

# Deep learning-based computer-aided diagnostic models *versus* other methods for predicting malignancy risk in CT-detected pulmonary nodules

## A systematic review and meta-analysis

Johnathan Watkins;<sup>a</sup> Abdullah Akram;<sup>b</sup> Janella Benemile;<sup>b</sup> Ruth Kathryn;<sup>b</sup> Filippo Croce;<sup>c</sup> Wahyu Wulaningsih.<sup>d</sup>

a. Optellum Ltd, Oxford, United Kingdom

b. Modamast Consulting Pte Ltd, Singapore

c. University Hospital of Wales, Cardiff, United Kingdom

d. The Royal Marsden, London, United Kingdom; Faculty of Life Sciences & Medicine, King's College London, London, United Kingdom

### Key points

**Question:** How effective are image-based, computer-aided diagnostic models that use deep learning methods to predict the malignancy risk of pulmonary nodules as compared with other methods used in clinical practice?

**Findings:** This systematic review and meta-analysis identified 20 observational studies (47,832 patients; 87,976 pulmonary nodules) from which pooled analyses found deep learning-based models to have a sensitivity of 0.87, specificity of 0.80, and summary area under the curve of 0.90 in predicting malignancy in pulmonary nodules. This was superior or comparable to other methods routinely used in clinical practice.

**Meaning:** Deep learning-based models are already being used in clinical practice in certain settings for nodule management. The results show their diagnostic performance justifies wider and more routine deployment.

### Abstract

**Importance.** There has been growing interest in the use of artificial intelligence (deep learning) to help achieve early diagnosis of prevalent diseases. None moreso than in lung cancer, where a combination of factors, including the high prevalence of nodules, the low prevalence of malignant nodules, and the indeterminacy of many nodules mean that it is fertile ground for the deployment of accurate, high-throughput deep learning (DL)-based tools.

**Objective.** To survey the landscape of externally validated DL-based CADx models, and assess their diagnostic performance for predicting the risk of malignancy in computed tomography (CT)-detected pulmonary nodules.

**Data sources.** An electronic search was performed in the MEDLINE (PubMed), EMBASE, Science Citation Index, Cochrane Library databases (from inception to 10 April 2023).

**Study selection.** Studies were deemed eligible if they were peer-reviewed experimental or observational articles that analysed the diagnostic performance of externally validated DL-based CADx models for the prediction of malignancy risk, with a direct comparison to models widely used in clinical practice.

**Data extraction and synthesis.** PRISMA guidelines were followed for the identification, screening, and selection process. A bivariate random-effect approach for the meta-analysis on the included studies was used. Quality Assessment of Diagnosis Accuracy Studies 2 (QUADAS-2) was used to assess risk of bias and applicability.

**Main outcomes and measures.** Main outcomes included sensitivity, specificity, and AUC.

**Results.** After screening, 20 studies were included, comprising 47,832 patients and 87,976 nodules, of which 4,147 were malignant. DL-based CADx models were 17.6% more sensitive than physician judgement alone, and 33.8% more than clinical risk models alone. They had a similar pooled specificity as physician judgement alone (0.80 [95% CI: 0.72–0.86] *v* 0.82 [95% CI: 0.76–0.87], respectively), but were 9.6% more specific than clinical risk models alone. Accounting for threshold effects, DL-based CADx models had superior summary areas under the receiver operating characteristic curve (sAUROC), with relative sAUROCs of 1.06 (95% CI: 1.03–1.09) and 1.22 (95% CI: 1.19–1.25), as compared to physician judgement and clinical risk models alone, respectively.

**Conclusions and relevance.** DL-based models show superior or comparable diagnostic performance when externally validated against widely used methods, such as the Brock and Mayo models. They have the potential to fulfil an unmet clinical-management need alongside experienced physician image readers. The included studies reported a high degree of heterogeneity, with threshold effects particularly prominent. Future research may consider more prospective studies and human-experimental studies.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

## **Introduction**

Five-year survival rates in the US for lung cancer fall from 73% at stage I to just 9–13% at stages III B and IV.<sup>1</sup> Hence, diagnosing lung cancer early is critical to reducing lung cancer mortality rates. Lung cancer is predominantly asymptomatic in its early stages, with pulmonary nodules often being the first sign.<sup>2</sup> Pulmonary nodules are discrete lung lesions, measuring  $\leq 30$  mm in size (average of axial diameters). Approximately 5% of these nodules are malignant.<sup>3</sup> Nodules, both benign and malignant, are detected in approximately 1.6 million people in the US each year.<sup>3</sup> The majority of these are detected by computed tomography (CT) scans, of which more than 12 million are performed in the US each year.<sup>4</sup> These facts combine to show that the detection and discrimination of pulmonary nodules are the most important means of diagnosing lung cancer early, and CT scans the most important modality.

Pulmonary nodules are easy to detect, but difficult to discriminate. Assessing the risk of indeterminate nodules – nodules without obvious signs of benignity (such as calcification) or obvious signs of malignancy (such as spiculation) – pose a particular challenge.<sup>5</sup> Multiple clinical risk prediction models are used to aid the physician in assessing nodules in order to diagnose lung cancer or refer high-risk cases for further, more invasive investigation. Most validated clinical risk models apply multivariable logistic regression to clinico-demographic predictors (such as age, family history, and smoking history) and radiological predictors (such as nodule size, morphology, and location). Externally validated models in clinical use include the Mayo model, Brock (or PanCan) model, and the Peking University People's Hospital (PKUPH) model.<sup>6–8</sup>

Recently, image-based artificial intelligence (AI) models that use deep learning (DL) methods have emerged to predict this malignancy risk.<sup>9</sup> One of the advantages of image-based computer-aided diagnostic (CADx) models is their ease and speed of use *versus* traditional risk tools, which require manual entry of input data. This manual entry leads to low adoption rates, which may fall even further as the number of patients entering the care pathway increases. As such, adding DL capability to image-based CADx models has the potential to fulfil an unmet clinical-management need, providing they produce comparable diagnostic performance.

The objective of this systematic review and meta-analysis was to survey the landscape of externally validated DL-based CADx models, and assess their diagnostic performance for predicting the risk of malignancy in CT-detected pulmonary nodules. Two previous systematic reviews have been conducted on studies of DL-based CADx models that diagnose lung cancer from pulmonary nodules,<sup>10,11</sup> but none have conducted pooled analysis on models that have been externally validated against models currently used in clinical practice. External validation in populations other than the populations used to develop the new model is essential to ensuring they are sufficiently robust to stand alongside existing, validated, and widely used models. This is the first systematic review to provide such a pooled analysis of studies, in that it considers only those studies that directly compare DL-based CADx models with physician judgement, clinical risk models, or Lung-RADS-based models.

## Methods

### Search strategy and screening

An electronic search was performed in the MEDLINE (PubMed), EMBASE, Science Citation Index, and Cochrane Library databases (from inception to 10 April 2023). Relevant English-language studies only were sought. Duplicate studies, case reports & series, non-systematic review articles, non-peer reviewed studies, non-human studies, meeting abstracts & proceedings, and unpublished studies were all excluded. The full set of keyword search terms may be found in eTable 1 in the Supplementary Material. Reference lists of key studies and domain-related systematic reviews were investigated for further studies that the search may have missed. This study followed the Preferred Reporting Items for Systematic Reviews and Meta-analyses (PRISMA) reporting guidelines.<sup>12</sup>

A total of 7,116 studies were found after removing duplicates (Figure 1). After screening out ineligible studies from their title and abstract, the full text of 69 studies were retrieved for final screening. Two reviewers (JB and RK) independently reviewed each text.

Studies that met the following criteria were included:

### Data extraction and quality assessment

From the included studies, the following information was extracted and tabulated independently by two reviewers (JB and RK): author; publication date; funding source; study type; study population country(ies); setting; outcome type(s); index test(s); reference test; number of participants in each validation dataset; number of nodules in each validation dataset; prevalence of malignancy among nodules; age range; sex; patient exclusions; proportion of smokers (current and former); nodule diameter range; median nodule diameter; nodule type(s); threshold (operating cut-off point); route of detection; and the outcomes reported. The data were subsequently checked for quality (AA, WW, and JW).

