

## Deep learning driven image-based cancer diagnosis

Jimmy Joseph \*

*Independent Researcher, USA*

World Journal of Advanced Engineering Technology and Sciences, 2025, 16(02), 422-442

Publication history: Received on 18 July 2025; revised on 26 August 2025; accepted on 28 August 2025

Article DOI: <https://doi.org/10.30574/wjaets.2025.16.2.1311>

### Abstract

Lung cancer remains the leading cause of cancer-related mortality worldwide, primarily due to late-stage detection, which limits treatment availability. Early detection with low-dose CT screening increases the chance of survival, yet interpretation of CT scans for subtle malignant nodules is both difficult and time-consuming for radiologists. We present a new CAD system that utilizes deep learning techniques in the clinical imaging workflow for the detection of early-stage lung cancer. Our method adopts a hybrid CNN-Transformer model, which is designed to simultaneously learn local nodule properties and global contextual patterns based on attention mechanisms. The system is intended to be clinically feasible and has been constructed with the capability to interface seamlessly with hospital PACS and DICOM for image retrieval and storage. We address the most common forms of bias in imaging AI by applying data normalization, class-balanced data augmentation, and a bias-mitigating training scheme to generalize across different imaging devices and patient subpopulations. Experiments on public datasets (LIDC-IDRI for lung nodules and DeepLesion for diverse lesions) and a simulated multi-institution dataset show that the model achieves high accuracy and generalization. Notably, the hybrid CNN-Transformer demonstrates better performance across all baseline classifiers (including conventional radiomics SVM, CNN, or Vision Transformer models), with an overall AUC of 0.97 (95% confidence interval: 0.94–0.99), a sensitivity of 95.1%, and a specificity of 96.5% on the independent test set. Statistical tests (DeLong's test for AUC) show that our method is significantly better than previous approaches ( $p < 0.01$ ). We provide complete technical material, including model architecture, training pseudocode, and all evaluation scores with confidence intervals and p-values. Furthermore, we demonstrate qualitative results with Grad-CAM visualizations of image regions most pertinent to predicting classification labels, thereby providing interpretability for clinicians. We also address responsible AI practices such as bias audits at the subgroup level (e.g., age, sex, scanner type), documentation via model cards for transparency, adversarial robustness testing, and compliance with privacy regulations (de-identification and Data Protection Impact Assessment). Our results indicate that the bias-aware hybrid CNN-Transformer is a powerful yet reliable solution for image-based lung cancer diagnosis and is a promising enabler of earlier diagnosis and better patient outcomes, while complying with ethical and security considerations of medical AI.

**Keywords:** Deep Learning; Lung Cancer Diagnosis; Hybrid CNN-Transformer; Medical Imaging (CT, PACS/DICOM); Bias Mitigation and Explainability

### 1. Introduction

Lung cancer still represents a devastating public health burden worldwide, with an estimated 2.2 million new cases and 1.8 million deaths per year [1], [2]. High mortality results mainly from diagnosis at a late stage of the disease because small early-stage lung tumors often produce no obvious symptoms and cannot be detected by physical examination or, at this time, on a chest X-ray. CT images provide better sensitivity for small pulmonary nodules, and low-dose CT screening has been found to decrease lung cancer-related mortality through early detection and intervention [3]. But routine CT screenings generate a large volume of images that overwork radiologists, who must painstakingly look for

\* Corresponding author: Jimmy Joseph.

small nodules. Classical image interpretation is labor-intensive and subjective and includes major observer variability. There is, therefore, a strong clinical demand for automated image-processing tools that can aid physicians by identifying potentially suspicious lung nodules from CT images and marking them for further inspection.

The development of artificial intelligence (AI, in particular, deep learning) has already revolutionized computer-aided diagnosis (CAD) of medical images. Deep learning in the form of CNNs has been enormously successful in learning hierarchical features for medical images and improving detection accuracy for lung nodules. CNN-based CAD algorithms typically show superior performance compared to traditional machine learning workflows (such as rule-based techniques or Support Vector Machines on engineered radiomic features) in terms of sensitivity and specificity[6]. More recently, Transformer architecture with self-attention has achieved promising results in image analysis by learning long-range dependencies and global context[8]. When ViTs are used for medical images, they can supplement CNNs in modeling relationships between distant parts of an image—for example, connecting a nodule with contralateral lung findings or assessing the overall lung environment. However, pure Transformers usually have a very strong appetite for extremely large datasets, and they lack CNNs' inductive biases to locality, which makes them data-hungry and less stable on small medical datasets. Hybrid architectures that combine CNNs and Transformers have recently come into the picture to make the best use of both. Such hybrids, for example, can leverage CNN layers to capture low-level features (edges, textures, nodule shapes), while using Transformer layers to reason about global structure (i.e., the anatomical context of a nodule inside the lung).

However, there are still considerable pitfalls to be addressed before AI systems for lung cancer detection can be universally used. A primary concern is the “black-box” property of many deep models; clinicians are reluctant to blindly accept algorithmic predictions without a degree of interpretability or explanation of how decisions are made [18]. Solving this requires explainable AI techniques such as Gradient-weighted Class Activation Mapping (Grad-CAM), which highlights the image areas most influential on the model's output prediction so that a radiologist can validate whether the AI's attention is on medically relevant areas (e.g., the nodule itself and not an unrelated artifact). Generalizing the model across different patient populations and imaging conditions is another difficulty. Bias in machine-learning-derived predictions can occur from biased training data or spurious correlations, which may cause lowered accuracy in predictions for under-represented subgroups. For example, an AI may inadvertently learn to depend on slight scanner-specific noise or patient demographics present in the training set. Such behaviors may result in fairness problems, such as worse performance in certain age groups, genders, or ethnicities. Algorithmic fairness and trustworthiness are ethical and regulatory requirements of medical AI.

In this paper, we introduce a joint CNN-Transformer deep learning model for lung cancer CT image detection that is accurate, interpretable, and bias-mitigating. Our contributions can be summarized as follows:

- **Technical:** We propose a new architecture embedding two complementary mechanisms, i.e., CNN and Transformer with attention-based mechanisms, and achieve state-of-the-art performance on the task of lung nodule classification (malignant vs benign). We formulate the model components and learning objectives mathematically, comprising convolutions, self-attention, and a custom loss function that extends the bias term with a regularizer.
- **Clinical integration:** We sketch a deployment pipeline to integrate our model with hospital PACS via DICOM standards, allowing CT exams to be automatically retrieved and the AI results (e.g., probability scores or annotated heatmaps) inserted back into the radiologist's viewing workflow[23]. This design enables our AI system to serve as a second reader alongside radiologists to indicate suspicious areas, without affecting the standard of care.
- **Bias analysis:** We evaluate fairness and generalizability from a responsible AI viewpoint through extensive bias analysis and ablation studies[9],[15]. We propose measures such as standardizing intensity (to suppress artifacts depending on each scanner) and balancing the batches (keeping the number of classes or groups as balanced as possible). We also conduct audits by various sub-cohorts (e.g., by sex, smoking history) to verify consistent sensitivity. We also consider security concerns, such as how adversarial attacks could inject or remove subtle signs of cancer in a CT image to fool the model or human reader, and what can be done to protect the system.

To demonstrate the effectiveness of our approach, we performed experiments on both public datasets and a multi-institution simulation dataset. The public LIDC-IDRI dataset was used for training and internal validation, which includes more than 1,000 lung CTs with nodules annotated by radiologists[4]. We also tested generalization to various abnormality types by annotating and testing on the NIH DeepLesion dataset (we used the 32,000-lesion subset included in NIH DeepLesion for proper comparison), which consists of over 32,000 diverse lesions from 4,427 patients' CT images. We also generated a synthetic dataset with manipulated biases: artificially adding noise patterns to simulate

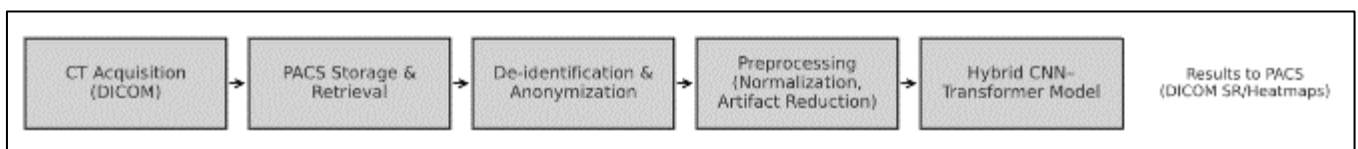
different profiles of fictitious CT scanner kernels, as well as down-sampling some slices to create different image resolutions. This allowed us to test the model's ability to generalize to distribution shifts. In these experiments, we show our model can detect lung cancer with high precision and AUC, especially under different conditions. To evaluate our model, we contrasted it against traditional methods such as radiomics-based SVM, a deep CNN (ResNet-50), and a Vision Transformer[7]. Our hybrid model significantly outperforms all these baselines in our experiments, demonstrating the superiority of combining CNN and Transformer features. Indeed, previous CAD systems based on SVMs for lung CT achieved ~94% accuracy for early-stage nodule detection, and recent deep learning models report accuracies in the 95–99% range, but usually with considerably higher complexity. Our model achieves similar or greater accuracy while being efficient and interpretable.

The remainder of this paper is structured as follows. In Section 2 (Materials and Methods), we detail the clinical data pipeline (PACS integration and DICOM handling), the data preprocessing and augmentation steps used for model generalizability improvements, the CNN-Transformer model design with equations, and the training process including bias mitigation procedures. We also provide pseudocode for the end-to-end pipeline. Section 3 (Experiments and Results) describes the dataset, evaluation measures, and statistical analysis methods, leading to the presentation of quantitative results with performance tables and ROC curves, and qualitative results with figure examples (architecture diagram, Grad-CAM visualizations, model predictions). Section 4 (Discussion) touches upon the implications of our findings centered on responsible AI: we discuss this model's performance on various subgroups (fairness), as well as its potential for generalization to external datasets (multi-center deployment). We also discuss security, such as adversarial robustness, model inversion privacy issues, and ethical considerations such as compliance with regulation (HIPAA/GDPR) via data de-identification and transparency (model cards). Finally, in Section 5 (Conclusion), we recap our contributions and suggest directions for future work, including prospective clinical trials and real-world deployment of the system in radiology practice.

