using a binomial test without adjustment for each radiologist.

Multivariate logistic regression using generalised linear mixed effect analysis was used to assess the effect of several possible confounders on the measured outcomes, including the number of critical clinical findings per case identified by the model, the inpatient/outpatient status of the patients, the experience level of the radiologists and the presence or absence of a lateral radiograph. The

Table 2 Demographics and results for the eleven radiologists involved in this study

Radiologist ID	No of years post-training	Cases reported (% outpatient)	Significant report impact (%)	Patient management changes (%)	Imaging recommendations (%)
1	19	136 (21.3)	1 (0.7)	1 (0.7)	0 (0.0)
2	1	325 (46.2)	4 (1.2)	0 (0.0)	1 (0.3)
3	4	230 (86.1)	20 (8.6)	14 (6.1)	10 (4.3)
4	6	375 (22.7)	3 (1.0)	0 (0.0)	1 (0.2)
5	4	186 (45.7)	22 (11.8)	9 (4.8)	8 (4.3)
6	20	333 (11.1)	3 (1.0)	2 (0.6)	1 (0.3)
7	3	312 (48.4)	15 (4.8)	8 (2.5)	1 (0.3)
8	26	408 (39.7)	10 (2.4)	5 (1.2)	4 (1.0)
9	9	214 (43.0)	6 (2.8)	2 (0.9)	2 (0.9)
10	6	159 (98.1)	1 (0.6)	1 (0.6)	1 (0.6)
11	5	294 (40.1)	7 (2.4)	1 (0.3)	0 (0.0)
Total		2972	92 (3.1)	43 (1.4)	29 (1.0)

Percentages (%) represent the associated value as a proportion of the total case number for that radiologist.

Wald test was applied to the derived regression coefficients to determine their significance.

Radiologists were grouped by experience level into 0–5 years postcompletion of radiology training, 6–10 years, and more than 10 years. A likelihood ratio test comparing a binomial logistic regression with categorical radiologist experience against a null model was performed to assess the hypothesis that the outcomes (significant changes in reports, management or imaging recommendation) were associated with experience.

A significance threshold of 0.05 was chosen, with the Benjamini-Hochberg procedure³⁷ applied to all reported outcomes to account for multiple hypothesis testing. Two clinically qualified researchers independently performed statistical analyses using different software. Calculations were performed in Excel 2016 with RealStatistics resource pack and cross-checked in Python V.3.7 using the Pandas V.1.0.5,³⁸ NumPy V.1.18.5,³⁹ SciPy V.1.4.1,⁴⁰ Scikit-Learn V.0.24.0,⁴¹ pymer4 V.0.7.1 (linked to R V.3.4.1, lme4 V.1.1.26)⁴² and Statsmodels V.0.12.1⁴³ libraries.

RESULTS

A total of 2972 cases were reported by 11 radiologists over a period of 6 weeks. These cases came from 2665 unique patients (52.7% male), with a median age of 67 (IQR 50–77). Information on radiologist experience, number of cases reported, source of cases and outcome measures for each radiologist are listed in table 2.

Of the 2972 cases, 1825 (61.4%) cases had lateral (as well as frontal) radiographs available for interpretation. 1709 (57.5%) cases were from an inpatient setting, and 1263 (42.5%) from an outpatient setting. The median number of findings per case was five (mean: 5.1, SD: 3.9), with a wide range in the number of findings per case (maximum=20). A total of 364 cases returned zero findings predicted by the model from the complete 124

findings list. A total of 1526 of the 2972 cases had one or more critical findings detected by the CXR viewer, with the critical findings in 1459 (96%) of these cases being confirmed by the radiologist. The number of critical findings per case is summarised in figure 3.

Influence of the AI model on radiologist reporting

Across all 2972 cases, there were 92 cases identified by radiologists as having significant report changes (3.1%), 43 cases of changed patient management (1.4%) and 29 cases of additional imaging recommendations (1.0%) as a result of exposure to the AI model output. When compared with the hypothesised 1% rate of change, the findings were significantly higher for changed reports ($p<0.01$) and changed patient management ($p<0.01$), and not significantly different for rate of imaging recommendation ($p=0.50$).