- **Study type.** Studies should be human experimental studies or human observational studies
- **Index tests.**
  - Index test being described and investigated should use AI or DL methods (DL-based model) – defined as the self-reported use of AI or DL – to classify or otherwise predict the risk of malignancy in pulmonary nodules detected via CT scans
  - External validation of the DL-based CADx model should be performed on datasets not used for the initial development of the DL-based CADx model and compared with other methods that are in widespread clinical use, the categories of which are:
    - Physician judgement (radiological image readers)
    - Clinical risk models (multivariable statistical models that use clinico-demographic and radiological variables as inputs)
    - Lung-RADS-based models (that allow computers or humans to automatically classify nodules based on nodule size, type, and stability over time)<sup>13</sup>
- **Reference tests.** Studies should confirm malignancy diagnosis via histopathological (biopsy) within the follow-up period after initial nodule detection
- **Target condition and population.** Study participants should be  $\geq 18$  years old, with at least one solid or part-solid pulmonary nodule (0–30 mm), as identified via CT scan (i.e. studies on ground-glass nodules [GGNs] only are excluded)
- **Outcomes.** Studies should report at least one of: sensitivity, specificity; areas under the curve (AUC); diagnostic odds ratios; or the number of true-positive, false-negative, true-negative, or false-positive cases (as confirmed by histopathological analysis)

Risk of bias and applicability was independently assessed by AA, JB, and JW using the Quality Assessment of Diagnostic Studies 2 (QUADAS-2) tool (eFigure 1 in the Supplementary Material).<sup>14</sup>

### Statistical analysis and quantitative synthesis

A meta-analysis of all included studies reporting diagnostic performance outcomes was conducted. For each of the index test types (DL-based CADx models; physician judgement alone; clinical risk models alone; Lung-RADS-based models alone), pooled estimates of sensitivity, specificity, and AUC were calculated using a bivariate,

random-effects approach, along with their respective 95% confidence intervals (CIs). Summary AUROC curves were plotted.

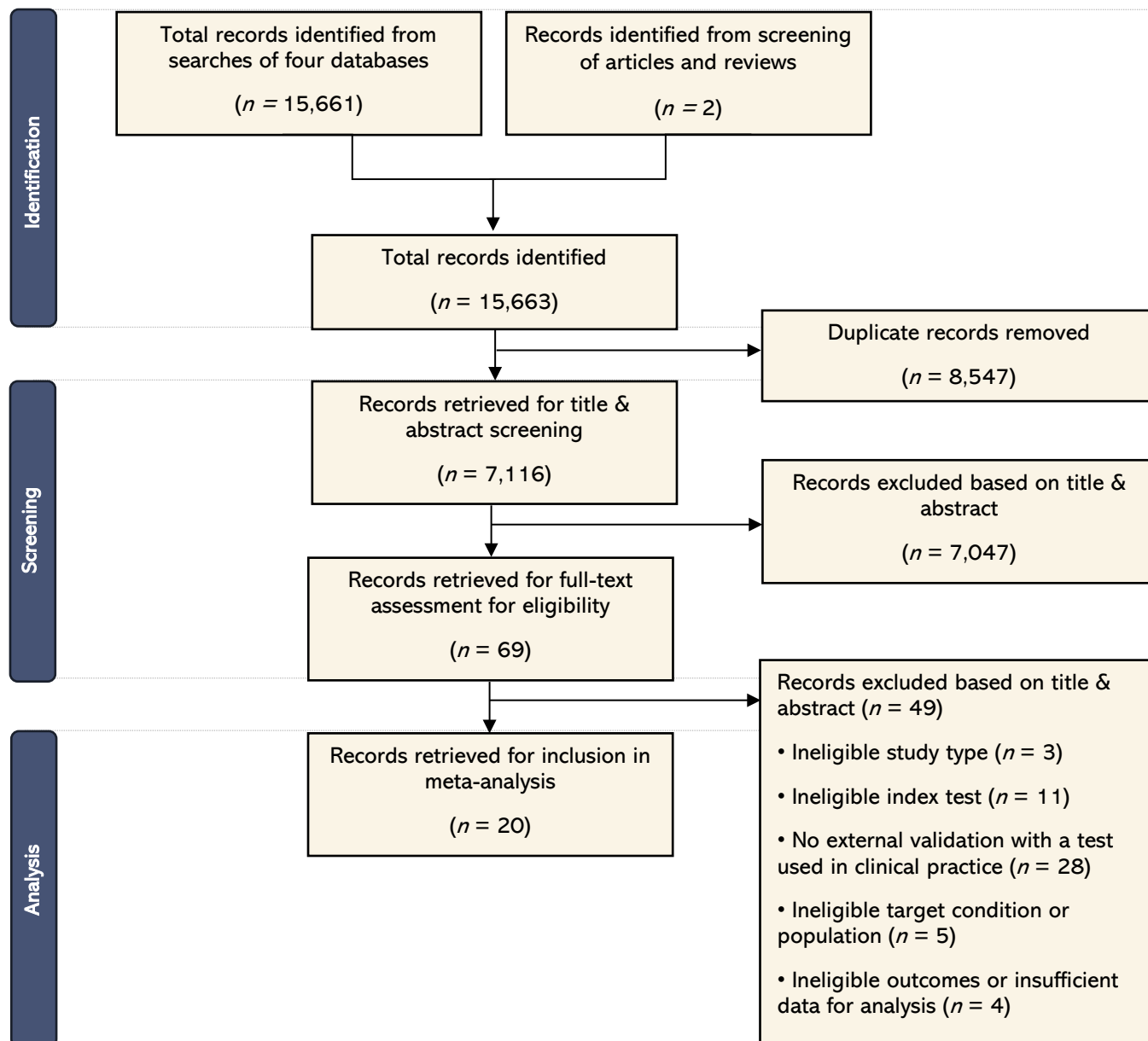
To assess heterogeneity and inconsistency among the studies,  $\chi^2$  statistic and  $I^2$  index values were calculated. An  $I^2$  value greater than 75% was considered indicative of substantial heterogeneity.

The Deeks funnel plot asymmetry test was performed to test for publication bias. A two-sided  $p < 0.05$  result was assumed to be statistically significant.

For the data available for extraction, we explored the sensitivity of the main pooled estimates to the study and population-outcome characteristics by conducting sub-group analyses. Once again, pooled sensitivity, specificity, and AUC estimates, along with  $\chi^2$  statistic and  $I^2$  index values, were calculated. The data were stratified by sub-group such as prevalence, route of detection, and median nodule size. This also helped uncover sources of heterogeneity.

Review Manager (RevMan) version 5.4 and Stata software, version 18 (StataCorp LLC) were used to conduct the statistical analyses.<sup>15,16</sup>

Figure 1. Literature search flow diagram



## Results

### Study characteristics

The literature search identified 20 studies for inclusion (Figure 1), comprising 38 validation datasets, representing 47,832 patients and 87,976 pulmonary nodules. Of these nodules, 4,147 were confirmed to be malignant (histopathological ground truth) within the follow-up period (on average, 24 months). A summary of the included studies is provided in Table 1.

All the studies save two were retrospective cohorts, with one prospective cohort and one case-control study.<sup>17,18</sup> The studies spanned continents and regions, with datasets taken from populations in North America (12 studies), Europe (7 studies), and East Asia (5 studies) (Table 1). Six of the studies used datasets taken from more than one country.

The outcome types in the studies included primarily assessed diagnostic performance. Some studies reported clinical utility measures, such as diagnostic re-classification.<sup>19,20</sup> However, due to the inconsistency in the outcomes reported, it was not possible to conduct a meta-analysis on clinical utility outcomes. The main outcomes sought were sensitivity and specificity, for which 17 of the 20 studies reported outcomes, and AUC, for which 13 studies reported outcomes (Table 1).

Seventeen DL-based CADx models were identified from the included studies. The commonest type of learning algorithm used in the DL-based models was a Convolutional Neural Network (CNN). Fourteen of the 17 models and 16 of the 20 included studies used a CNN algorithm as the basis for their DL malignancy prediction score.