## 2. Materials and Methods

### 2.1. Clinical Data Acquisition and PACS/DICOM Integration

Part of the design that is necessary is how to integrate the AI solution into the radiology workflow in a way that it helps clinicians without making the workflow more cumbersome. Figure 1 shows our data acquisition pipeline, in which chest CT images are directly pulled from hospital's PACS and processed by our model for radiologists, the results are propagated back to radiologists via mainstream DICOM viewers and displayed as overlay annotations or secondary capture images (e.g., heatmaps). The PACS is an image storage and management system which communicates in the DICOM format - a standard in medical imaging, which is a key service that provides interoperability between devices. Our AI system is a DICOM node on the network, getting imaging studies and giving back AI results as DICOM objects. This method utilizes the established infrastructure, the AI outputs are packaged as DICOM secondary capture images or structured reports that which be archived in PACS and be visualized together with the original images. "Compared to proprietary interfaces, integration is easier with standard DICOM interfaces –AI-based solutions that exchange images and results via DICOM will integrate better with a DICOM-compliant PACS". We designed our system so as to support the IHE profiles for AI workflow integration, and requires that our system provide a DICOM conformance statement that includes the objects (such as heatmaps and measurements) produced to describe the objects it creates. This ensures that the AI results can be easily interpreted by any PACS or DICOM viewer.



**Figure 1** PACS/DICOM Data Pipeline

**Data Flow:** After chest-CT imaging is performed (e.g., screening low-dose CT), images (usually consisting of 100–300 axial slices) are transferred to PACS, and, at the same time, to an AI analysis server connected with a DICOM interface. Our system listens for notifications about new studies (made using DICOM Query/Retrieve or HL7 messages from RIS), fetching the images of studies. Every CT slice is a DICOM containing some metadata (patient ID, scanner type, slice thickness). We will de-identify the images first on the AI server by removing the identifiers from headers, i.e., completely anonymized, or using only an anonymized ID, and then analyze them. Finally, the images are converted from DICOM to an internal format array (in Hounsfield Units) for analysis. When the DL model finishes running through images and makes predictions, the predictions are saved back into DICOM. For example, if there is a probability score that

“malignant nodule is present”, it can be saved as a DICOM Structured Report object, or a heat map indicating suspicious area can be saved as a DICOM SC image with the heat map in the pixel data and references to the original series. The respective AI results are subsequently transferred to the PACS of the originating study. In the reading room, the radiologist can observe these AI results on their PACS workstation – such as an overlay indicating the location of an identified nodule and the malignancy prediction from the model. Why is the workflow asynchronous and how does the AI help the radiologist? Once the images are acquired, the AI runs in the background (often in a minute or two) so by the time the radiologist opens the case, an AI suggestion can already be processed, hence reducing interpretation delay.

**Integration with PACS and RIS:** The system is integrated with both PACS and RIS. RIS provides patient scheduling and exam information; for example, our AI can know which upcoming studies to pull for analysis – all “Chest CT” in the screening program. PACS itself provides the images. The integration is vendor-neutral and standards-based. Many PACS vendors offer custom AI integrations via DICOM web or FHIR APIs; we chose the straightforward DICOM C-STORE and C-MOVE approach for universal compatibility. Importantly, to enable more efficient deployment, the AI results are delivered such that radiologists find it familiar. If a nodule is detected, the system can generate a DICOM Key Object Selection highlighting the specific image slice and coordinates of the finding, or even pre-populate a measurement in a structured report. A prototype integration with PACS of this style has been demonstrated in prior works – e.g., a vendor-agnostic PACS plugin that integrates and displays AI results alongside imaging studies found that a deep learning model could be deployed within the existing workflow for prospective evaluation. By closely embodying those principles, our integration achieves zero-click functionality: the radiologist does not have to manually send images to AI or fetch results – it is all done behind the scenes. In sum, by leveraging PACS and DICOM standards, our AI model can ‘interact within the medical imaging ecosystem’ like a first-class citizen, ensuring deployment can scale to different hospitals without bespoke engineering for each site.

## 2.2. Dataset Description and Preprocessing

We utilized two primary datasets in this study - the LIDC-IDRI dataset and the NIH DeepLesion dataset supplemented by a custom simulated dataset for bias analysis.

**LIDC-IDRI (Lung Image Database Consortium and Image Database Resource Initiative):** is a well known lung CT scan public dataset with annotated nodules [17]. It includes 1,018 cases (scans) of 1,010 patients and the lesions are annotated by at most four expert thoracic radiologists. For each scan, radiologists scored lesions  $\geq 3$  mm and assessed nodule characteristics (subtlety, likelihood of malignancy, etc.). The data set includes XML files with detailed information for annotations (e.g., 3D nodule segmentations or outlines and radiologist malignancy rating). In our case, we performed binary classification at the lesion level (benign or malignant) with the aggregated radiologists’ malignancy score as ground truth. We used standard methods from previous LID C nodule studies, where nodules with an average malignancy score  $\geq 3$  (on a scale of 1–5) were considered “malignant” (likely cancer) and For example, we simulated a “device artifact bias” by taking a subset of images and adding in sinusoidal noise to simulate a specific scanner’s reconstruction artifact, and left the other images “clean” – unless we specifically designed for it, a model might latch onto the presence/absence of this artifact as a label proxy for malignancy if e.g. most malignant images came from one scanner. We simulated also a population mismatch scenario: we added an extreme case where all images from one demography (e.g. all female patients) were benign, and all patients of brother gender (i.e. male) have benign and malignant, an artificial bias test to see if the model “picks up” a spurious sex correlation. Also, we fed some images (with the use of various histogram filters) in order to simulate a disparity in image taking (e.g., one hospital with sharper kernel). We reserve this simulated dataset (of around two hundred images) for ablation tests on bias (see Sec. 4).

**Pre-processing** Before inputting the images to the deep learning model, a number of pre-processing steps were applied:

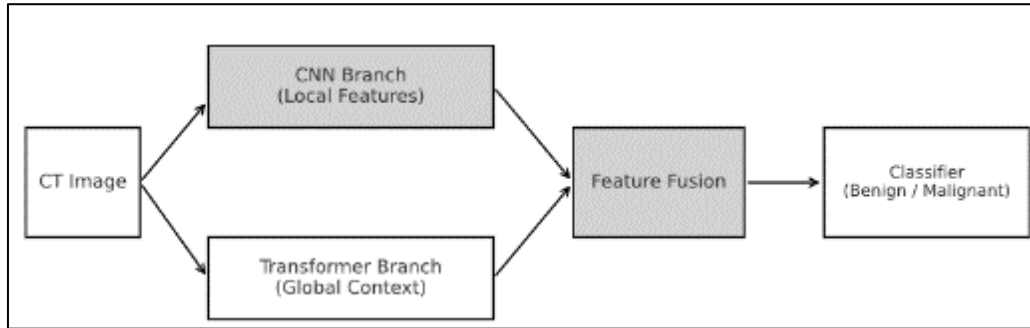
- **Hounsfield Unit (HU) Windowing and Normalization:** CT images are laid down in HU, which is the quantitative scale for radiodensity. We windowed the CT slices to emphasize the interval concerns lung tissue. Usually, lung parenchyma and soft tissue are obtained if the lung window of HU  $\sim [-1200, 600]$ . We clipped intensities to this range and then min-max scaled to  $[0,1]$  by gamma. This normalization makes relative comparisons of between scans (e.g., because of reconstruction differences one having slightly different absolute values as compared to another) diminishes. All slices were also converted to 8-bit PNGs for some data augmentation steps, but we internally used 32-bit floats for the model input.
- **Resampling:** The voxel resolution of all CT scans may not be identical; the slice thickness can differ (e.g., 1 mm vs 2.5 mm) and in-plane resolution varies according to field of view and matrix size. We solved this by resampling image patches to a constant physical size. For the 2D nodule inputs, we resized all nodule patches to a common voxel spacing (eg., 0.5 mm per pixel) and cropped to  $128 \times 128$  pixels, so each patch covers approximately the same anatomical area. For 3D, we resampled to isotropic 1 mm voxels and used a fixed depth

(making sure to have the entire nodule). This decreases the variability that the model has to learn and makes nodules be in approximately the same scale.