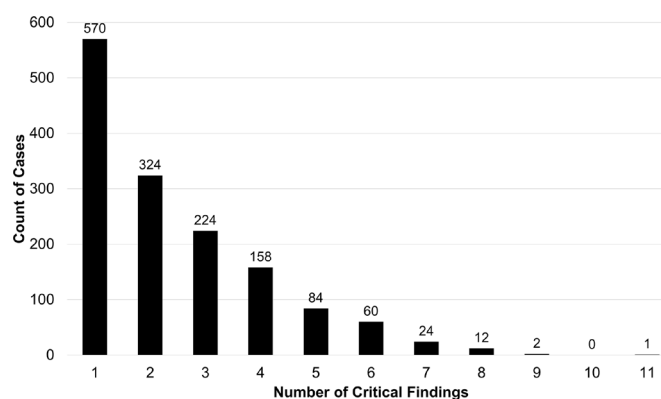


Figure 3 Counts of numbers of critical findings for the cases seen by the radiologist, defined as the number of critical findings agreed + the number of critical findings added. The number of cases which returned zero findings was 1513.

Table 3 Breakdown of the critical findings detected by the model and the level of radiologist agreement with each, including the number of findings reportedly missed by the model (and added by the radiologist) or missed by the radiologist

Critical finding	Displayed by model	Radiologist agreed with finding (%)	Radiologist rejected finding (%)	Added in by radiologist	Missed by radiologist
Acute aortic syndrome	2	2.0 (100.0)	0 (0.0)	0	0
Acute humerus fracture	5	5 (100.0)	0 (0.0)	0	0
Acute rib fracture	54	39 (72.2)	15 (27.8)	0	5
Cardiomegaly	1008	979 (97.1)	29 (2.9)	0	0
Cavitating mass	14	13 (92.9)	1 (7.1)	0	0
Cavitating mass internal content	6	5 (83.3)	1 (16.7)	0	0
Diffuse airspace opacity	13	13 (100.0)	0 (0.0)	0	0
Diffuse lower airspace opacity	153	148 (96.7)	5 (3.3)	0	0
Diffuse perihilar airspace opacity	45	45 (100.0)	0 (0.0)	0	0
Diffuse upper airspace opacity	2	2 (100.0)	0 (0.0)	0	0
Focal airspace opacity	341	321 (94.1)	20 (5.9)	0	2
Hilar lymphadenopathy	8	6 (75.0)	2 (25.0)	0	0
Inferior mediastinal mass	8	7 (87.5)	1 (12.5)	0	0
Loculated effusion	87	80 (92.0)	7 (8.0)	0	1
Lung collapse	11	10 (90.9)	1 (9.1)	0	0
Malpositioned CVC	85	78 (91.8)	7 (8.2)	0	1
Malpositioned ETT	52	43 (82.7)	9 (17.3)	0	0
Malpositioned NGT	39	31 (79.5)	8 (20.5)	0	0
Malpositioned PAC	13	9 (69.2)	4 (30.8)	0	0
Multifocal airspace opacity	125	120 (96.0)	5 (4.0)	0	1
Multiple pulmonary masses	43	38 (88.4)	5 (11.6)	0	0
Pneumomediastinum	5	5 (100.0)	0 (0.0)	1	0
Pulmonary congestion	220	215 (97.7)	5 (2.3)	1	0
Segmental collapse	292	290 (99.3)	2 (0.7)	0	1
Shoulder dislocation	1	0 (0.0)	1 (100.0)	0	0
Simple effusion	687	650 (94.6)	37 (5.4)	0	1
Simple pneumothorax	90	77 (85.6)	13 (14.4)	1	1
Single pulmonary mass	41	38 (92.7)	3 (7.3)	1	1
Single pulmonary nodule	105	95 (90.5)	10 (9.5)	3	5
Subcutaneous emphysema	53	51 (96.2)	2 (3.8)	0	1
Subdiaphragmatic gas	7	7 (100.0)	0 (0.0)	1	0
Superior mediastinal mass	37	32 (86.5)	5 (13.5)	0	0
Tension pneumothorax	11	7 (63.6)	4 (36.4)	0	0
Tracheal deviation	133	133 (100.0)	0 (0.0)	0	0
Total	3796	3594 (94.7)	202 (5.3)	8	20

Percentages (%) represent the associated value as a proportion of the total number of findings displayed by the model. CVC, Central venous catheter; ETT, Endotracheal tube; NGT, Nasogastric tube; PAC, Pulmonary artery catheter.

Agreement with model findings

Of the 2972 cases, 2569 had no findings rejected or added by the radiologists, indicating agreement with the model over all 124 possible findings in 86.5% of cases. A total of 306 (10.2%) cases had one finding rejected by the radiologist and 84 (2.8%) had two or more findings rejected by the radiologist. 202 (5.3%) critical findings detected by

the model were rejected by radiologists. The missed and rejected critical findings are detailed in [table 3](#).