For the external validation, the commonest validation was physician readers (17 of 38 datasets, from 12 studies). The majority of these readers were radiologists with  $\geq 3$  years' experience. Among the non-radiologists, thoracic surgeons comprised the majority.

With the clinical risk models, the Brock model was the commonest validation method (13 datasets from nine studies), followed by the Mayo model with eight datasets from four studies. This accords somewhat with the clinical risk models used in clinical practice. The Mayo model is considered the most externally validated model,<sup>21</sup> but the Brock model has been shown to perform better than the Mayo model in screening populations.<sup>22,23</sup>

The majority of studies considered participants in the 50–75 age bracket, with very few examples of younger participants. All studies included both female and male participants. For the five studies that conducted external validation on datasets from the US National Lung Screening Trial (NLST), participants were all current or former heavy smokers.

Two studies excluded calcified nodules, and two studies excluded GGNs.<sup>19,20,24</sup> The studies spanned the range of nodule sizes, with one study focussing only on malignancy risk prediction for large nodules  $>15$  mm.<sup>25</sup>

In terms of prevalence of nodular malignancy, datasets ranged from as low as 1% up to 67%, with an average prevalence across all datasets of 4.7%, which coincidentally roughly accords with real-world baseline prevalence in CT-detected nodules among US patients.<sup>3</sup> A number of studies adjusted their dataset populations so that the number of malignant nodules matched the number of benign nodules. Despite this matching, most incidentally detected nodule populations had prevalence at or higher than 20%, whereas most screening populations had prevalence under 20%. This is reflected in real-world populations, where screening populations tend to have lower rates of nodular malignancy as compared to nodules that have been incidentally detected.<sup>3</sup>

### Diagnostic performance

For the DL-based CADx models, sensitivity ranged from 0.37 (95% CI: 0.25–0.50) for a 0.98 (95% CI: 0.95–0.99) specificity,<sup>19</sup> to 1.00 (95% CI: 0.98–1.00) for a 0.28 (95% CI: 0.26–0.31) specificity (Figure 2A and Figure 3A).<sup>26</sup> Pooled receiver operating characteristic (ROC) analysis of all DL-based CADx model results gave a pooled AUC of 0.90 (95% CI: 0.87–0.93), sensitivity of 0.87 (95% CI: 0.81–0.92) and specificity of 0.80 (95% CI: 0.72–0.86) (Figure 2A, Figure 3A, and Figure 4A). Pooled studies had an  $I^2$  index of 95.62% (95% CI: 94.63–96.60) for sensitivity and 99.38% (95% CI: 99.31–99.45), corresponding to very high statistical heterogeneity. The Deeks funnel plot showed no significant asymmetry, indicating no evidence of publication bias (eFigure 2).

Separate pooled analysis for physician readers gave a pooled AUC of 0.85 (95% CI: 0.82–0.88), sensitivity of 0.74 (95% CI: 0.65–0.81) and specificity of 0.82 (95% CI: 0.76–0.87) (Figure 2B, Figure 3B, and Figure 4B). Pooled studies had an  $I^2$  index of 83.53% (95% CI: 75.18–91.88) for sensitivity and 98.01% (95% CI: 97.47–98.54), indicating high statistical heterogeneity.

**Table 1. Characteristics of included studies**

ID	Citation	Country	Deep learning-based index test	Non-deep learning-based index test	Number of participants in validation	Number of nodules in validation	Prevalence of malignant nodules / scans	Route of detection	Outcomes reported
01	Adams et al 2021 <sup>27</sup>	USA	Patient management informed by DL CNN model fed by 3D CNN full-volume detection and ROI detection output	<ul style="list-style-type: none"> <li>Physician readers (6 radiologists)</li> </ul>	3,197	3,197	<ul style="list-style-type: none"> <li>1.1%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>
02	Adams et al 2023 <sup>28</sup>	USA Canada	RevealAI-Lung supervised ML classifier CADx (mSI) and Lung-RADS	<ul style="list-style-type: none"> <li>LUNG-RADS criteria retrospectively applied</li> <li>Mayo model</li> <li>Brock model</li> </ul>	963	1190	<ul style="list-style-type: none"> <li>49.4%</li> <li>11.7%</li> <li>32.1%</li> </ul>	Screening; Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>
03	Ardila et al 2019 <sup>29</sup>	USA	DL CNN model fed by image-based 3D CNN full-volume detection and ROI detection output	<ul style="list-style-type: none"> <li>LUNG-RADS criteria retrospectively applied</li> <li>Physician readers (6 radiologists)</li> </ul>	7,531	7,531	<ul style="list-style-type: none"> <li>1.3%</li> <li>16.4%</li> <li>13.0%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>
04	Baldwin et al 2020 <sup>26</sup>	UK	DL CNN model (LCP)	<ul style="list-style-type: none"> <li>Brock model</li> </ul>	1,187	1,397	<ul style="list-style-type: none"> <li>17%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>
05	Chae et al 2020 <sup>30</sup>	China	DL CNN model (CT-LungNet)	<ul style="list-style-type: none"> <li>Physician readers (2 radiologists)</li> <li>Physician readers (4 non-radiologists)</li> </ul>	208	60	<ul style="list-style-type: none"> <li>50%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>
06	Chen et al 2021 <sup>18</sup>	China South Korea	DL non-CNN XGBoost-based model (PKU-M)	<ul style="list-style-type: none"> <li>Brock model</li> <li>Mayo model</li> <li>Physician readers (radiologist and 3 thorax surgeons)</li> </ul>	520	783	<ul style="list-style-type: none"> <li>55%</li> <li>63%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>
07	Chen et al 2022 <sup>24</sup>	China	DL CNN model (Deepwise Healthcare)	<ul style="list-style-type: none"> <li>Physician readers (2 radiologists)</li> </ul>	104	148	<ul style="list-style-type: none"> <li>57%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>
08	Çoruh et al 2021 <sup>31</sup>	Turkey	DL CNN model	<ul style="list-style-type: none"> <li>Physician readers (2 radiologists)</li> </ul>	158	158	<ul style="list-style-type: none"> <li>49%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>
09	Gao et al 2021 <sup>32</sup>	USA	Co-learning model fed by image-based DL model output and CDE	<ul style="list-style-type: none"> <li>PLCO<sub>M2012</sub></li> <li>Brock model</li> </ul>	23,652	64,898	<ul style="list-style-type: none"> <li>2%</li> <li>18%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>AUC</li> </ul>
10	Gao et al 2022 <sup>33</sup>	USA	Co-learning model fed by image-based DL model	<ul style="list-style-type: none"> <li>Mayo model</li> <li>Brock model</li> </ul>	387	387	<ul style="list-style-type: none"> <li>50%</li> <li>49%</li> <li>54%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>AUC</li> </ul>