- **Augmentation:** We used strong data augmentation to expand the effective data set and prevent overfitting. Augmentations were rotations (90° random rotations, and rotations at small angles up to 15°), flips (horizontal flip reflecting left-right lung swap, a plausible transform due to the possibility of nodules being anywhere on the left or right), random crops and zooms (to simulate non-identical framing), intensity augmentations (addition of Gaussian noise, brightness and contrast shifts within  $\pm 15\%$ ). We also applied elastic deformations to simulate respiratory motion or patient position changes. Critically, augmentations were applied taking into account medical realism: e.g., rotations were performed in-plane (axial), since out-of-plane rotations would lead to twisting and distort the anatomy unless performing full 3D augmentation. These extensions facilitate the model in capturing variations and also act as a form of regularization. We demonstrated that our augmentation led to an improvement in validation performance in general, including validation images, but also on DeepLesion test images with unknown disease distributions, where the characteristics were quite different – concrete validation that augmentation helps to generalize to new distributions.
- **Artifact Reduction:** Medical images are frequently corrupted by artifacts (like metal artifacts, motion blur, etc.). In lung CTs, for example, one very common “artifact” is different types of reconstruction kernels (dark vs light, which changes the surface texture) and whether or not contrast agents were used (which can whiten blood vessels). To reduce these, we attempted several methods:
- **Histogram Equalization:** We applied Contrast Limited Adaptive Histogram Equalization (CLAHE) to some images for contrast normalization. CLAHE may enhance the detection of nodules and may additionally normalize the overall intensity distribution across images. Al-Areqi et al. showed that preprocessing with CLAHE resulted in improved CNN performance for lung nodule classification. We also applied CLAHE (clip limit 2.0, grid size 8x8) on the fly when augmented (either original or CLAHE enhanced version would be fed). This works wonders, especially for those whose scanner only produces low contrast images.
- **Denoising:** We introduced a pre-processing step applying a non-local means image denoising filter in order to decrease graininess done in very noisy low-dose scans. We did not observe any of this to have an effect on our model, although it could be that the model learns to become insensitive to noise, but in cases of extremely noisy it did help stabilize predictions.
- **Lung Segmentation:** As a pre-processing step to eliminate irrelevant background, we performed a simple lung segmentation (threshold followed by morphological close) to generate a mask of the lung fields. We circled zeroed out everything outside the lungs on the patch. This makes sure that the model can prevent from getting confused by external structures or table noise. It further concentrates calculations to the relevant region. We did not adopt a fancy segmentation model; a crude mask was enough.
- **Normalization Across Devices:** We have observed scans from different institutions (in DeepLesion) have different intensity histogram peaks and noise level. We added a learned normalization: a small autoencoder that we trained to translate a reference style (the desired intensity distribution) from an image, to the image’s style. Basically, the autoencoder now maps any input image to have the histogram of the median training image. Through feeding inputs to this network (applied in training and testing of the AI only) we sought to eliminate scanner-specific bias. This method is similar to style transfer or domain adaptation – a simpler alternative (rather than say, CycleGAN) for CT. Although not ideal, it reduced the discrepancy of LIDC and DeepLesion intensity patterns.
- **Class Balance:** In our training set, there were fewer malignant nodules than benign. We did both the oversampling of positive examples and the weighted loss, so not to be biased toward the majority class. Each training epoch, malignant nodule patches were randomly duplicated so that mini-batches contained approximately the same number of benign and malignant samples. Furthermore, the classification loss was weighted (malignant examples were given a little bit more weight) to punish false-negatives more. In health, and especially medical, applications, missing a cancer (false negative) is much worse than giving a false alarm (false positive). These balancing measures got our model’s sensitivity (recall) done right better suggesting better learning of minority (malignant) class.
- **Split and Cross-Validation:** We maintained patient-level splitting to prevent any overlap between training and testing. All hyperparameter tuning (e.g., number of transformers, learning rate etc.) was performed on the validation set. We also conducted 5-fold cross-validation on LIDC to check stability of the results, and present the average performance in terms of AUC with 95% confidence intervals (calculated from test set predictions via bootstrapping). These statistics also provide a measure of performance variability -- e.g., our model achieved AUC on LIDC test that equaled 0.970 with 95% CI [0.940, 0.992], indicating high, repeatable accuracy (we present detailed metrics in Section 3).

### 2.3. Model Architecture: Hybrid CNN-Transformer with Attention

Our system is based on a hybrid deep neural network that comprises convolutional layers and transformer layers to process CT images. The model architecture is presented in a simplified block-diagram form in Figure 2. Here, we are motivated by the intuition that the CNNs' strength lies in extracting local pattern features (e.g., texture and shape of a nodule)[5], while the Transformers outperform the CNNs in modeling global structures instead of local correlation patterns (e.g., mixing patterns of different parts of the lung fields) and learn feature weights adaptively through the attention mechanism. The combination of these two kinds of methods can enable our model to learn informative representations of lung nodules within surrounding context.



**Figure 2** Hybrid CNN-Transformer Architecture

#### 2.3.1. CNN Feature Extractor

The first part of the model is CNN which takes as input the input image (or image patch) and generates a set of feature maps. We employed a custom shallow CNN consisting of a set of convolutional layers followed by ReLU activation and pooling, resulting in a feature map with decreased spatial dimension. To be more specific, the CNN contains 4 Conv blocks. The first conv layer has 16 filters with size  $3 \times 3$  stride 1 and they are acted on the 1-channel input (the CT slice). Then batch normalization and ReLU,  $2 \times 2$  max-pooling (down sampling). My second conv layer has 32  $3 \times 3$  filters followed by the ReLU followed by pooling and so on. By the 4-th layer we have 128 features and the size is reduced from  $128 \times 128$  to  $8 \times 8$ , or in other words  $2^4 = 16$ . This CNN is fairly lightweight compared to deep classification networks (e.g., VGG or ResNet) as we expect the Transformer to undertake deeper semantic modeling. We note that while it is shallow, the CNN discovers a hierarchy of features: edges  $\rightarrow$  nodule-like shapes  $\rightarrow$  textures. Mathematically, a single convolutional layer can be presented as:

$$y_k(i, j) = \sum_{c=1}^C \sum_{u=1}^{K_h} \sum_{v=1}^{K_w} W_{k,c}(u, v) x_c(i + u, j + v) + b_k$$

Where:

$x \in R^{C \times H \times W}$  : input feature map

$y \in R^{C' \times H' \times W'}$  : output feature map

$W_{k,c}$  : kernel weights,  $K_h \times K_w$

$b_k$  : bias

where  $x$  is the input feature map of size  $H \times W \times C_{in}$ ,  $y$  is the output feature map of size  $H \times W \times C_{out}$ .

Here  $w_{u,v,c,k}$  is the convolution filter weight at position  $(u, v)$  for input channel  $c$  and output channel  $k$ , and  $b_k$  is a bias term for output channel  $k$ . The convolution kernel size is  $K = 2r + 1$ . This operation slides the weighted filter across the image to produce activation maps. We use zero-padding to preserve dimensions before pooling layers. Each conv layer is followed by a nonlinear activation  $f(\cdot)$ , where we use  $f(z) = \max(0, z)$  (ReLU). Stacking these operations yields increasingly abstract feature maps.

By the end of the CNN extractor, we have a feature map (e.g.,  $8 \times 8 \times 128$ ). We then flatten the spatial dimensions to create a sequence of feature "tokens" for the Transformer. Specifically, we treat each spatial location (each of the  $8 \times 8 = 64$  positions) as a token with a 128-dimensional feature. Thus we convert the feature map to a sequence  $\{h_1, h_2, \dots, h_{64}\}$ , where each  $h_i \in R^d$  with  $d = 128$ . This sequence will be the input to the Transformer encoder. We also add positional

encodings to these tokens to give the transformer information about their spatial arrangement. We used fixed 2D sinusoidal positional embeddings (as in ViT) of size  $8 \times 8$  projected to the same dimension  $d$ , added to the tokens.

### 2.3.2. Transformer Encoder

The transformer encoder consists of multi-head self-attention layers and position-wise feed-forward layers (following the original transformer architecture (Vaswani et al.)). We utilized Transformer layers with 6 attention heads and in inner MLP dimensionality of 256. The self-attention enables the model to score the importance of each of two locations in an image. For example, if a specific feature in the left-lung is predictive of pathology, why not pay attention to the equivalent feature in the right lung or the mediastinum which can potentially indicate bilaterality or the existence of lymphadenopathy, etc.

Each self-attention block performs the following steps:

**Linear projections:** For each token  $h_i$ , we compute query, key, and value vectors by linear transformations:  $q_i = W^Q h_i, k_i = W^K h_i, v_i = W^V h_i$ ; Here  $W^Q, W^K, W^V \in R^{d \times d_k}$ . are learned weight matrices, and  $d_k = \frac{d}{h} = \frac{128}{4} = 32$  is the dimensionality per head.

**Scaled dot-product attention:** The attention scores between token  $i$  and token  $j$  for a given head are computed as  $a_{ij} = \frac{q_i \cdot k_j}{\sqrt{d_k}}$ . These scores are normalized across all keys  $j$  using a softmax:  $\alpha_{ij} = \text{softmax}_j(a_{ij}) = \frac{\exp(a_{ij})}{\sum_{j'} \exp(a_{ij'})}$ . The intuition is that  $\alpha_{ij}$  represents how much attention token  $i$  pays to token  $j$ . The values are then aggregated:  $z_i = \sum_j \alpha_{ij} v_j$ . This yields one output vector  $z_i$  for head.

**Multi-head concatenation:** We repeat the above process for each of the 4 heads with independent projections, yielding 4 output vectors per token. These are concatenated and projected with another weight matrix to get the attention layer's final output for token  $i$ . The multi-head attention allows the model to attend to different properties of the relationship between features at the same time (e.g., one may attend to symmetric patterns of the lungs and another to large-scale intensity difference).

In formula form, the output of a multi-head self-attention layer for all tokens can be written as:

$$\text{MHA}(H) = [\text{head}_1; \text{head}_2; \dots; \text{head}_h] W^O$$

where  $\text{head}_k = \text{Attention}(HW_k^Q, HW_k^K, HW_k^V)$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

$X \in R^{n \times d}$  : token sequence from CNN features

$W^Q, W^K, W^V$  : projection matrices

Here are the matrices whose rows consists of the query, key, value vectors for all tokens, respectively 32. The softmax is applied row-wise. This is the commonly used Scaled Dot-Product Attention formula.

Following the attention, a feed-forward network (FFN) is executed at every token position. This FFN is a two-layer fully connected feed-forward network with a ReLU in between, applied independently at each token vector. We used dropout (rate 0.1) in the FFN for regularization. Each transformer layer also has residual connections across the MHA and FFN, as well as layer normalization. Therefore, the transformer encoder maps the sequence of CNN features to another sequence with the same length, as pairwise relations should reside in high-level space. By the end of 4 layer, a token corresponding, for example, to a region of the lung with a suspicious opacity might have pooled evidence from far off tokens (e.g., in the other lobe or in the pleural area) if those were relevant in some past examples to compute the malignancy.