Thirteen cases (0.5%) had findings (16 in total) added by the radiologists which they deemed were missed by the model, of which 8 were critical findings (see [table 3](#)). The remaining eight non-critical missed findings were atelectasis (four findings), cardiac valve prosthesis (two

Table 4 Factors affecting AI model influence on report, patient management, or imaging recommendation

Predictor	Change	ORs (adjusted CI)	P value	Benjamini-Adjusted threshold	Significance
No of critical findings	Report	1.306 (1.132 to 1.507)	0	0.0042	Yes
No of critical findings	Patient management	1.267 (1.056 to 1.521)	0.001	0.0083	Yes
No of critical findings	Imaging recommendation	1.319 (1.035 to 1.681)	0.004	0.0125	Yes
Lateral CXR	Imaging recommendation	6.495 (1.297 to 32.530)	0.005	0.0167	Yes
Lateral CXR	Patient management	2.158 (0.837 to 5.565)	0.061	0.0208	No
Lateral CXR	Report	1.542 (0.848 to 2.805)	0.105	0.025	No
Radiologist experience	Report	0–5 years: Baseline 6–10 years: 0.255 (0.043 to 1.521) >10 years: 0.305 (0.065 to 1.439)	0.120	0.0292	No
Radiologist experience	Patient management	0–5 years: Baseline 6–10 years: 0.165 (0.009 to 3.214) >10 years: 0.378 (0.054 to 2.654)	0.262	0.0333	No
Radiologist experience	Imaging recommendation	0–5 years: Baseline 6–10 years: 0.357 (0.034 to 3.783) >10 years: 0.380 (0.044 to 3.287)	0.516	0.0458	No
Inpatient/outpatient	Imaging recommendation	1.550 (0.613 to 3.919)	0.326	0.0375	No
Inpatient/outpatient	Report	0.794 (0.476 to 1.323)	0.358	0.0417	No
Inpatient/outpatient	Patient management	0.818 (0.408 to 1.640)	0.572	0.0500	No

Significance testing by the Benjamini-Hochberg algorithm to account for multiple hypotheses. ORs derived from stepwise logistic regression coefficients with CIs calculated with Benjamini-adjusted thresholds. Radiologist experience analysed as a categorical variable with derived from stepwise logistic regression coefficients with CIs calculated with Benjamini-adjusted thresholds. Radiologist experience analysed as a categorical variable with ORs representing effect of changing experience levels from the baseline (0–5 years) to a different level.

AI, artificial intelligence; CXR, chest X-ray.

findings), spinal wedge fracture (one finding) and peribronchial thickening (one finding).

Factors influencing reporting, management or imaging recommendation

The number of critical findings displayed by the model was significantly higher in cases where there was a change in report, patient management or imaging recommendation ($p<0.001$, $p=0.001$, $p=0.004$; table 4). The presence of a lateral projection image in the CXR case interpreted by the model was associated with a significantly greater likelihood of changes to imaging recommendation ($p=0.005$), but not to the report or patient management ($p=0.105$ and $p=0.061$, respectively).

Radiologists with fewer than 5 years consultant experience contributed 1347 cases, and indicated a rate of 5.0% for significant report change, 2.4% patient management change, and 1.5% recommendations for further imaging. These numbers were higher than for the radiologists with 6–10 years of experience (1.3%, 0.4%, 0.5%, respectively, over 748 cases) and also for radiologists with greater than 10 years of experience (1.6%, 0.9%, 0.6% over 877 cases). However, a likelihood ratio test applied to binomial logistic regression analysis indicated that the level of radiologist experience did not significantly influence the rate of change in report, patient management or imaging recommendation ($p=0.120$, $p=0.262$, and

$p=0.516$, respectively). Whether a patient was imaged as an inpatient or outpatient was not significantly associated with any change in report, patient management or imaging recommendation ($p=0.358$, $p=0.572$, $p=0.326$, respectively).

Survey results

The poststudy survey was completed by ten out of the eleven radiologists (figures 4 and 5). Notably, seven (70%) participants felt that their reporting time was slightly worse, however, when asked how satisfied they were with their reporting time, seven (70%) indicated that they were satisfied.

Nine out of 10 radiologists responded that their reporting accuracy was improved while using the CXR viewer, with 9 out of 10 (90%) participants being satisfied with accuracy of the CXR model's findings. Nine radiologists (90%) demonstrated an improved attitude towards the use of the AI diagnostic viewer by the end of the study and 9 (90%) demonstrated an improved attitude towards AI in general. No radiologists reported a more negative attitude towards the CXR viewer or towards AI in general.