output, biomarker output and CDE (M3Net)										
11	Huang et al 2018 <sup>17</sup>	USA	DL model	<ul style="list-style-type: none"> <li>Physician readers (3 radiologists)</li> </ul>	186	46	<ul style="list-style-type: none"> <li>43%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>	
12	Huang et al 2019 <sup>34</sup>	Canada	DL model (DeepLR precursor)	<ul style="list-style-type: none"> <li>LUNG-RADS criteria retrospectively applied</li> </ul>	2,294	2,294	<ul style="list-style-type: none"> <li>4%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
13	Hunter et al 2022 <sup>25</sup>	UK	DL- and radiomics-based model (LN-RPV)	<ul style="list-style-type: none"> <li>Physician readers (3 radiologists)</li> <li>Brock model</li> <li>Herder model</li> </ul>	252	369	<ul style="list-style-type: none"> <li>63%</li> <li>63%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
14	Jacobs et al 2021 <sup>35</sup>	USA Denmark Canada	<ul style="list-style-type: none"> <li>DL CNN model (grt123; Liao et al 2019)<sup>36</sup></li> <li>DL CNN model (JWDH)</li> <li>DL CNN model (Aidence)</li> </ul>	<ul style="list-style-type: none"> <li>Physician readers (11 radiologists)</li> </ul>	300	300	<ul style="list-style-type: none"> <li>33%</li> <li>33%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>AUC</li> </ul>	
15	Kim et al 2022 <sup>20</sup>	USA UK	DL CNN model (LCP)	<ul style="list-style-type: none"> <li>Physician readers (6 radiologists)</li> <li>Physician readers (6 pulmonologists)</li> </ul>	300	600	<ul style="list-style-type: none"> <li>50%</li> <li>50%</li> </ul>	Screening; Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
16	Liu et al 2020 <sup>37</sup>	China	DL CNN model	<ul style="list-style-type: none"> <li>Physician readers (6 radiologists)</li> </ul>	153	168	<ul style="list-style-type: none"> <li>67%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
17	Massion et al 2020 <sup>19</sup>	USA UK	DL CNN model (LCP)	<ul style="list-style-type: none"> <li>Brock model</li> <li>Mayo model</li> </ul>	2,505	926	<ul style="list-style-type: none"> <li>14%</li> <li>14%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
18	Trajanovski et al 2021 <sup>38</sup>	USA	DL CNN model (N-Net)	<ul style="list-style-type: none"> <li>Brock model</li> <li>Physician readers (6 radiologists)</li> </ul>	3,286	854	<ul style="list-style-type: none"> <li>2%</li> <li>21%</li> <li>12%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
19	Venkadesh et al 2021 <sup>39</sup>	USA Netherlands Belgium	DL CNN model	<ul style="list-style-type: none"> <li>Brock model</li> <li>Physician readers (11 physicians)</li> </ul>	599	1,235	<ul style="list-style-type: none"> <li>7%</li> <li>34%</li> <li>33%</li> </ul>	Screening	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> <li>AUC</li> </ul>	
20	Zhang et al 2019 <sup>40</sup>	China	DL 3D CNN model	<ul style="list-style-type: none"> <li>Physician readers (25 physicians)</li> </ul>	50	50	<ul style="list-style-type: none"> <li>50%</li> </ul>	Incidental	<ul style="list-style-type: none"> <li>Sensitivity</li> <li>Specificity</li> </ul>	



Pooled analysis for clinical risk models gave a pooled AUC of 0.74 (95% CI: 0.70–0.78), sensitivity of 0.65 (95% CI: 0.28–0.90) and specificity of 0.73 (95% CI: 0.36–0.93) (Figure 2C, Figure 3C, and Figure 4C). Pooled studies had an  $I^2$  index of 98.17% (95% CI: 97.64–98.70) for sensitivity and 99.47% (95% CI: 99.38–99.57), indicating very high statistical heterogeneity.

Lastly, pooled analysis for Lung-RADS-based models gave a pooled AUC of 0.67 (95% CI: 0.63–0.71), sensitivity of 0.57 (95% CI: 0.37–0.75) and specificity of 0.69 (95% CI: 0.51–0.82) (Figure 2D, Figure 3D, and Figure 4D). Pooled studies had an  $I^2$  index of 95.62% (95% CI: 93.15–98.09) for sensitivity and 99.75% (95% CI: 99.70–99.80), indicating very high statistical heterogeneity.

Sub-group analyses revealed that DL-based CADx models displayed significantly higher sensitivity on incidentally detected nodules than on screening-detected nodules, 0.91 (95% CI: 0.81–0.96) *versus* 0.84 (95% CI: 0.76–0.90), respectively (eFigure 3). This increased reliability in detecting lung cancer came at the cost of specificity with screening-detected nodules having 0.86 (95% CI: 0.79–0.90) as compared to incidentally detected nodules at 0.70 (95% CI: 0.55–0.81). There was no significant difference for physician readers between screening and incidental detection. However, clinical risk models showed significantly poorer specificity for incidentally detected nodules as compared to screening-detected nodules, 0.86 (95% CI: 0.77–0.92) *versus* 0.59 (95% CI: 0.11–0.95) (eFigure 4).

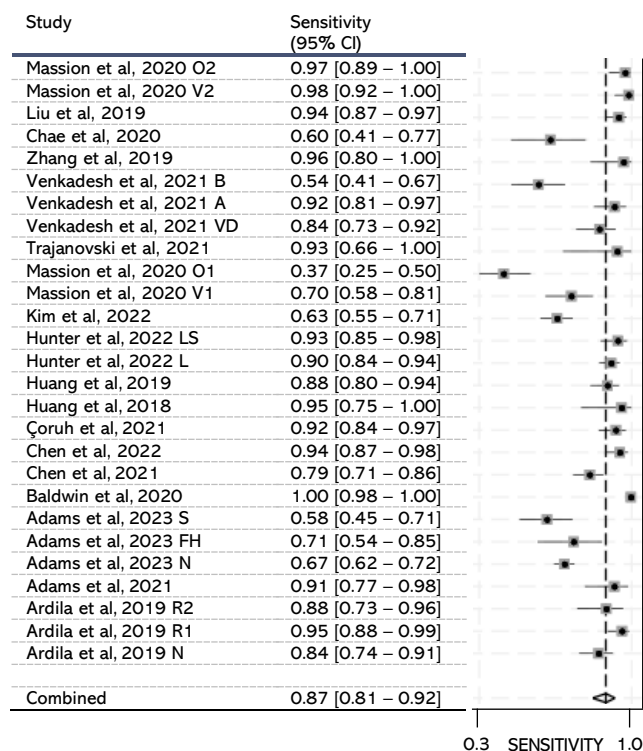
Further sub-group analyses were carried out on prevalence. To do this, we took the baseline prevalence of malignancy in CT-detected pulmonary nodules in the US, ~5%,<sup>3</sup> and multiplied it by a factor of 4, and used this as the threshold for classifying a study's prevalence as high or low. Our reasoning was that if a study's prevalence was 4 times as high as the baseline population prevalence, it could safely be considered high. Thus, our threshold for defining studies as having high or low prevalence was 20%. DL-based CADx models performed significantly better in low-prevalence studies than in high-prevalence studies: sensitivity of 0.90 (95% CI: 0.79–0.95) and specificity of 0.83 (95% CI: 0.69–0.91) as compared to sensitivity of 0.85 (95% CI: 0.77–0.91) and specificity of 0.76 (95% CI: 0.68–0.83), respectively.

#### **Quality assessment**

The results of the quality assessment using QUADAS-2 are shown in eTable 2 in the Supplementary Material. Overall, a low-to-moderate risk of bias was found in most studies.

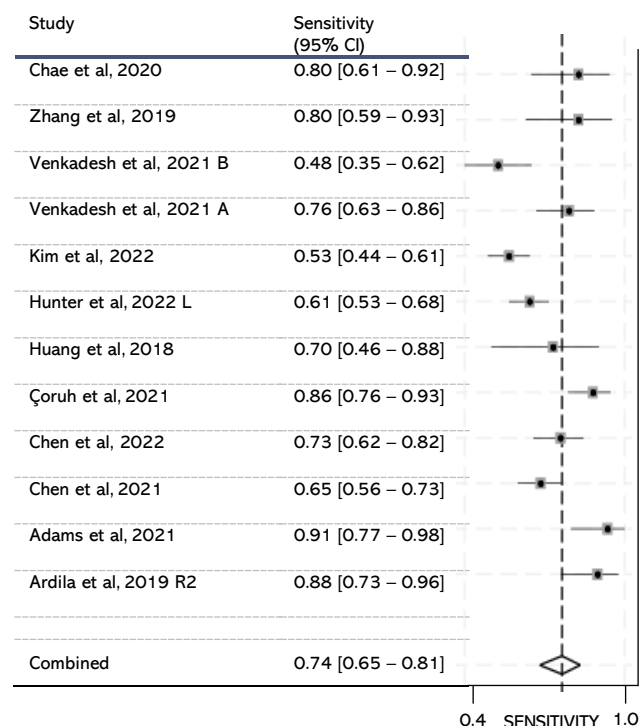
**Figure 2. Pooled sensitivity analyses of the included studies and their datasets**