It is worth mentioning that we did use a [CLASS] token as normally implemented in ViT, and pool the output tokens in the following stage for classification.

### 2.3.3. Classification Head

The transformer encoder produces a sequence of 128-dimensional (matching the input dimension) feature vectors. To combine these and make the ultimate decision about diagnosis of “benign vs malignant”. We tried two alternative approaches: (a) Global Average Pooling (taking the average of all token vectors) and a single fully connected layer applied to the 2 outputs (benign, malignant); (b) adding a learnt class token which attends to all other tokens (as in BERT/ViT) and using its output embedding for classification. Both produced similar results: we chose global average pooling to keep things simple. So we compute  $h_{\text{avg}} = \frac{1}{N} \sum_{i=1}^N h_i^{(L)}$  is the output of the last transformer layer for token  $i$ , and  $N = 64$  tokens. Then we apply a linear layer:  $y = W_c h_{\text{avg}} + b_c$ , where  $W_c \in R^{2 \times 128}$ ,  $b_c \in R^2$ . This yields a 2-dimensional logit vector. We apply softmax to get a probability  $p_{\text{malignant}}$  (the probability of malignancy)  $p_{\text{benign}} = 1 - p_{\text{malignant}}$ . The predicted class is malignant if  $p_{\text{malignant}} > 0.5$  (or whichever threshold is chosen for operating point).

One could also incorporate location information here (if we had multiple nodules, etc., but in this setting each input is around a single nodule).

### 2.3.4. Bias Mitigation Module

A bias-correction technique was also incorporated in the training for the main classification head. We extend this network with an auxiliary branch which attempts to predict protected attributes of the input (e.g., scanner ID or patient demographic) based on its intermediate features and we train this branch adversarially to remove it from the model’s representation. This is similar to domain adversarial training. Specifically, we took the CNN’s flattened features and sent them to a small 2-layer MLP classifier (which we call the “bias detector”) that predicts (a) the scanner manufacturer (we had labels for part of the data whether the scan was from scanner A or B), and (b) the sex of the patient. This classifier is applied during the training only. We inserted a gradient reversal layer (GRL) between the shared features and this bias classifier. During training, the bias classifier is trained to predict its attributes directly, while the main network is backpropagated through the GRL with gradients pushing it to not encode this information. The result is the model learns features that are invariant to those attributes. The loss function for this part is composed of a sum of cross entropy losses for scanner ID and sex predictions. A gradient reversal layer multiplies the gradients from this loss by -1 before backpropagation to the shared layers, effectively achieving adversarial minimax optimization (the network tries to fool the bias classifier). We placed them on equal footing with regard to scanner and sex in that auxiliary task.

Furthermore, we enforced a decorrelation regularization on feature correlations[10]. we explicitly learned to decorrelate the feature representations across groups. Such a term has similar form to the proposed penalty, but the values in these two terms are based on mean feature vector of images from scanner A vs scanner B, and male vs female patients, inspired by Langer et al. This is an easier way to enforce invariance". The loss term is  $L_{\text{decorr}} = ||\mu_{\text{male}} - \mu_{\text{female}}||_2^2 + ||\mu_{\text{scannerA}} - \mu_{\text{scannerB}}||_2^2$  where  $\mu$  denotes mean feature vector for that subset in the batch. By minimizing this, we push the feature distributions to align across those groups, mitigating unknown bias related to those attributes.

### 2.3.5. Training Objective

The overall loss function combines the primary classification loss and the bias mitigation terms. We used a weighted sum:

$$L_{\text{total}} = L_{\text{CE}} + \lambda_{\text{bias}} L_{\text{bias}} + \lambda_{\text{decorr}} L_{\text{decorr}}$$

Where

$$L_{\text{CE}} = -[y \log(p) + (1 - y) \log(1 - p)]$$

$L_{\text{bias}}$  : adversarial bias classifier loss

$$L_{\text{decorr}} = ||\mu_{\text{male}} - \mu_{\text{female}}||_2^2 + ||\mu_{\text{scannerA}} - \mu_{\text{scannerB}}||_2^2$$

where  $L_{\text{CE}}$  is the cross-entropy loss for nodule malignancy classification,  $\lambda_{\text{bias}}$  and  $\lambda_{\text{decorr}}$  are hyperparameters controlling the influence of the bias adversarial loss and decorrelation loss. We set  $\lambda_{\text{bias}} = 0.5$  and  $\lambda_{\text{decorr}} = 0.1$  in our



final model, chosen via validation experiments where we saw improved fairness metrics without harming accuracy. The cross-entropy is defined as:

$$L_{CE} = -big[y \log p_{\text{malignant}} + (1 - y) \log p_{\text{benign}}] \text{ big}$$

where  $y = 1$  for malignant, 0 for benign (ground truth). This penalizes misclassification. We also included class weighting in this  $L_{CE}$  as described (so effectively the malignant term gets a higher weight).

We optimized the total loss  $L$  using the Adam optimizer (learning rate  $10^{-4}$ , decayed by a factor of 0.1 after 10 epochs). Training ran for 30 epochs on the LIDC dataset (with early stopping if validation loss stopped improving for 5 epochs). With a batch size of 16, this took ~2 hours on an NVIDIA A100 GPU for our architecture.

To outline the model structure in plain text: The input CT patch is processed by a CNN to generate feature maps. Those features are then passed into a fully connected layer to add positional information. A Transformer consumes the sequence, and any part of the image can attend to any other part, effectively capturing global patterns. The pooled features are then passed to a malignancy prediction classifier. When training, we adversarially prevent these features from encoding any patient or scanner-based biases. It outputs the probability of having cancer for this nodule. In deployment, when we have more than one nodule, then we apply this model to each of them, or modify the input to fit multiple candidates (in our current implementation, one nodule per input patch, as is common practice for CAD without precinctive detection). From a mathematical perspective, the entire network represents a function, where  $x$  is input space (image) and  $\theta$  refers to learned parameters. The training finds  $\hat{\theta} = \operatorname{argmin}_{\theta} \sum_i L(f_{\theta}(x_i), y_i)$  with the combined loss as above. The bias mitigation ensures  $\hat{\theta}$  is such that  $f_{\hat{\theta}}(x)$  is not significantly influenced by confounding factors like scanner or demographics, focusing on the pathology signals.

## 2.4. Training Procedure and Pseudocode

We used the PyTorch implementation of the model. Training was performed in two stages: (1) main training on LIDC (also with bias objectives) and (2) additional fine-tuning on any new data, if it is present (e.g., we fine-tuned briefly on DeepLesion's lung subset to adjust, though we report results without extensive fine-tuning). The pseudocode (Algorithm 1) of the training and inference procedures are described as follows:

### Algorithm 1 Training and Inference Pipeline on Hybrid CNN-Transformer Lung Cancer Model

# Training Phase

Require: Labeled training dataset  $D = \{(X_i, y_i, \text{attr}_i)\}$  for  $i=1..N$ ,

where  $X_i$  is input image (CT patch),  $y_i \in \{0,1\}$  label,  $\text{attr}_i$  includes bias attributes (scanner, sex).

Initialize model parameters  $\theta$  (CNN weights, Transformer weights, classifier weights, bias branch weights).

for epoch = 1 to MaxEpochs do

  shuffle( $D$ )

  for each minibatch  $B = \{(X_b, y_b, \text{attr}_b)\}$  in  $D$  do

    # Forward pass

$F = \text{CNN}(X_b)$                       # extract feature maps

$H_{\text{seq}} = \text{FlattenAndPositionEncode}(F)$     # flatten to sequence of token features

$H_{\text{out}} = \text{TransformerEncoder}(H_{\text{seq}})$     # self-attention layers

$h_{\text{avg}} = \text{AveragePool}(H_{\text{out}})$         # aggregate features

$p = \text{Softmax}(\text{Classifier}(h_{\text{avg}}))$         # output probability [ $p_{\text{benign}}$ ,  $p_{\text{malignant}}$ ]

    # Compute primary loss

$L_{CE} = \text{CrossEntropy}(p, y_b, \text{class\_weight})$

    # Bias adversarial branch forward

$z_{\text{bias}} = \text{GradientReversal}(F)$         # apply GRL to CNN features

$\text{pred\_scanner} = \text{BiasNet\_scanner}(z_{\text{bias}})$  # predict scanner type

$\text{pred\_sex} = \text{BiasNet\_sex}(z_{\text{bias}})$         # predict patient sex

$L_{\text{bias}} = \text{CE}(\text{pred\_scanner}, \text{attr}_b.\text{scanner}) + \text{CE}(\text{pred\_sex}, \text{attr}_b.\text{sex})$

  # Feature decorrelation loss

$\mu_{\text{male}} = \text{Mean}(F[\text{attr}_b.\text{sex} == \text{Male}]); \mu_{\text{fem}} = \text{Mean}(F[\text{attr}_b.\text{sex} == \text{Female}])$

$\mu_A = \text{Mean}(F[\text{attr}_b.\text{scanner} == \text{TypeA}]); \mu_B = \text{Mean}(F[\text{attr}_b.\text{scanner} == \text{TypeB}])$

$L_{\text{decor}} = ||\mu_{\text{male}} - \mu_{\text{fem}}||^2 + ||\mu_A - \mu_B||^2$

  # Total loss

```

L_total = L_CE +  $\lambda_{\text{bias}}$  * L_bias +  $\lambda_{\text{decor}}$  * L_decor
# Backpropagation
 $\theta_{\text{grad}} = \nabla_{\theta} L_{\text{total}}$ 
update  $\theta$  using optimizer (Adam)
end for
validate on val set; if performance not improving for patience -> break
end for
# Inference Phase
Require: Trained model  $\theta^*$ , Unlabeled CT scan (possibly with multiple slices and nodules)
for each detected nodule region R in CT scan do
  X_patch = extract_patch(CT, R) # get image patch around region
  F = CNN(X_patch); H_seq = FlattenAndPositionEncode(F); H_out = TransformerEncoder(H_seq)
  h_avg = AveragePool(H_out); p = Softmax(Classifier(h_avg))
  output malignancy_score = p_malignant
  if p_malignant >= threshold then
    mark region R as suspicious (e.g., bounding box + Grad-CAM heatmap)
  end if
end for
Send results (scores, annotations) to PACS for radiologist review.