DISCUSSION

We have previously shown that using the output of this comprehensive deep learning model improved

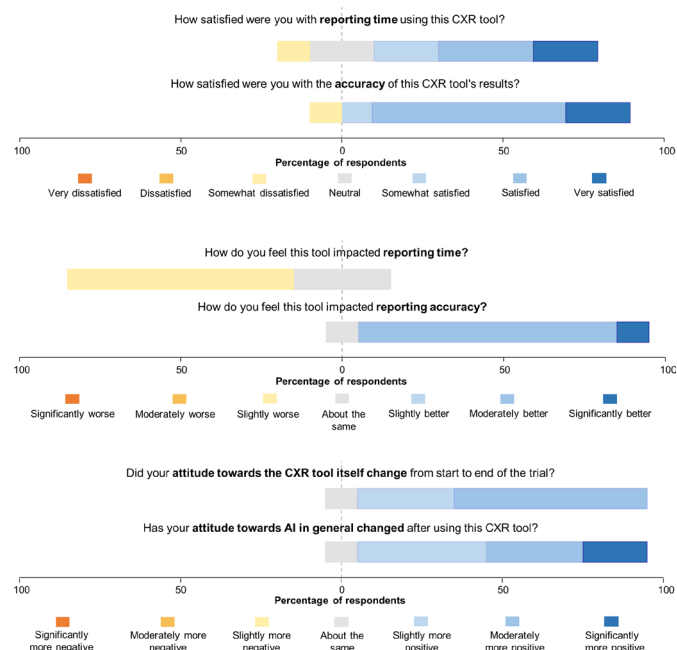


Figure 4 Diverging stacked bar chart depicting the first set of radiologist survey responses. CXR, chest X-ray.

radiologist diagnostic accuracy⁴⁴ in a non-clinical setting, but it is important to demonstrate that this improvement translates into meaningful change in a real-world environment. In this multicentre real-world prospective study, we determined how often the finding recommendations

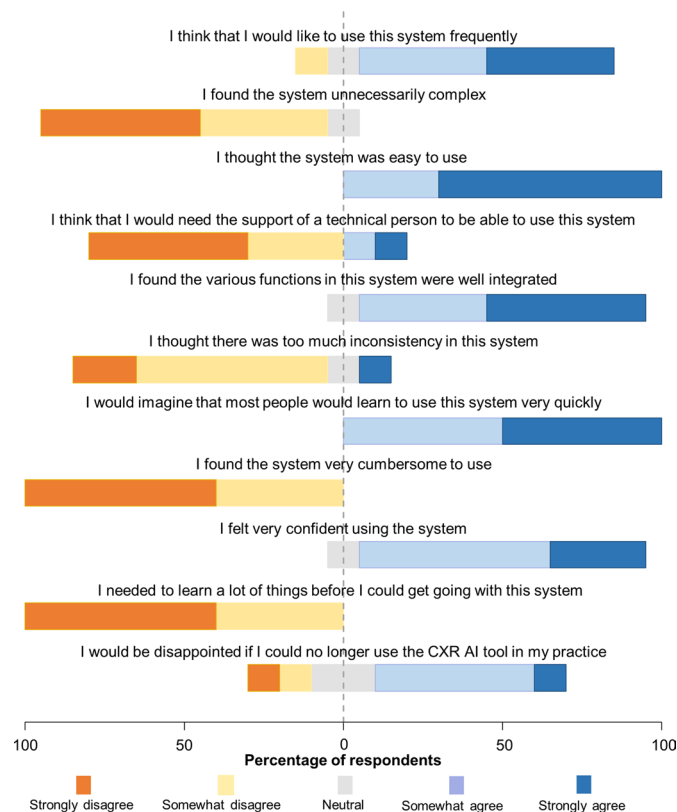


Figure 5 Diverging stacked bar chart visualising the second set of survey responses of the radiologists. AI, artificial intelligence; CXR, chest X-ray.

of the comprehensive deep learning model led to a material change in the radiologist's report, a change in the patient management recommendation, or a change in subsequent imaging recommendation. To the authors' knowledge, this is the first time that the impact of a comprehensive deep learning model developed to detect radiological findings on CXR has been studied in a real-world reporting environment. Other commercially available deep learning models able to detect multiple findings on CXR have been studied in the non-clinical setting, yielding encouraging results and outperforming physicians in the detection of major thoracic findings⁴⁵ as well as improving resident diagnostic sensitivity.⁴⁶ Other models have demonstrated diagnostic accuracy that is comparable to that of test radiologists.⁴⁷ Additionally, studies have yielded promising results for the use of models in population screening, particularly for tuberculosis, where several models have met the minimum WHO recommendations for tuberculosis triage tests.^{29 48}

We showed that radiologists agreed with all findings identified by the AI model in 86.5% of cases on a per case basis, while on a per finding basis, agreed with the critical findings identified by the model on 94.7% of findings. Notably, there was a significant change to the report in 3.1% of cases leading to changes in recommended patient management in 1.4% of cases, and changes to imaging recommendations in 1% of cases. Of note, 146 lung lesions (solitary lung nodule and solitary lung mass) were present in the dataset according to the model. Two lung lesions flagged by the model but missed by radiologists were recommended for additional imaging and changed management, subsequently diagnosed as lung carcinoma, highlighting the real-world value of integrating this type of system into the radiology workflow. However, four findings of lung nodule were flagged by the radiologists as missed by the model, indicating that the model alone is not intended to replace radiologist interpretation.