**A Included studies for deep learning-based models**



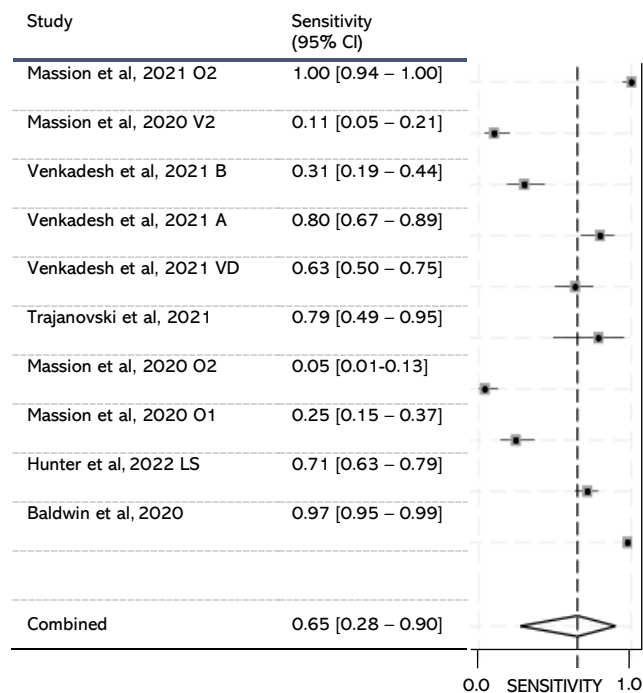
$Q = 593.33, df = 26.00, p = <0.05$   
 $I^2 = 95.62$  [95% CI: 94.63–96.60]

**B Included studies for physician reader models alone**



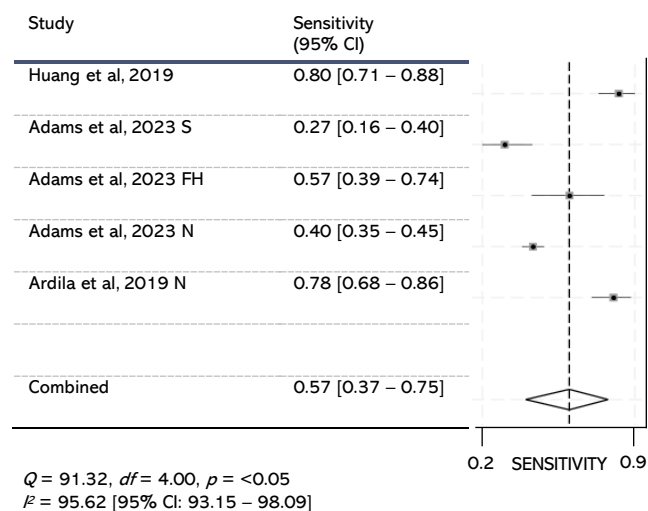
$Q = 66.78, df = 11.00, p = <0.05$   
 $I^2 = 83.53$  [95% CI: 75.18 – 91.88]

**C Included studies for clinical risk models alone**



$Q = 491.14, df = 9.00, p = <0.05$   
 $I^2 = 98.17$  [95% CI: 97.64 – 98.70]

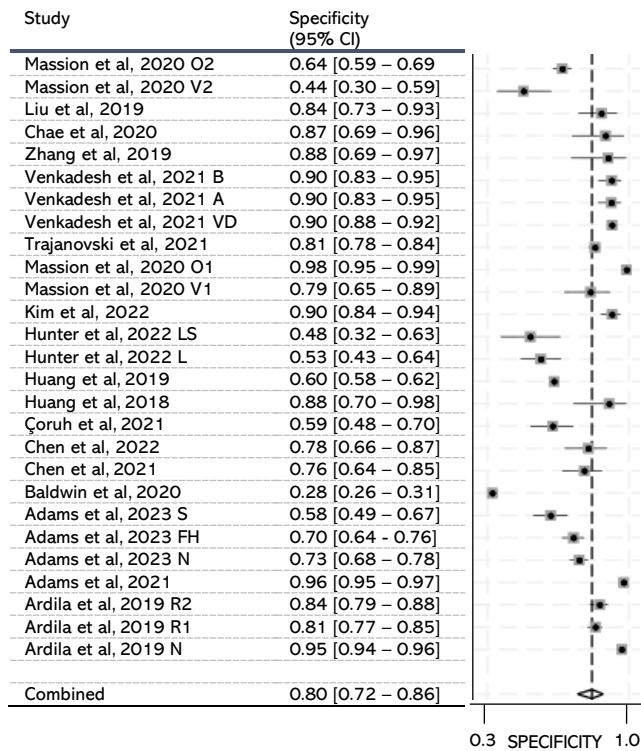
**D Included studies for Lung-RADS-based models alone**



$Q = 91.32, df = 4.00, p = <0.05$   
 $I^2 = 95.62$  [95% CI: 93.15 – 98.09]

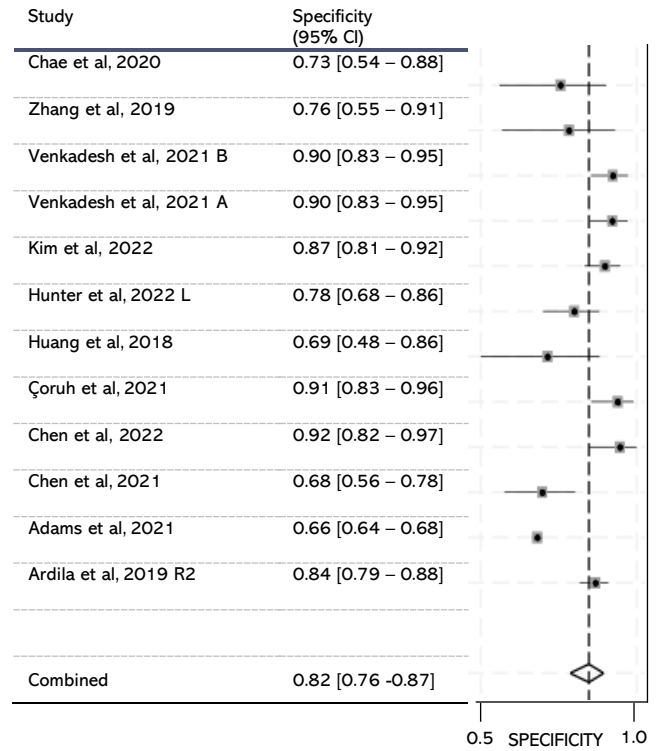
**Figure 3. Pooled specificity analyses of the included studies and their datasets**

**A Included studies for deep learning-based models**



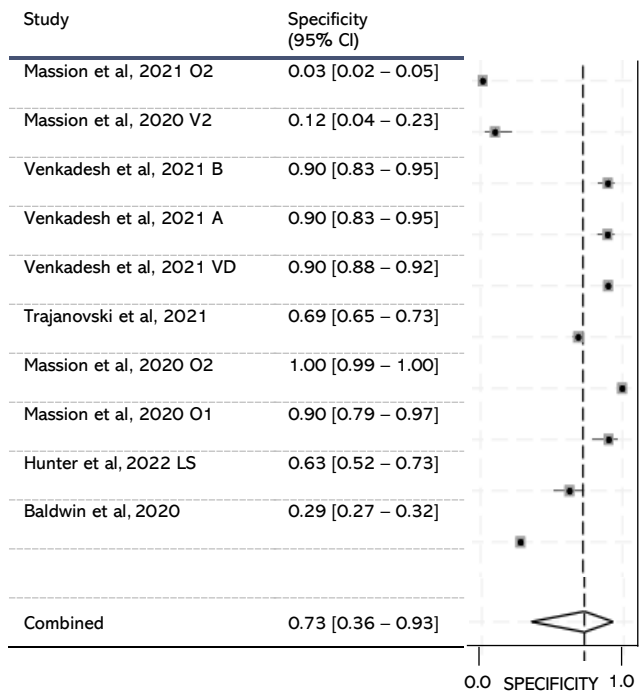
$Q = 4198.50$ ,  $df = 26.00$ ,  $p = <0.05$   
 $I^2 = 99.38$  [95% CI: 99.31 – 99.45]