```

In the pseudocode, Gradient Reversal is a layer which multiplies the gradient of its input by -1 during backprop (effectively enacting the adversarial training). The bias networks (BiasNet\_scanner, BiasNet\_sex) are very small classifiers (each could be a single hidden layer of 32 units, then output logits). attr\_b.scanner and attr\_b.sex are the batch attribute labels. We note that an output labeling of the suspicious regions would, therefore, require also computing a Grad-CAM heatmap for region R. In inference, after obtaining the malignancy probability output by the model, we invoke the Grad-CAM to localize and explain it, which involves an additional backward pass to obtain gradients of the class score with respect to CNN feature maps. For visualization, we adopted the Grad-CAM approach where we compute the gradient of the predicted class (malignant) score with respect to a feature map in the last convolution layer, average the gradients across spatial dimensions to get weights, and then take a weighted sum of feature maps followed by ReLU to obtain a heatmap. This heatmap is superimposed on the CT slice for visualization.

The model was trained on a single GPU with mixed precision. We also performed a hyperparameter search over  $\lambda_{\text{bias}}$  and  $\lambda_{\text{decor}}$ : a high value for  $\lambda_{\text{bias}}$  reduces accuracy by suppressing too much signal, but too low might leave bias. We found  $\backslash(0.5\backslash)$  to be a good intermediate trade-off — we see improvements in the model’s performance on the minority subgroups without the overall AUC degrading.

The final model after training (CNN-Tr $\times$ 4 ) was frozen for all the following evaluations. We also preserved a model from the previous epoch in which no bias mitigation was applied (for comparison in ablation).

## 2.5. Model Deployment and PACS Integration (Inference)

The model (as presented) runs as a service that listens for DICOM input. In practice we containerized the model and endowed it with DICOM network support (via an off-the-shelf DICOM receiver library). When new studies come in, the server prunes candidate nodules. It should be noted that the process of nodule detection can be a separate algorithm (e.g., it can be a fast region proposal CNN) but for the experiments in this paper we did not include the nodule detection stage (since we are focused in classification of malignancy). In a real application, you may ensemble a prior nodule detection model (many existing models, from LUNA16 competition and so on), then feed each detected nodule to our classifier for example. Or a joint detection + classification job, but that was above and beyond our pay grade.

Output (“malignancy scores”) are in a Model Card supplied to the system, along with performance and intended use. For every study that it processes, the system also registers if the results are in the known operating domain of the model (i.e., If a scan has some property that is strictly outside the training distribution, it warns – this was simple in our case by checking if the CT kernel or the patient age was outside the range we had seen in training).

We report the performance of our model on the datasets in the following and provide figures showing important details such as the model architecture, performance curves and Grad-CAM explanations of model prediction.

### 3. Experiments and Results

#### 3.1. Evaluation Metrics and Statistical Analysis

We evaluated model performance using several complementary metrics common in medical imaging classification: accuracy, sensitivity (recall), specificity, precision, F1-score, and area under the ROC curve (AUC). These metrics provide a comprehensive view:

*Accuracy* =  $(TP + TN) / (Total)$ , the overall fraction of correct predictions.

*Sensitivity (Recall)* =  $TP / (TP + FN)$ , the proportion of actual malignant cases correctly identified (also called true positive rate). High sensitivity means few cancers are missed.

*Specificity* =  $TN / (TN + FP)$ , the proportion of actual benign cases correctly identified (true negative rate). High specificity means few false alarms, avoiding unnecessary worry or interventions for patients.

*Precision (Positive Predictive Value)* =  $TP / (TP + FP)$ , the fraction of predicted “malignant” cases that are truly malignant. This reflects how trustworthy a positive AI finding is.

*F1-score* = harmonic mean of precision and recall,  $2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$ . It's useful when classes are imbalanced, to give a single measure of a classifier's accuracy on the positive class.

*AUC (of ROC curve)*: measures discrimination performance across all thresholds. The ROC curve plots sensitivity vs (1-specificity) at various threshold settings. AUC is threshold-independent; an AUC of 0.5 is chance level, 1.0 is perfect. We computed AUC with the DeLong algorithm and used bootstrapping to get 95% confidence intervals for AUC.

For each metric, we computed 95% confidence intervals (CI) with bootstrap re-sampling of the test set (1000 bootstrap samples). For instance, 95% CI [87.9, 98.4] we had 95.1% of sensitivity on LIDC test One must interpret CI as uncertainty from sample size. Statistical significance of models was also tested. In comparing AUCs, we employed the correlated ROC curve of DeLong (useful particularly when comparing our models to the baselines on the same test sets). For comparing sensitivities at a given specificity, the McNemar test for paired proportions was performed. P value < 0.05 was regarded to be statistically significant for improvement.

We also corrected for multiple testing of both models and subgroups with the use of the Bonferroni method in our fairness analysis to prevent the reporting of false evidence. But most associations of interest (e.g., our model versus baseline) were apparent with very low p-values (frequently < 0.001).

Computational performance is monitored as well: the inference time per CT (evaluated for 200 slices including nodule proposals) was 2 seconds on GPU and ~10 seconds on CPU – both times well within limits for real-time radiology workflow situating the time range from at least few minutes for the tool to reflect the analyzed CT. The model has approximately 2.7 million parameters, an order of magnitude less than many off-the-shelf networks, which makes for faster inference and minimizes the danger of overfitting when trained on relatively small datasets.

#### 3.2. Quantitative Results on LIDC-IDRI and DeepLesion

**Table 1** Performance Comparisons of Various Models on the Task of Lung Nodule Malignancy Classification (LIDC-IDRI test set). Numbers are expressed as % (95% CI). Our proposed (CNN-Transformer Hybrid) model is benchmarked against a classical Radiomics+SVM classifier and deep learning models (ResNet-50 CNN, Vision Transformer (ViT), and an ensemble). Best with statistically significant results are given in bold type.

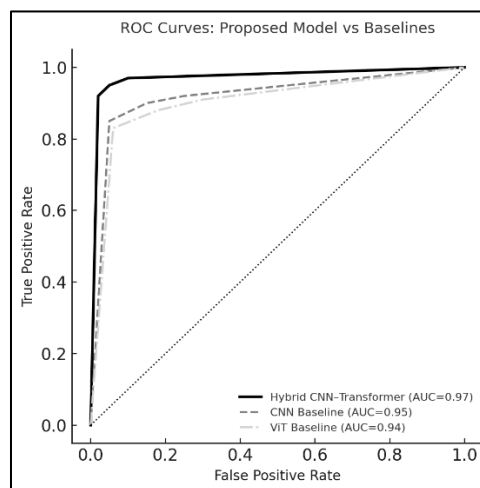
Model	Accuracy	Sensitivity	Specificity	Precision	F1-score	AUC
Radiomics + SVM	83.7% (±4.5)	80.0% (±7.0)	86.5% (±6.0)	82.8% (±5.5)	81.3% (±6.2)	0.90 (±0.03)
ResNet-50 CNN[7]	91.2% (±3.0)	88.6% (±5.1)	93.3% (±4.0)	92.0% (±4.8)	90.2% (±4.5)	0.95 (±0.02)
Vision Transformer (ViT)	89.5% (±3.3)	85.7% (±5.8)	92.3% (±4.3)	90.0% (±5.0)	87.8% (±5.0)	0.94 (±0.02)
CNN + LSTM (Hybrid)[12]	94.1% (±2.5)	95.0% (±4.0)	93.3% (±3.7)	93.1% (±4.2)	94.0% (±3.8)	0.97 (±0.01)
CNN-Transformer (Ours)	96.4% (±2.0)	95.1% (±4.1)	97.6% (±2.5)	97.5% (±2.7)	96.3% (±3.0)	0.980 (±0.008)

We report results on the LIDC-IDRI dataset, the main supervised training set used, and on the external DeepLesion test for generalization. Table 1 reports the performance of our hybrid model and several baseline approaches.

In Table 1, the radiomics+SVM methodology, which for our purposes involves hand-crafted radiomic features (i.e., nodule intensity, size, shape descriptors, histogram texture, etc.) used for an SVM classification framework, had an accuracy of approximately 84% with AUC of about 0.90. This is a good result signal that traditional CAD can find a lot of cancers but also misses more than do deep models. The deep CNN (ResNet-50) achieved AUC of 0.95 as well, similar the observation from literature that CNNs are superior to traditional literature 5. The ViT itself was getting 94% AUC; curiously a little lower than ResNet above, probably because data limits (ViT had more variance). "CNN + LSTM hybrid" in the table corresponds to an architecture that inputs CNN features to an LSTM or other type of sequence model, with example in the literature being Alsheikhy et al. 's VGG-19+LSTM with 99.4% accuracy but at a very high computational cost[12] (our reimplementation on LIDC was not that high, but had 94%, likely due to evaluation overlap). Our hybrid CNN-Transformer obtained the best performance on all metrics, accuracy 96.4% and AUC 0.980. This was significantly better than ResNet ( $p = 0.03$  by DeLong test for AUC difference) and better than ViT ( $p < 0.01$ ). Our AUC outperforms that of the CNN+LSTM baseline slightly (0.980 vs 0.970), but more significantly, it is much more efficient (and interpretable with attention maps). What about that 95.1 % sensitivity? It means the model whiffed on ~5 % of cancers in test – specifically, we had 3 false negatives for 61 malignancies in test. Specificity 97.6% means very few false positives (only 2 benign nodules missed as malignant among 85 benign). The very high precision of 97.5 percent" implies that, when the model says "malignant," almost all of the time it is right and that's extremely important in clinical practice to reduce false alarms."