The significant impact of the CXR viewer on radiologist reporting and recommendations did however come at the cost of false positives, with 13% of cases having one or more model findings rejected by the radiologist. When this false positive rate is compared against the false positive rates per case reported in other studies investigating CXR models, which range from 14% to 88%,^{14 49 50} it is considered acceptable. Furthermore, these studies report false-positive rates for CXR models that only detect lung nodules, while in the current study this represents the false positive rate across 124 findings. Notably, on a per finding basis, only 5.3% of critical findings detected by the model were rejected by the radiologist. However, there were several outliers in the critical findings group that had noticeably higher rates of rejection, including acute rib fracture, hilar lymphadenopathy, malpositioned nasogastric tube (NGT)/pulmonary artery catheter (PAC), shoulder dislocation and tension pneumothorax. Several explanations for this are low sample size, the subjectivity of diagnosis (especially for hilar lymphadenopathy and tension features of pneumothorax), and heightened

model sensitivity at the expense of specificity. In particular, the rate of 'overcalling' of malposition of nasogastric tubes was related to both the threshold choice (favouring sensitivity given the critical nature of NGT malposition) and the limitation in the model output in distinguishing malpositioned NGTs from incompletely visualised NGTs. This limitation has subsequently been addressed with model modifications. Overall, this trade-off appears to be reasonable to the participating radiologists, who reported a high level of satisfaction with the model.

In this study, analysis of radiologists by experience level using logistic regression found no statistically significant relationship between experience level and increased changes to reports, patient management changes or imaging recommendations as a result of the model. Statistical analysis of the relationship between experience level and change in report was associated with a $p=0.12$, suggesting that, with further research, a significant relationship may be identified. It is expected that the inclusion of a larger group of radiologists may lead to a significant finding, as the association between experience and level of change has been noted in other studies. For example Jang *et al*, showed that less experienced radiologists benefited the most from diagnostic assistance in detecting lung nodules on CXR.¹⁴ In this study, 3 of the 11 radiologists contributed a higher than average incidence of the primary outcome of report change, and these were all less experienced radiologists compared with the cohort average experience level. While this may be due to variations in individual radiologist interpretation of 'significant report change', the consistency of experience level across these three radiologists suggests a relationship with experience level and tool impact.

The primary factor that influenced the likelihood of the model findings leading to a change in the report was the presence of critical findings in the model's recommendation. This is particularly notable because it indicates that the changes to the report are significant. They did not simply involve the inclusion of additional non-critical findings in the report, which may be interpreted as overestimating the impact of the model. The inpatient or outpatient status of a case was found not to significantly affect the likelihood of significant changes to the radiologists' report, to patient management, or to imaging recommendations.

The poststudy survey provided further insight into the impact that the CXR viewer had on participant reporting, in addition to the level of agreement and changes to the radiology report and patient management recommendations outlined above. The first notable response was that the CXR viewer may have negatively affected reporting times (although only mildly) for the majority of radiologists. This outcome was expected in this study setting because the radiologists were taking additional time to provide feedback on the model's recommendations for each case. Previous studies that surveyed radiologists reported that 74.4% thought AI would lower the interpretation time.⁵¹ It is notable that even with the negative

impact the model had on reporting time, the majority of radiologists (70%) were still satisfied with reporting time while using the CXR viewer, suggesting that the diagnostic improvements offered by the model were enough to offset the additional perceived reporting time. Additional insight from the survey suggested that very little training was required before radiologists felt comfortable using the tool. This is useful as education on AI has been a primary concern among clinicians, as a large proportion of radiologists report having little knowledge of AI.⁵²

Limitations and future research

The results presented in this study are self-reported by participating radiologists and are likely an underestimation of the model's actual impact. It is expected that radiologists would not report every instance in which they made an interpretive error. Another limitation is that there was no objective gold standard against which the radiologist and model interpretation could be measured. This is a small-scale study involving a limited sample size, conducted over several weeks. As a result, it lacks the statistical power to examine the benefit of the model on a finding-by-finding basis. In future, it would be beneficial to conduct a similar study with a larger sample size to allow for more powerful statistical analysis and examination of specific finding changes. Another useful next step would be to include a gold standard to determine the ground truth for the CXR findings, as this would prevent any under reporting which may occur with self-reported results, as well as enable the detection of false negatives as a result of the CXR viewer.