**B Included studies for physician reader models alone**



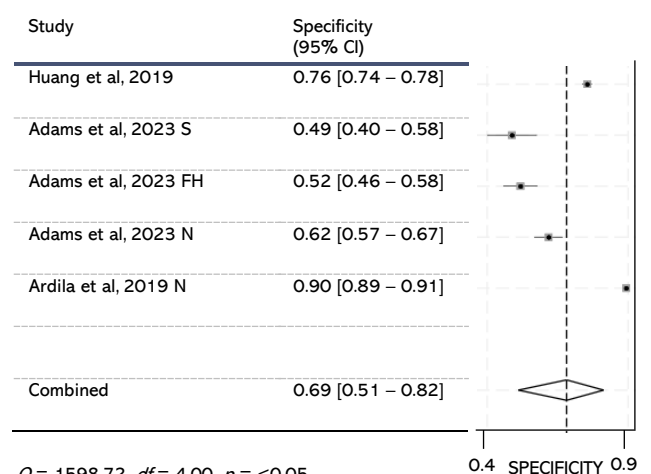
$Q = 552.32$ ,  $df = 11.00$ ,  $p = <0.05$   
 $I^2 = 98.01$  [95% CI: 97.47 – 98.54]

**C Included studies for clinical risk models alone**



$Q = 1713.10$ ,  $df = 9.00$ ,  $p = <0.05$   
 $I^2 = 99.47$  [95% CI: 99.38 – 99.57]

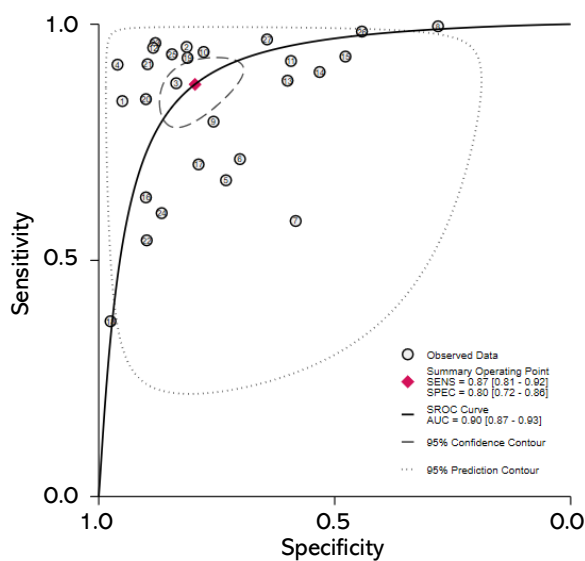
**D Included studies for Lung-RADS-based models alone**



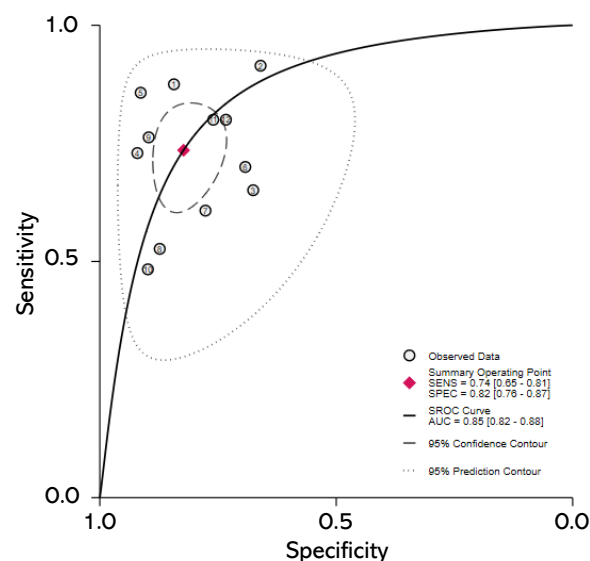
$Q = 1598.72$ ,  $df = 4.00$ ,  $p = <0.05$   
 $I^2 = 99.75$  [95% CI: 99.70 – 99.80]

**Figure 4. Summary ROC curve analyses of included diagnostic models**

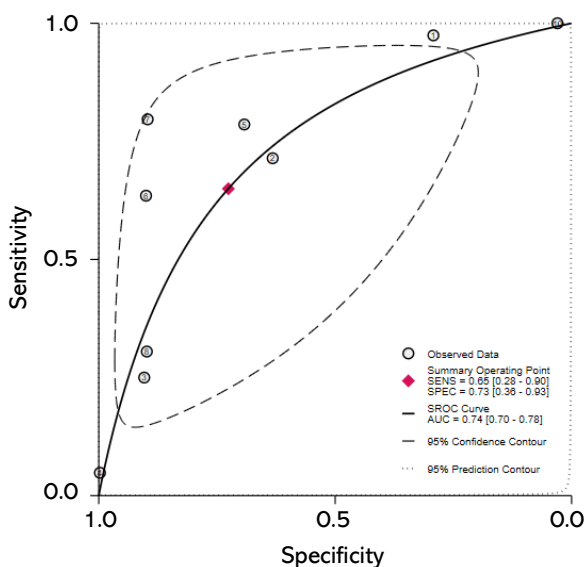
**A. Included studies for deep learning-based models**



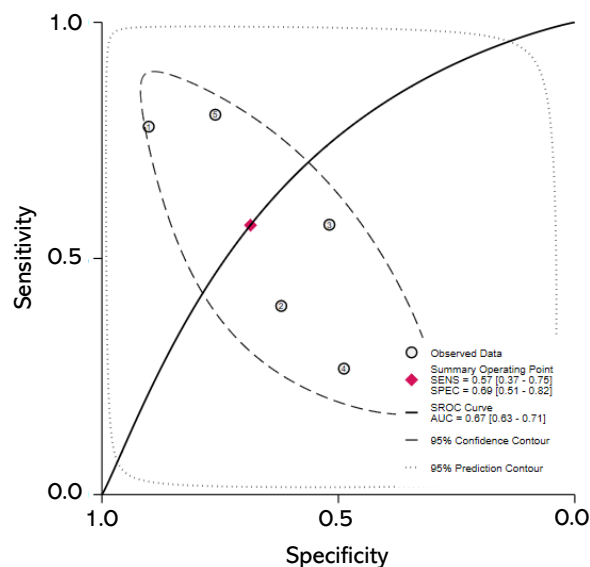
**B. Included studies for physician reader models alone**



**C. Included studies for clinical risk models alone**



**D. Included studies for Lung-RADS-based models alone**



## Discussion

A systematic review and meta-analysis that investigated the diagnostic performance of DL-based CADx models in predicting the risk of malignancy in pulmonary nodules *versus* methods currently used in clinical practice was performed. DL-based CADx models were significantly more sensitive than physician judgement alone, 17.6% (0.87  $\nu$  0.74), and clinical risk models alone, 33.8% (0.87  $\nu$  0.65). They had approximately the same pooled specificity as physician judgement alone (0.80 [95% CI: 0.72–0.86]  $\nu$  0.82 [95% CI: 0.76–0.87], respectively), but were 9.6% (0.80  $\nu$  0.73) more specific than clinical risk models alone.

Accounting for threshold effects, DL-based CADx models had significantly superior sAUROC, with relative sAUROCs of 1.06 (95% CI: 1.03–1.09), 1.22 (95% CI: 1.19–1.25), and 1.34 (95% CO: 1.31–1.37) as compared to physician judgement alone, clinical risk models alone, and Lung-RADS-based models alone, respectively.

This study attempted to exhaustively search the literature for all studies and models relevant to the research question. During screening, two of the commonest reasons for ineligibility were that the study did not conduct any direct external validation of the DL-based CADx model being analysed (28 studies with no direct external validation at the final screening stage) and that the study's model was for detection of pulmonary nodules, not classification or diagnosis of them (11 studies with ineligible index tests at the final screening stage) (Figure 1).

In order to evaluate their performance when applied to different populations, it is crucial to conduct external validation of DL-based risk prediction models in populations independent of those used during model development.<sup>41</sup> The majority of studies considered were on the development of the DL-based model, and not on external validation.