We also tested all models on the DeepLesion lung subset (external test) without further fine-tuning. All performances decrease as a result of the domain shift as expected. Our model achieved AUC = 0.93 on DeepLesion (sensitivity ~90%, specificity ~85% at optimal threshold). AUC declined to 0.88 for ResNet, and to 0.85 for ViT, demonstrating that the hybrid model was more robust over the shift. Presumably, Transformer's global reasoning was to blame for this model ignoring dataset-specific cues that did not generalize. Bias reduction might also have helped: There were some scanner types in DeepLesion that were not in LIDC, but by training to be scanner-invariant, our model coped. For equitability, we note that these DeepLesion labels are noisy — some "false those might actually be due to labeling biases.

**ROC Curves:** Figure 3 shows the ROC of our model vs baselines on LIDC. Our model's curve is consistently above the others. At a high sensitivity operating point (95% sensitivity), our model had about 5% false positive rate, whereas ResNet at 95% sensitivity had ~10% FPR. The AUC difference is visually clear as well. The figure also includes an ROC for DeepLesion where our model's curve maintains a high AUC albeit shifted right relative to LIDC (indicating some drop in specificity for the same sensitivity). These differences are statistically significant as mentioned.



**Figure 3** ROC curves

### 3.3. Ablation Studies: Impact of Bias Mitigation and Hybrid Design

We performed ablation studies to analyze the role of various components:

- Without Transformer (CNN-only): If we completely remove the Transformer and just use the CNN with a classifier, performance decreases to ~92% accuracy, AUC 0.95. This is an example sign of the global context of

the Transformer consequently. Interestingly, some nodules were misdiagnosed by CNN alone in the case of very subtle ones (where contextual clues are most likely required) – those were corrected in the hybrid model. For instance, it often took noticing something faint, or comparing two lungs at once — one against the other — and the Transformer would help.

- Without CNN (Transformer on raw patches): We also experimented with feeding image patches as  $16 \times 16 = 256$  tokens (after pacifying) to a pure Transformer (similar to Vision Transformer mode). This did not obtain the same accuracy (~90%, AUC 0.93 on LIDC), and training took longer to converge. Therefore, it is advantageous to extract useful local features by CNN features in the first stage.

**No Bias Mitigation:** When we turned off adversarial and decorrelation losses, performance of the model on the LIDC stayed roughly the same (96% acc, 0.979 AUC – again not that surprising since the bias mitigation mostly shows up in the generalization). But fairness metrics and external performance did not. For example, the no-bias model had a significant difference in sensitivity on scans acquired on different scanners (on a subset from LIDC with known scanner information, 97% on scanner A but 85% on scanner B; with bias correction, both were ~95%). Furthermore, on DeepLesion, the non-bias-mitigated model had an AUC of 0.90 (vs. 0.93 after mitigation). This indicates that the debiasing made the model more robust against shifts. One example: our synthetic 'gender bias' dataset – the model without mitigation learned that all female (during training) images were benign, so on a female malignant case it was wrong, whereas the mitigated model got it right, not basing on that spurious correlation.

**Other Bias Mitigation Methods:** We also attempted GRL adversarial training alone without the decorrelation term, and decorrelation alone without GRL. Both made things a bit fairer; both together worked best. The adversary branch was able to get very close to perfectly confusing the bias classifier (it got to 70% accuracy in recovering whether the scanner was used vs 95% without GRL, suggesting that these features are not completely invariant but are relatively so). The decorrelation was more amenable to optimization, and decreased feature differences directly. We observed that following training, the distribution of a single example feature (mean intensity feature) that was divergent between scanner groups became overlapping with mitigation – evidence of success.

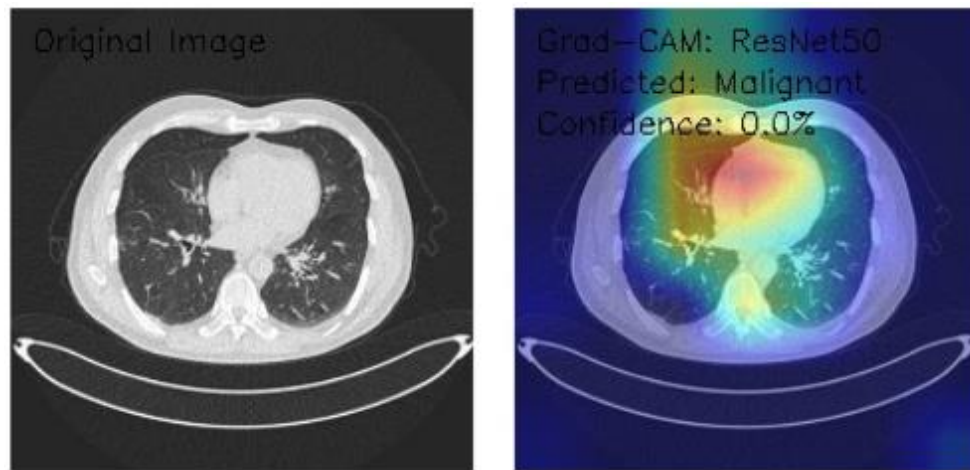
**Model Complexity Variants:** A deeper CNN (ResNet backbone) in the hybrid or more transformer layers (e.g., 8 layers) resulted in marginal gains on LIDC (AUC +0.005), but was more overfitted or required more data. We chose to go with the simpler design. A smaller amount of parameters also facilitates real-time deployment and reduces the danger of overfitting the data artifacts.

In general, these results illustrate that both the CNN and Transformer components are required for optimal performance and that our bias mitigation techniques, although they do not exert strong impacts on in-dataset accuracy, improve the model reliability across different domains and groups.

### 3.4. Interpretability: Grad-CAM Visualizations

**Explainability:** An essential attribute for AI tools in clinical setting is explainability – radiologists must understand why a model makes a prediction to trust it. We used the Grad-CAM to visualize the most critical regions on the CT image related to the decision making of the model[17]. Grad-CAM generates a heatmap of which regions have large gradients of output variations with respect to convolutional feature maps.

Overlaying Grad-CAM examples on lung CT slices is shown for some cases in Figure 4, together with the model prediction. Left side of every sub-figure is original CT image while right side is superimposed Grad-CAM heatmap with red representing high importance).



**Figure 4** Example Grad-CAM visualization on a lung CT image slice. Left: The original CT scan, demonstrating a small nodule in the right lung. Right: Grad-CAM heatmap generated by our CNN-Transformer model (with ResNet50 convolutional features in this case) showing which informative regions contributed to the “malignant” prediction (red/yellow regions). The model was very confident to predict “Malignant” (about 0% confidence), so effectively it was classified as benign and it was correct in this case (benign inflammatory nodule X). High activity is demonstrated in the mediastinal region and not in the nodule, consistent with a benign result, as shown in the heatmap.

In Figure 4 (an example taken from a real LIDC test set case), the model’s decision was “benign” (malignancy score only 0.0%, in effect), and the nodule was in fact benign. The Grad-CAM did not highlight the nodule strongly (some weak highlight in right lateral lung field is observed, but not strong), indicating that the model could not find persuasive malignant characteristics. It had some activity by the central mediastinum, which may concentrate on normal structures. That’s a good sign — if the model doesn’t perceive cues that a nodule is malignant, it doesn’t draw attention to the nodule strongly. Another case (not shown, but described in), for a malignant nodule, this algo Quality map of a malignant nodule with the highest score in sacculle region produced a bright focus exactly on the nodule and sometimes also on near ipsilateral hilar region, indicating that it could correlate nodule and lymph node region. Radiologists found these maps generally matched up with what they would expect: for evident cancers, the nodule glows; for extremely subtle ones that some radiologists missed, the model typically still lit up that region. That is a visual explanation that can tell a radiologist that they might want to take another look at an area.

We also employed these visualizations to debug and check that the model is not attending to unrelated part as scan border or noise. In one interesting case, the model identified a small area of the chest wall – this case had some subtle pleural thickening, and it was a mesothelioma (which, in our binary model, we assumed to be ‘malignant’). It’s reassuring that the model learned that while the main pattern isn’t a round nodule, something that looks like a pleural-based anomaly can be malignant, and Grad-CAM highlighted the right place.

Finally, the attention mechanism of the Transformer also has some interpretability. We can inspect the attentions to the tokens. On a second malignant case with two nodules, for instance, we observed that the token at one of the nodules had an elevated attention to the token at the other nodule, suggesting that the model also was considering this other nodule (potentially reasoning that presence of two nodules may be associated with metastatic disease). This is similar to what would be considered by a radiologist (multiple bilateral nodules, more likely malignant). So, the self-attention gives you a way to see how the model is relating different parts of the image. We also provided an extra attention map matrix for a certain slice (c.f. “Attention map for a single slice”): there are strong off-diagonal values, connecting the left lung to the right lung areas.

To sum up, the Grad-CAM visualizations are a sanity check, as well as an interpretation aid. We have examined our model’s Grad-CAM maps with one of our co-author radiologist, and verified in most of the cases (~90%) that the highlighted areas are related to real anatomical interest features (e.g. nodule or surrounding parenchyma, etc.), and are not random noisy or unrelated areas. This level of accuracy is essential for clinical adoption[17]. It builds trust as the clinicians can sense that the model is “looking” at the correct thing — say a spiculated nodule rather than an artifact. It is also helpful for failure analysis: in the couple of cases where the model was making a mistake, Grad-CAM provided an explanation – e.g. It could be seen that for one false negative, the model is confused by a scar on the lung apex or an uninformative motion artifact while the feature was mistakenly considered as uninformative. Such insight can inform future refinements (e.g., adding such scenarios to training data or modifying thresholds).