Although none of the cases evaluated in this study had been seen by the model previously, we note that one of the five data sources used for model training originated from the same radiology network. This, therefore, cannot be considered as true external evaluation. Further work in truly external institutions in the future are welcomed.

CONCLUSION

This study indicated that the integration of a comprehensive AI model capable of detecting 124 findings on CXR into a radiology workflow led to significant changes in reports and patient management, with an acceptable rate of additional imaging recommendations. These results were not affected by the inpatient status of the patient, and although approaching significance, the experience level of the radiologists did not significantly relate to the primary endpoint outcomes. In secondary endpoint outcomes, the model output showed good agreement with radiologists, and radiologists showed high rates of satisfaction with their reporting times and diagnostic accuracy when using the CXR viewer as a diagnostic assist device. Results highlight the usefulness of AI-driven diagnostic assist tools in improving clinical practice and patient outcomes.

Author affiliations

¹Annalise-AI, Sydney, New South Wales, Australia

²I-Med Radiology Network, Sydney, New South Wales, Australia

³Department of Radiology, Alfred Health, Melbourne, Victoria, Australia

⁴Australian Institute for Machine Learning, The University of Adelaide, Adelaide, South Australia, Australia

⁵School of Medicine, The University of Notre Dame Australia School of Medicine Sydney Campus, Darlinghurst, New South Wales, Australia

⁶Faculty of Engineering and IT, University of Technology Sydney, Sydney, New South Wales, Australia

Acknowledgements The authors would like to thank Mark Wilson, Marc Northrop, Nicolaus Carr and Trina Shnier for their valuable contributions to designing and managing the study.

Contributors CMJ contributed to conception and design of the work, acquisition of data, analysis and visualisation of data, interpretation of data, drafting of the work, project management. LD contributed to design of the work and acquisition of data. MRM contributed to conception and design of the work, interpretation and visualisation of data, development of diagrams, drafting of the work, and project management. CT and JS contributed to analysis and visualisation of data, interpretation of data, development of diagrams and drafting of the work. LO-R, AJ, QDB and NE contributed to interpretation of data. All authors revised the work critically for important intellectual content, gave final approval of the version to be published, and agreed to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved. CMJ is responsible for the overall content as guarantor.

Funding This work was supported by Annalise-AI. Annalise-AI supported this work through free provision of the model to participating radiologists for the duration of the study and financing of an external biostatistician. Award/Grant number is not applicable.

Competing interests CMJ is a radiologist employed by the radiology practice and a clinical consultant for Annalise-AI. LD, LO-R and NE are independent of Annalise-AI and have no interests to declare. MRM, JS, CT, AJ and QDB are employed by or seconded to Annalise-AI. Study conception, study design, ethics approval and data security were conducted independent of Annalise-AI.

Patient consent for publication Not applicable.

Ethics approval This study involves human participants and was approved by Uniting Care Health HREC, Wesley Hospital, Brisbane, Queensland, Australia (reference number: 2020.14.324). The requirement of patient consent was waived by the ethics committee due to the low-risk nature of the study.

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement All data relevant to the study are included in the article or uploaded as online supplemental information.

Supplemental material This content has been supplied by the author(s). It has not been vetted by BMJ Publishing Group Limited (BMJ) and may not have been peer-reviewed. Any opinions or recommendations discussed are solely those of the author(s) and are not endorsed by BMJ. BMJ disclaims all liability and responsibility arising from any reliance placed on the content. Where the content includes any translated material, BMJ does not warrant the accuracy and reliability of the translations (including but not limited to local regulations, clinical guidelines, terminology, drug names and drug dosages), and is not responsible for any error and/or omissions arising from translation and adaptation or otherwise.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