On the second commonest exclusion, CAD solutions for pulmonary nodule management can be broadly categorised into two types: computer-aided detection (CADE) and CADx (diagnosis). CADE involves a module designed to detect suspicious lung nodules and segment them for further analysis. Its purpose is to assist in the identification of potential abnormalities. CADx, on the other hand, goes beyond detection. It provides a nodule-level and, possibly, patient-level classification of the risk of malignancy. Only CADx is considered in this systematic review and meta-analysis. In broad terms, detection of pulmonary nodules is relatively easy. Distinguishing malign nodules from benign ones is not.<sup>2,5</sup>

Two previous systematic reviews have studied this issue,<sup>10,11</sup> albeit without the direct comparison between DL-based CADx models and external validation with methods used in clinical practice. Forte et al 2022 was the only one to conduct a meta-analysis, considering six studies, all of which are also included in this review and analysis. Pooled sensitivity and specificity were 0.94 (95% CI: 0.86–0.98) and 0.69 (95% CI: 0.51–0.83), respectively, both with significant heterogeneity, while sAUROC was 0.90 (95% CI: 0.86–0.92).<sup>10</sup> No quantitative comparisons against physician reader or clinical risk models alone were performed, and nor were any sub-group analyses performed. The authors noted that DL-based CADx models performed well and that as non-invasive methods, they could provide support to clinics in detecting and diagnosing lung cancer early.

## Limitations

Although these results strongly support the use of externally validated DL-based CADx models, two primary limitations were noted. Only observational studies could be found, and of these only one was prospective. No randomised controlled trials or other interventional studies were found. This is to be expected given that evidentiary requirements for diagnostic tools are not set as high as therapeutic interventions (drugs and biologics), and the difficulty in conducting such studies with diagnostic tools.<sup>42</sup> The second limitation was the high heterogeneity found among studies.

Sources of heterogeneity were investigated by conducting sub-group analyses. Prevalence was found to be a weak source of heterogeneity for clinical risk models in particular. Age range of the study population was another moderate source of heterogeneity across all models. However, the strongest source of heterogeneity was likely the threshold or operating cut-off point used by researchers in testing the models. The types of thresholds used varied considerably from study to study. They included fixing the specificity of models to 0.90,<sup>39</sup> to setting rule-out (definite benignity) malignancy scores at 0.05 (out of 1.0) or rule-in (definite malignancy) malignancy scores to 0.65 (out of 1.0).<sup>19,20</sup> Sensitivity to threshold effects was not investigated due to these inconsistent methods. However, the inclusion of AUC and its concordance with the sensitivity and specificity for each index test type helped alleviate this concern. Additionally, the low-to-moderate risk of bias found in most studies, and no significant publication bias demonstrated the findings were robust, in spite of the high heterogeneity.

Placing it further in context, the high heterogeneity in DL-based CADx models makes sense given the very different models under consideration, and the further work required on calibrating these models. However, as more validation in clinical practice occurs, and certain models become standard use in clinical practice – as has

happened with the Mayo and Brock models for clinical risk models – heterogeneity may reduce.<sup>6,7</sup> On this point, it important to note that use in clinical practice is important. Other clinical risk models, such as the Gurney model and the Bayesian Inference Malignancy Calculator (BIMC),<sup>43–45</sup> both of which use Bayesian analysis of clinico-demographic variables rather than logistic regression, have undergone external validation but are not used in clinical practice. Taking excellent diagnostic performance into clinical practice is the next step to ensuring improved patient outcomes are fully realised.

## **Conclusion**

These results demonstrate that DL-based CADx models have superior or comparable diagnostic performance as compared to methods currently used in clinical practice. The results support the use of DL-based CADx models alongside physician readers in clinical practice, especially for the management of incidentally detected nodules. While further research is required before they become an essential and routine part of the physician’s toolkit, recommendation for use in clinical practice as an option in the physician’s toolkit is justified by our findings.

## References

1. Woodard GA, Jones KD, Jablons DM. Lung Cancer Staging and Prognosis. *Cancer Treat Res*. 2016;170:47-75. doi:10.1007/978-3-319-40389-2\_3
2. Loverdos K, Fotiadis A, Kontogianni C, Iliopoulou M, Gaga M. Lung nodules: A comprehensive review on current approach and management. *Ann Thorac Med*. 2019;14(4):226-238. doi:10.4103/atm.ATM\_110\_19
3. Mazzone PJ, Lam L. Evaluating the Patient With a Pulmonary Nodule. *JAMA*. 2022;327(3):264. doi:10.1001/jama.2021.24287
4. Mahesh M, Ansari AJ, Mettler FA. Patient Exposure from Radiologic and Nuclear Medicine Procedures in the United States and Worldwide: 2009–2018. *Radiology*. 2023;307(1). doi:10.1148/radiol.221263
5. Paez R, Kammer MN, Massion P. Risk stratification of indeterminate pulmonary nodules. *Curr Opin Pulm Med*. 2021;27(4):240-248. doi:10.1097/MCP.0000000000000780
6. Swensen SJ, Silverstein MD, Ilstrup DM, Schleck CD, Edell ES. The probability of malignancy in solitary pulmonary nodules. Application to small radiologically indeterminate nodules. *Arch Intern Med*. 1997;157(8):849-855.
7. McWilliams A, Tammemagi MC, Mayo JR, et al. Probability of Cancer in Pulmonary Nodules Detected on First Screening CT. *New England Journal of Medicine*. 2013;369(10):910-919. doi:10.1056/NEJMoa1214726
8. Li Y, Wang J. A mathematical model for predicting malignancy of solitary pulmonary nodules. *World J Surg*. 2012;36(4):830-835. doi:10.1007/s00268-012-1449-8
9. Lee JH, Hwang EJ, Kim H, Park CM. A narrative review of deep learning applications in lung cancer research: from screening to prognostication. *Transl Lung Cancer Res*. 2022;11(6):1217-1229. doi:10.21037/tlcr-21-1012
10. Forte GC, Altmayer S, Silva RF, et al. Deep Learning Algorithms for Diagnosis of Lung Cancer: A Systematic Review and Meta-Analysis. *Cancers (Basel)*. 2022;14(16). doi:10.3390/cancers14163856
11. Wu Z, Wang F, Cao W, et al. Lung cancer risk prediction models based on pulmonary nodules: A systematic review. *Thorac Cancer*. 2022;13(5):664-677. doi:10.1111/1759-7714.14333
12. Page MJ, McKenzie JE, Bossuyt PM, et al. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *BMJ*. Published online March 29, 2021:n71. doi:10.1136/bmj.n71
13. Chelala L, Hossain R, Kazerooni EA, Christensen JD, Dyer DS, White CS. Lung-RADS Version 1.1: Challenges and a Look Ahead, From the *AJR* Special Series on Radiology Reporting and Data Systems. *American Journal of Roentgenology*. 2021;216(6):1411-1422. doi:10.2214/AJR.20.24807
14. Whiting PF, Rutjes AWS, Westwood ME, et al. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Ann Intern Med*. 2011;155(8):529-536. doi:10.7326/0003-4819-155-8-201110180-00009
15. StataCorp LLC. Stata Statistical Software: Release 18. Published online 2023.
16. The Cochrane Collaboration. Review Manager (RevMan). Published online 2020. Accessed June 6, 2023. [revman.cochrane.org](http://revman.cochrane.org)
17. Huang P, Park S, Yan R, et al. Added Value of Computer-aided CT Image Features for Early Lung Cancer Diagnosis with Small Pulmonary Nodules: A Matched Case-Control Study. *Radiology*. 2018;286(1):286-295. doi:10.1148/radiol.2017162725
18. Chen K, Nie Y, Park S, et al. Development and Validation of Machine Learning-based Model for the Prediction of Malignancy in Multiple Pulmonary Nodules: Analysis from Multicentric Cohorts. *Clin Cancer Res*. 2021;27(8):2255-2265. doi:10.1158/1078-0432.CCR-20-4007
19. Massion PP, Antic S, Ather S, et al. Assessing the Accuracy of a Deep Learning Method to Risk Stratify Indeterminate Pulmonary Nodules. *Am J Respir Crit Care Med*. 2020;202(2):241-249. doi:10.1164/rccm.201903-0505OC
20. Kim RY, Oke JL, Pickup LC, et al. Artificial Intelligence Tool for Assessment of Indeterminate Pulmonary Nodules Detected with CT. *Radiology*. 2022;304(3):683-691. doi:10.1148/radiol.212182