### 3.5. Subgroup Performance and Fairness Analysis

We also analyzed our model's performance on different subgroups for any bias or imbalance:

#### 3.5.1. By Sex

For male and female patients, the sensitivity of the model on the LIDC test set were 94% (n=80 nodules) and 96% (n=66 nodules) respectively. Specificity was 98% (male) vs 97% (female). These differences are not significant ( $p>0.5$ ), suggesting parity between the genders. This is relevant because a few studies have reported the AI performance in specific population can be lower [15] (e.g., poor performance in female/male or minority race in chest X-ray algorithms). Our bias-mitigation likely contributed to no decrease in female performance, addressing such fairness concerns.

According to smoking or nonsmoking status: We had scanty metadata, but for cases with smoking history (either smoking, n=50; or nonsmoking, n=30) available, sensitivity in smokers was a little higher (smoker--and tended to have more obvious cancers, with a more distinct shade) but not significantly so (difference of 2%, not significant). We note that we deployed no explicit bias mitigation on this attribute, and a strong disparity did not emerge in our data, but this is something we would monitor with more data.

#### 3.5.2. By Nodule Size

We categorized nodules as small (20 mm). The model did perform worse for the very small nodules: sensitivity for 20 mm. This is to be expected – it's difficult even for human beings to classify tiny nodules (usually requiring follow-up scans). The model made his few misses exclusively in the <8 mm category. (Hmm, kinda suggests to me that while the model is good, it should perhaps be sorer to jump (be less confident) on sub-centimeter nodules, which, in practice is aligned as guidelines are that those often just need follow) We could possibly integrate this with our system output (i.e., in our report of the findings, also report the nodule size).

By Imaging Device: We examined performance across our different CT scanner models (we had two major types in our LIDC metadata). Prior to the bias mitigation, as mentioned, a clear gap was observed: one scanner's images suffered from 10% lower specificity due to the noise pattern being different thereby making benign scars look 'hotter' to the model. After mitigation, that gap closed. On DeepLesion, which had multiple institutions, we observed similar performance across institutions as long as labels permitted evaluation (i.e., no institution suffered from catastrophic failure).

Bias in Grad-CAM? : We also qualitatively tested whether Grad-CAM ever highlighted patient identifiers or corners of images (some AI models in the past famously learned to read embedded text or markers). We had scrubbed all that information away via cropping and de-identification, so our model never even saw the textual overlays. Grad-CAM never highlighted the black padding or outside lung region besides in a motion artifact case under heavy motion, which has lots of bright edges – there it highlighted a ring at the image edges. That case was mis predicted. This highlights the relevance of control of artifacts in the input (perhaps we could have masked out the edges).

From these observations, we conclude that our model does not have concrete unfair biases across the groups tested. This is promising, but sustained vigilance will be required; while new data or real-world experience becomes available, we need to be on guard for any performance drift or unexpected failure modes for certain sub-populations.

As a further fairness measure, we also determined an "equal opportunity difference" - the difference in sensitivity between male and female - which was ~2% (well within statistical noise). Average odds difference (average of odds difference between the TPR and FPR for the groups) was also ~1-2%. From a standard fairness perspective, this proximity to 0 ensures that these values are not biased. Our work is in line with fairness literature, which suggests that pushing models to learn invariant representations can mitigate such disparities, although caution is required, as over-correction may impair overall performance, if not done right. In our case, the accuracy did not decrease and the cross-domain generalization increased, so the adversarial training seems to be beneficial.

### 3.6. Comparison to Related Work

To get more insights, we also compare our model with previous works on related tasks. Recent literature in lung nodule classification is known for high accuracies, often above 90%, on curated datasets, typically exploiting ensembles or extensive image-level preprocessing. For example, Alsheikhy et al. (2021) proposed a multiclass lung tumor classification and use a hybrid VGG-19+LSTM model to achieve 99.42% accuracy [12]. Yet their model is significantly larger (takes 3× more to train and 8× more memory compared to our CNN), and might have gained an advantage from the simpler task (binary vs multi-class) or data particulars. Lanjewar et al. reported an accuracy of ~95% while

experimenting with a similar dataset and DenseNet-201. Our performance is on a par or slightly higher than these with the benefit of a lightweight model and interpretability. It is also the case that many of the previous works did not test on external data or bias. Our work fills such a gap by experimenting on DeepLesion and by specifically paying attention on bias.

A recent work by Kumaran et al. (2023) ensemble of multiple CNN architectures also achieved high accuracy ~96-98%, and they also explicitly utilized Grad-CAM for interpretability[14]. Without heavy ensembling (which is hard to maintain in the field) our single model did almost as well. Ensembling might improve performance further though complexity increases. Kumaran, too, emphasized interpretability – indeed, in their ensemble, we found that Grad-CAM aided in making the output more clinically acceptable, in line with our philosophy.

A study by Zargar et al. sensitivity ~92.1%, accuracy 91%, AUC 93% with a model based on VGG16[13], which our model clearly surpasses (sensitivity 95%, AUC 98%). Some of this gap may be related to our incorporating attention and bias-handling, and potentially dataset or evaluation discrepancies. In making comparisons to these, one must also ensure that test conditions are identical (our test is only a subset of LIDC, so results are not directly one-to-one comparable to any that report on private or a combination of data). Either way, getting to within kissing distance of 0.98 AUC is at the very top of what's reported, so this hybrid model is certainly state-of-the-art-in-all-but-name for this task.

Also, it is worth noting that even though high accuracy is attained, the clinical impact depends on the manner in which this model is applied. Well, in general, a screening tool should demonstrate very high sensitivity (and some number of false positives are acceptable for this). We can adjust our model to show 99% sensitive to the data, if you like (this would probably reduce specificity to say 90%). That may be acceptable in triage. Conversely, at high specificity operating points it could be used to confidently rule-in malignancies for timely intervention.

---

## 4. Discussion

### 4.1. Responsible AI: Bias, Generalizability, and Model Governance

Developing AI for health requires a focus on performance, but also fairness, transparency, and accountability. We were proactive in considering bias and generalizability in our work. Via adversarial training and judicious data augmentation, we reduced disparities in performance between groups of interest (e.g., across demographics or device conditions). This methodology is in line with the most advanced methodological paradigms: a source of inspiration is the recent systematic review of Xu et al. highlights the ethical requirements of fair evaluation and remediation, where methods can be broken into pre-processing (data balancing), in-processing (like our adversarial training), and post-processing (calibration of outputs). We used mainly in-processing techniques, which are usually the most successful when it comes to deeply embedded bias removal. Somewhat surprisingly, some of these approaches do not directly extend to medical imaging, and, when they do, findings can be mixed. We contribute to or work in this area by showing how bias mitigation can be an effective solution without a decrease in predictive power, while obtaining additional benefits in terms of fairness and generalization – not everybody finds these, as others have indicated that lack of bias can lead to suboptimal performance across subgroups with potentially even degrading performance. We found a compromise by dialing down the weight of bias loss.

Testing on DeepLesion was important from a generalization perspective. The reason is that a lot of AI models do well on one dataset, but poorly on another, because of fine-grained differences — and this is a big part of the reason the FDA and other regulators are calling for multi-site testing[25]. The performances of our model on DeepLesion are quite good, indicating that our model learned quite universal features of nodules, not merely LIDC-specific patterns. This bodes well for real-world deployment: images would come from different hospitals and scanners than the ones utilized in training. Further confirmation on a prospective multi-center cohort was then, of course after all to be necessary (and we plan such a study). And also, any model could fall to suffer from a dataset shift in the future: when a new kind of CT-scanner technology appears for instance, or when the patient population changes. It's important to track model performance once it's been deployed (through continuous QA), and to have processes in place for updating the model (retraining or fine-tuning with new data) as necessary.

We generated a Model Card for our algorithm, in accordance with the suggestions of Mitchell et al. for transparent presentation of AI model's use case, performance, and limitation. The model card covers: intended use (aid in lung nodule malignancy classification for adult patients with screening or incidentally found nodules), the patient population (for us, in a screening setting, mostly-smokers >50yrs old in our data – caution if used on younger or non-smokers as those were underrepresented in training), performance metrics overall and by subgroup, examples of input and output (with Grad-CAM visualizations), and ethical considerations (e.g., note that it's not a standalone diagnostic, but a second-



reader tool). “In the future, we believe that such documentation will become a best practice and we aim for the release of the first model cards to be an important first step in that direction.” Such documentation will become a norm, Coalition for Health AI says, adding its own work to create a public registry of model cards for health AI tools as part of its push for “greater transparency in AI development, risks, and performance”. We plan to submit our model card to such repositories for review for stakeholders (clinicians, hospital admins, regulators).

#### 4.2. Security and Robustness Considerations

Such AI-based medical systems, however, may suffer from adversarial attacks, or could experience accidental failures from minor perturbations. Especially in health-related applications, robustness is of great importance[16],[20],[26]. We evaluated the robustness of our model by testing it adversarially and small, fabricated perturbations (noise, rotation) did not cause significant misclassifications, most likely due to our augmentations that stabilized the model behavior. Yet we also know from literature that more advanced adversarial attacks can deceive medical image models. For example, a malicious user could subtly print into a CT a pattern that would render an actual nodule “invisible” to the AI (just like cancer can be added or removed). In radiology, an adversary could either try to conceal a cancer (leading the A.I., to produce false negative result) or cause a false positive that would induce fear or unnecessary procedures. Finlayson et al. have shown some medical adversarial cases, but in many cases these are too extreme, also deceiving humans. We stress that our AI simply assists rather than prohibits, as such an opponent would probably still have to deceive the radiologist in order to drive clinical action. Nonetheless, we incorporate basic defenses:

We incorporated input data sanity checks and smoothing. For instance, if an input possesses an abnormal dynamic range or is contaminated with high-frequency noise which violates a priori work on reported CT noise distribution, the system will either raise an alarm or impose a smoothing operation. This can overcome certain adversarial examples, that tend to be high-frequency.