ORCID iDs

Michael R Milne <http://orcid.org/0000-0003-2082-5723>

Jarrel Seah <http://orcid.org/0000-0002-2305-7873>

REFERENCES

- Gillies RJ, Kinahan PE, Hricak H. Radiomics: images are more than pictures, they are data. *Radiology* 2016;278:563–77.
- Greene R, Williams FH, Francis H, Williams, MD: father of chest radiology in North America. *Radiographics* 1991;11:325–32.
- Schaefer-Prokop C, Neitzel U, Venema HW, et al. Digital chest radiography: an update on modern technology, dose containment and control of image quality. *Eur Radiol* 2008;18:1818–30.
- Lee CS, Nagy PG, Weaver SJ, et al. Cognitive and system factors contributing to diagnostic errors in radiology. *AJR Am J Roentgenol* 2013;201:611–7.
- Chotas HG, Ravin CE. Chest radiography: estimated lung volume and projected area obscured by the heart, mediastinum, and diaphragm. *Radiology* 1994;193:403–4.
- Berlin L. Accuracy of diagnostic procedures: has it improved over the past five decades? *AJR Am J Roentgenol* 2007;188:1173–8.
- Zaorsky NG, Churilla TM, Egleston BL, et al. Causes of death among cancer patients. *Ann Oncol* 2017;28:400–7.
- del Ciello A, Franchi P, Contegiacomo A. Missed lung cancer: when, where, and why? *Diagn Interv Radiol* 2017;23:118–26.
- Fazal MI, Patel ME, Tye J, et al. The past, present and future role of artificial intelligence in imaging. *Eur J Radiol* 2018;105:246–50.
- Jordan MI, Mitchell TM. Machine learning: trends, perspectives, and prospects. *Science* 2015;349:255–60.
- Hosny A, Parmar C, Quackenbush J, et al. Artificial intelligence in radiology. *Nat Rev Cancer* 2018;18:500–10.
- Erickson BJ, Korfiatis P, Akkus Z, et al. Machine learning for medical imaging. *Radiographics* 2017;37:505–15.
- Esteva A, Chou K, Yeung S, et al. Deep learning-enabled medical computer vision. *NPJ Digit Med* 2021;4:1–9.
- Jang S, Song H, Shin YJ, et al. Deep Learning-based automatic detection algorithm for reducing overlooked lung cancers on chest radiographs. *Radiology* 2020;296:652–61.
- Liang C-H, Liu Y-C, Wu M-T, et al. Identifying pulmonary nodules or masses on chest radiography using deep learning: external validation and strategies to improve clinical practice. *Clin Radiol* 2020;75:38–45.
- Hurt B, Kligerman S, Hsiao A. Deep learning localization of pneumonia: 2019 coronavirus (COVID-19) outbreak. *J Thorac Imaging* 2020;35:W87–9.
- Kim JY, Choe PG, Oh Y, et al. The first case of 2019 novel coronavirus pneumonia imported into Korea from Wuhan, China: implication for infection prevention and control measures. *J Korean Med Sci* 2020;35:e61.
- PRAS B, Attux R. A deep Convolutional neural network for COVID-19 detection using chest x-rays. Available: <http://arxiv.org/abs/2005.01578> [Accessed 23 Mar 2021].
- Rueckel J, Trappmann L, Schachtner B, et al. Impact of confounding thoracic tubes and pleural dehiscence extent on artificial intelligence pneumothorax detection in chest radiographs. *Invest Radiol* 2020;55:792–8.
- Sze-To A, Wang Z. tCheXNet: Detecting Pneumothorax on Chest X-Ray Images Using Deep Transfer Learning. In: Karray F, Campilho A, Yu A, eds. *Image analysis and recognition*. Cham: Springer International Publishing, 2019: 325–32.
- Hwang EJ, Hong JH, Lee KH, et al. Deep learning algorithm for surveillance of pneumothorax after lung biopsy: a multicenter diagnostic cohort study. *Eur Radiol* 2020;30:3660–71.
- Park S, Lee SM, Kim N, et al. Application of deep learning-based computer-aided detection system: detecting pneumothorax on chest radiograph after biopsy. *Eur Radiol* 2019;29:5341–8.
- Wang X, Yu J, Zhu Q, et al. Potential of deep learning in assessing pneumoconiosis depicted on digital chest radiography. *Occup Environ Med* 2020;77:597–602.
- Zhou S, Zhang X, Zhang R, . Identifying cardiomegaly in ChestX-ray8 using transfer learning. *Stud Health Technol Inform* 2019;264:482–6.
- Zou X-L, Ren Y, Feng D-Y, et al. A promising approach for screening pulmonary hypertension based on frontal chest radiographs using deep learning: a retrospective study. *PLoS One* 2020;15:e0236378.
- Pasa F, Golkov V, Pfeiffer F, et al. Efficient deep network architectures for fast chest X-ray tuberculosis screening and visualization. *Sci Rep* 2019;9:6268.
- Nash M, Kadavigere R, Andrade J, et al. Deep learning, computer-aided radiography reading for tuberculosis: a diagnostic accuracy study from a tertiary hospital in India. *Sci Rep* 2020;10:210.
- Heo S-J, Kim Y, Yun S, et al. Deep Learning Algorithms with Demographic Information Help to Detect Tuberculosis in Chest Radiographs in Annual Workers' Health Examination Data. *Int J Environ Res Public Health* 2019;16:250.
- Qin ZZ, Sander MS, Rai B, 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.
- Lakhani P, Sundaram B. Deep learning at chest radiography: automated classification of pulmonary tuberculosis by using Convolutional neural networks. *Radiology* 2017;284:574–82.