21. Choi HK, Ghobrial M, Mazzone PJ. Models to Estimate the Probability of Malignancy in Patients with Pulmonary Nodules. *Ann Am Thorac Soc*. 2018;15(10):1117-1126. doi:10.1513/AnnalsATS.201803-173CME
22. González Maldonado S, Delorme S, Hüsing A, et al. Evaluation of Prediction Models for Identifying Malignancy in Pulmonary Nodules Detected via Low-Dose Computed Tomography. *JAMA Netw Open*. 2020;3(2):e1921221. doi:10.1001/jamanetworkopen.2019.21221
23. White CS, Dharaiya E, Campbell E, Boroczky L. The Vancouver Lung Cancer Risk Prediction Model: Assessment by Using a Subset of the National Lung Screening Trial Cohort. *Radiology*. 2017;283(1):264-272. doi:10.1148/radiol.2016152627
24. Chen Y, Tian X, Fan K, Zheng Y, Tian N, Fan K. The Value of Artificial Intelligence Film Reading System Based on Deep Learning in the Diagnosis of Non-Small-Cell Lung Cancer and the Significance of Efficacy Monitoring: A Retrospective, Clinical, Nonrandomized, Controlled Study. *Comput Math Methods Med*. 2022;2022:2864170. doi:10.1155/2022/2864170
25. Hunter B, Chen M, Ratnakumar P, et al. A radiomics-based decision support tool improves lung cancer diagnosis in combination with the Herder score in large lung nodules. *EBioMedicine*. 2022;86:104344. doi:10.1016/j.ebiom.2022.104344
26. Baldwin DR, Gustafson J, Pickup L, et al. External validation of a convolutional neural network artificial intelligence tool to predict malignancy in pulmonary nodules. *Thorax*. 2020;75(4):306-312. doi:10.1136/thoraxjnl-2019-214104
27. Adams SJ, Mondal P, Penz E, Tyan CC, Lim H, Babyn P. Development and Cost Analysis of a Lung Nodule Management Strategy Combining Artificial Intelligence and Lung-RADS for Baseline Lung Cancer Screening. *Journal of the American College of Radiology*. 2021;18(5):741-751. doi:10.1016/j.jacr.2020.11.014
28. Adams SJ, Madtes DK, Burbridge B, et al. Clinical Impact and Generalizability of a Computer-Assisted Diagnostic Tool to Risk-Stratify Lung Nodules With CT. *Journal of the American College of Radiology*. 2023;20(2):232-242. doi:10.1016/j.jacr.2022.08.006
29. Ardila D, Kiraly AP, Bharadwaj S, et al. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nat Med*. 2019;25(6):954-961. doi:10.1038/s41591-019-0447-x
30. Chae KJ, Jin GY, Ko SB, et al. Deep Learning for the Classification of Small ( $\leq 2$  cm) Pulmonary Nodules on CT Imaging: A Preliminary Study. *Acad Radiol*. 2020;27(4):e55-e63. doi:10.1016/j.acra.2019.05.018
31. Gürsoy Çoruh A, Yenigün B, Uzun Ç, et al. A comparison of the fusion model of deep learning neural networks with human observation for lung nodule detection and classification. *Br J Radiol*. 2021;94(1123):20210222. doi:10.1259/bjr.20210222
32. Gao R, Tang Y, Khan MS, et al. Cancer Risk Estimation Combining Lung Screening CT with Clinical Data Elements. *Radiol Artif Intell*. 2021;3(6):e210032. doi:10.1148/ryai.2021210032
33. Gao R, Li T, Tang Y, et al. Reducing uncertainty in cancer risk estimation for patients with indeterminate pulmonary nodules using an integrated deep learning model. *Comput Biol Med*. 2022;150:106113. doi:10.1016/j.compbiomed.2022.106113
34. Huang P, Lin CT, Li Y, et al. Prediction of lung cancer risk at follow-up screening with low-dose CT: a training and validation study of a deep learning method. *Lancet Digit Health*. 2019;1(7):e353-e362. doi:10.1016/S2589-7500(19)30159-1
35. Jacobs C, Setio AAA, Scholten ET, et al. Deep Learning for Lung Cancer Detection on Screening CT Scans: Results of a Large-Scale Public Competition and an Observer Study with 11 Radiologists. *Radiol Artif Intell*. 2021;3(6):e210027. doi:10.1148/ryai.2021210027
36. Liao F, Liang M, Li Z, Hu X, Song S. Evaluate the Malignancy of Pulmonary Nodules Using the 3-D Deep Leaky Noisy-OR Network. *IEEE Trans Neural Netw Learn Syst*. 2019;30(11):3484-3495. doi:10.1109/TNNLS.2019.2892409
37. Liu J, Zhao L, Han X, Ji H, Liu L, He W. Estimation of malignancy of pulmonary nodules at CT scans: Effect of computer-aided diagnosis on diagnostic performance of radiologists. *Asia Pac J Clin Oncol*. 2021;17(3):216-221. doi:10.1111/ajco.13362



38. Trajanovski S, Mavroeidis D, Swisher CL, et al. Towards radiologist-level cancer risk assessment in CT lung screening using deep learning. *Comput Med Imaging Graph.* 2021;90:101883. doi:10.1016/j.compmedimag.2021.101883
39. Venkadesh KV, Setio AAA, Schreuder A, et al. Deep Learning for Malignancy Risk Estimation of Pulmonary Nodules Detected at Low-Dose Screening CT. *Radiology.* 2021;300(2):438-447. doi:10.1148/radiol.2021204433
40. Zhang C, Sun X, Dang K, et al. Toward an Expert Level of Lung Cancer Detection and Classification Using a Deep Convolutional Neural Network. *Oncologist.* 2019;24(9):1159-1165. doi:10.1634/theoncologist.2018-0908
41. Ramspek CL, Jager KJ, Dekker FW, Zoccali C, van Diepen M. External validation of prognostic models: what, why, how, when and where? *Clin Kidney J.* 2021;14(1):49-58. doi:10.1093/ckj/sfaa188
42. Mazumdar M, Zhong X, Ferket B. Diagnostic Trials. In: *Principles and Practice of Clinical Trials.* Springer International Publishing; 2021:1-28. doi:10.1007/978-3-319-52677-5\_281-1
43. Gurney JW. Determining the likelihood of malignancy in solitary pulmonary nodules with Bayesian analysis. Part I. Theory. *Radiology.* 1993;186(2):405-413. doi:10.1148/radiology.186.2.8421743
44. Soardi GA, Perandini S, Motton M, Montemezzi S. Assessing probability of malignancy in solid solitary pulmonary nodules with a new Bayesian calculator: improving diagnostic accuracy by means of expanded and updated features. *Eur Radiol.* 2015;25(1):155-162. doi:10.1007/s00330-014-3396-2
45. Gurney JW, Lyddon DM, McKay JA. Determining the likelihood of malignancy in solitary pulmonary nodules with Bayesian analysis. Part II. Application. *Radiology.* 1993;186(2):415-422. doi:10.1148/radiology.186.2.8421744

## Supplementary Material

**eTable 1.** Keyword search terms

Term	Component	Operator
Computer-aided	Index test(s) set #1	AND
Computer aided		OR
Computer-assisted		OR
Computer assisted		OR
CADx		OR
Artificial intelligence		OR
Machine intelligence		OR
Co-learning		OR
Colearning		OR
Machine learning		OR
Deep learning		OR
Predict*	Index test(s) set #2	AND
Diagnos*		OR
Classif*		OR
Estimat*		OR
Evaluat*		OR
Risk		OR
Compute* tomograph*	Index test(s) set #3	AND
Axial tomograph*		OR
CT scan*		OR
CAT scan*		OR
Cancer*	Target condition set #1	AND
Carcinoma*		OR
Neoplas*		OR
Tumour*		OR
Tumor*		OR
Malignan*		OR
Nodule*		OR
Lung*	Target condition set #2	AND
Pulmonary		OR