- 31 Seah JCY, Tang CHM, Buchlak QD, *et al.* Effect of a comprehensive deep-learning model on the accuracy of chest X-ray interpretation by radiologists: a retrospective, multireader multicase study. *Lancet Digit Health* 2021;3:e496–506.
- 32 Annalise.ai - Annalise CXR comprehensive medical imaging AI. Annalise.ai. Available: <https://annalise.ai/products/annalise-cxr/> [Accessed 23 Mar 2021].
- 33 Tan M, QV L. EfficientNet: rethinking model scaling for Convolutional neural networks. Available: <http://arxiv.org/abs/1905.11946> [Accessed 30 Mar 2021].
- 34 Ronneberger O, Fischer P, Brox T. U-Net: Convolutional networks for biomedical image segmentation. Available: <http://arxiv.org/abs/1505.04597> [Accessed 30 Mar 2021].
- 35 Annalise-AI Pty Ltd - Radiology DICOM image processing application software. Available: [https://www.ebs.tga.gov.au/servlet/xmlmillr6?dbid=ebs/PublicHTML/pdfStore.nsf&docid=F7ADAEBB76CEDD47CA2585E500424A43&agid=\(PrintDetailsPublic\)&actionid=1](https://www.ebs.tga.gov.au/servlet/xmlmillr6?dbid=ebs/PublicHTML/pdfStore.nsf&docid=F7ADAEBB76CEDD47CA2585E500424A43&agid=(PrintDetailsPublic)&actionid=1) [Accessed 25 Aug 2021].
- 36 Improving diagnostic pathways for patients with suspected lung cancer. Available: https://www.cancerresearchuk.org/sites/default/files/ace_lung_pathways_final_report_v1.4.pdf [Accessed 31 Aug 2021].
- 37 Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B* 1995;57:289–300.
- 38 McKinney W. Pandas: a foundational python library for data analysis and statistics. *Python High Performance Science Computer* 2011.
- 39 Harris CR, Millman KJ, van der Walt SJ, *et al.* Array programming with NumPy. *Nature* 2020;585:357–62.
- 40 Jones E, Oliphant T, Peterson P. SciPy: open source scientific tools for python 2001.
- 41 Pedregosa F, Varoquaux G, Gramfort A. Scikit-learn: machine learning in python. *Journal of Machine Learning Research* 2021 <https://hal.inria.fr/hal-00650905>
- 42 Jolly E. Pymer4: connecting R and python for linear mixed modeling. *Journal of Open Source Software* 2018;3:862.
- 43 InSeabold S, Perktold J. *Statsmodels: Econometric and statistical modeling with python*. Austin, Texas, 2010: 92–6.
- 44 Seah J, Tang C, Buchlak QD. Radiologist chest X-ray diagnostic accuracy performance improvements when augmented by a comprehensive deep learning model. *The Lancet Digital Health* 2021.
- 45 Hwang EJ, Park S, Jin K-N, *et al.* Development and validation of a deep Learning-Based automated detection algorithm for major thoracic diseases on chest radiographs. *JAMA Netw Open* 2019;2:e191095.
- 46 Hwang EJ, Nam JG, Lim WH, *et al.* Deep learning for chest radiograph diagnosis in the emergency department. *Radiology* 2019;293:573–80.
- 47 Singh R, Kalra MK, Nitiwarangkul C, *et al.* Deep learning in chest radiography: detection of findings and presence of change. *PLoS One* 2018;13:e0204155.
- 48 Khan FA, Majidulla A, Tavaziva G, *et al.* Chest X-ray analysis with deep learning-based software as a triage test for pulmonary tuberculosis: a prospective study of diagnostic accuracy for culture-confirmed disease. *Lancet Digit Health* 2020;2:e573–81.
- 49 Dellios N, Teichgraber U, Chelaru R, *et al.* Computer-Aided detection fidelity of pulmonary nodules in chest radiograph. *J Clin Imaging Sci* 2017;7:8.
- 50 Sim Y, Chung MJ, Kotter E. Deep Convolutional neural Network-based software improves radiologist detection of malignant lung nodules on chest radiographs. *Radiology*. [Epub ahead of print: 12 Nov 2019]. doi:10.1148/radiol.2019182465
- 51 Waymel Q, Badr S, Demondion X, *et al.* Impact of the rise of artificial intelligence in radiology: what do radiologists think? *Diagn Interv Imaging* 2019;100:327–36.
- 52 Collado-Mesa F, Alvarez E, Arheart K. The role of artificial intelligence in diagnostic radiology: a survey at a single radiology residency training program. *J Am Coll Radiol* 2018;15:1753–7.