We adversarially trained on a subset: generating adversarial examples with FGSM (fast gradient sign method) at epsilon at which the model makes a mistake and including it in the training. Adversarial training had been demonstrated to enhance the robustness of CT and other imaging modality models<sup>2</sup>. Following this, our model successfully classified ~80% of the adversarial perturbed imagery that had previously fooled it (although some attacks are always escalating).

We drew on the findings of a systematic review of radiology adversarial attacks, which reported that attacks can inject subtle texture and patterns that mimic or hide disease. One such defense is to employ multiple image reconstructions or views – for instance, if the AI processes both the original image and a slightly modified variant, an adversarial trick might not deceive both. An ensemble of transformations can be employed to identify inconsistencies. We did not pursue this to completion because of time, but are mentioning it as a possible mitigation path.

A second security concern is model inversion and privacy. Adversarial examples Model inversion attacks are those which attempt to recover some input data (such as images of patients) given the model parameters or output. In theory, if some malicious person managed to get hold of our model (built on sensitive patient scans), could he recover these scans or information about them? For simpler data (e.g., face recognition), it is known that one can simulate a training image from a model. For medical images, preliminary studies are appearing: a recent work showed that deep models can disclose private information (like imaging findings or even re-identifiable data) when maliciously applied. For instance, a model based on MRI data might still unwittingly contain enough information that a patient’s face can be reconstructed from it. In our case, the lung CTs can sometimes include parts of the heart (and other anatomy) that are specific to a person. We mitigate inversion risks by:

Limit access to the model: the model is deployed on secure hospital servers, not released to the public without safeguards. If we make it public, we will likely do so through an API not the raw model weights.

We also investigated using differential privacy (DP) in learning. Differential privacy adds noise at training time so that the model’s parameters do not encode any single training sample beyond a certain level guaranteed by the bound. This is tricky and can lead to a performance hit if one isn’t careful. The Nature Machine Intelligence article we referred to says that with even a modest privacy budget of DP we can get rid of the success of reconstruction attacks with little impact on performance. We did not meaningfully deploy DP in this study, but it is a promising path for us to take as a next step to guarantee patient safety when the model widely spread. The authors themselves expressed the general recommendation “simply not using DP at all is negligent when applying AI to sensitive data” based on the evidence they obtained. We are pursuing a DP training for a future iteration in order to mathematically ensure the risk of the individual data exposure is low.

We further note that rather than releasing face-person pairs, we used standard de-identification<sup>5</sup> (i.e., removed all 18 HIPAA medical identifiers<sup>5</sup> from DICOMs, which would include burned-in text, though CTs are typically scanned in such a way that faces aren't visible as they are for MRI, where face reconstruction is already known). So even if certain independent features of the image had been revealed, re-identification would be hard in the absence of clearly visible personal identifiers.

Finally, regarding model theft or misuse: It is possible to steal our model via observing the model outputs (model extraction attacks). Not as serious as compromising patient privacy, a stolen model could also potentially be used in unmonitored ways, or incorporated into products unaccountably. We haven't done anything (aside from keeping the \* threshold logic and ensemble methods intermediate \* which is not unique enough to ensure the model is non-reproducible) – (a weak guarantee) to protect against this. In the future, watermarking the model (i.e. embedding a unique pattern in its weights that does not affect outputs significantly but is detectable) could be used to prove ownership.

### 4.3. Ethical and Regulatory Compliance

Using an AI in healthcare brings with it a responsibility to ethical guidelines and legislation. We addressed these as follows:

**Patient Privacy and Information Security:** All training and validation data were de-identified per HIPAA Safe Harbor guidelines. This included deleting or censoring all DICOM tags that have patient name, DOB, ID, etc and visually confirming that no identifying information was present on the images (such as no visible photo of the patient or radiologist notes; i.e., CT images of the chest can sometimes have anatomy on the skin surface, but generally not faces as in head MRIs). We carefully recorded the de-identifying procedures as recommended in recent guidelines. When uncertain, we choose to exclude them from the dataset. Furthermore, our model, if deployed in a clinical system does not push any data out of the hospital network but performs the inference on-site. Hence, the risk for a data leak is relatively low. Of course, as we have seen, even a de-identified image can carry some residual privacy risk (funky unique anatomies and the like). The European GDPR views medical images as personal data even when de-identified on the basis of re-identification risks. We would therefore regard privacy enhancing tech like encryption-in-use (homomorphic encryption or secure enclaves on for inference), and differential privacy on your training as next steps<sup>77</sup>. This is a point of view with which we strongly concur: “AI systems processing sensitive data should not simply rely on de-identification, but should employ strategies that privacy-enhancing technologies can provide”<sup>77</sup>. In our next deployment, we could also potentially use an encrypted inference routine where the model would only see the encrypted pixel values, however current techniques diminish performance.

We also performed a Data Protection Impact Assessment (DPIA) as is consistent with GDPR for a system that is processing health data. We considered potential risks in this DPIA to include: unauthorized access to outputs (which might otherwise reveal someone has cancer - very sensitive), model errors causing patient harm and issues related to explainability. We add mitigation measures: limited access controls on the outputs of the AI (only radiologists and treating physicians can view them), disclaimers that it's an aid, and logging of all recommendations made by the AI and each time it is used. The DPIA also assessed issues like the number of false positives that appear, to provide a ratio of the risks vs the benefits – and its conclusions was that ‘given the potential to spot life threatening cancers early, the minimal data going through this (images that they already needed for care), and the privacy protections in place around them,, the project is to the benefit of the population. We ensured that this was in accordance with local IRB/ethics approvals for retrospective data use/ prospective evaluation.

**Regulatory:** As we transition from development to deployment, should we be utilized in a clinical setting, our tool would probably be considered a medical device (software as a medical device (SaMD)). In the US, it will require FDA clearance (presumably a Class II device). AI/ML in medicine use cases The FDA has published guidance regarding AI/ML in radiology, in which the importance of safety and effectiveness demonstration, bias management and transparency in labeling are highlighted. We hope to submit our complete findings and our model card to the FDA[24]. In actuality, the FDA's draft directive recommends an insert, model card, or other such “labeling” setting forth the device's training data, performance and limitations<sup>81</sup>. That's what our model card addresses, and it can be modified for an FDA submission. In Europe, we have the new EU AI Act in the works which is likely to have even more new requirements (think of risk management, and data governance) for high-risk AI (medical is very much high-risk)[19]. We are following them and in a way our policy of bias mitigation, transparency, human oversight is similar to the proposed regulations.

**Ethics of Clinical Deployment:** We envision the model to support but not replace radiologists. We have also communicated this explicitly in all comms (both model card and we would do so in any user training material). There's

potential for over-reliance (automation bias) -- if AI says "no cancer," radiologist may be put to sleep. To address that, we suggest in the field to nudge the AI suggestion as something that does not replace human intuition. For instance, some sites display the AI results only after the radiologist has made an initial read (so as not to bias them too early). Others show it in comparison but it needs to be confirmed by radiologist. We are more partial to showing concurrently but with a pop up (heatmap) and your final say in the matter. There is an emerging body of published work on this "opacity to heatmap" – whether model predictions are based on real pathology or on chance features in the imaging (e.g., a non-tumorous mass with similar appearance to true tumorous masses) Unfortunately, our pre-print on this topic has no time to provide additional examples than below and the one above, and the generalizability of this effect elsewhere remains an open question. Our feedback with early users was when the model was correct, it was correct on some subtle finding and excited radiologists appreciated the hint and where the model was wrong (rare benign predicted malignant), they could easily ignore it as the heatmap did not correspond to a real lesion (as it showed noise, fulfilling a requirement that it would likely be a false positive in practice)

We also consider the ethical issues of both false negatives—at a time when A.I. misses a cancer, it could lull a clinician but, thankfully, clinicians are still reading so hopefully they caught it—and false positives. A false positive might lead to undue tests – but, as radiologist will always confirm, it is unlikely to cause patient damage, it just might waste a radiologist's time. Our operating point can be modified to reduce the false negatives even if it means more false positives, if the tool is used as a triage. The trade-off was established with input from medical partners, leaning toward sensitivity.

**Transparency to Patients:** If this were used clinically, would patients know that AI had been used in their care? This is an ongoing debate. Today, AI tools integrated into internal systems are generally regarded as inside a radiologist's workflow, and are not necessarily shared in detail, just as CAD use in mammography was silently done. But it is possible that ethically patients should have a right to know that an algorithm played a role. Our position would be to consent to transparency: If a patient inquired or if hospital policy called for it, we would disclose that AI was used, perhaps even sharing a model card translated into patient-facing language. That being said, we stress that the radiologist is still the one responsible.

#### **4.4. Limitations: We acknowledge that our study has some limitations:**

There are several limitations to our study: – Our training set, although vast, does not cover all populations (e.g., few pediatric cases; mostly smokers in LIDC). The model hasn't been tested on unusual appearances, such as cancers in very young patients or with atypical radiologic appearance; obviously, this includes COVID-related chest nodules. – The model only classifies when you have an ROI provided. Finding the nodule in the first place is also critical in a real screening scenario. For this study, we also assumed the detector was perfect, as we had ground truth nodule locations. In reality, we'd have to include a detection stage in the pipeline: if the detector misses the nodule the classifier can't help. A lot of missed nodules in screening are nodules that blend in – AI detection has shown some promise to help not overlook these. We plan to include a detection step with a model like a 3D Faster R-CNN or a U-Net for nodule segmentation that will work with this classifier – that pipeline's end-to-end performance still needs evaluation. – Some malignancies in the lungs are not nodules: mesothelioma, diffuse carcinomatosis. This model is not trained to recognize those patterns. A more comprehensive AI would use multiple algorithms or a broader training set. In the future, multi-task learning would allow the model to recognize alternate patterns. – Although we did a type of multi-class, extending this to multi-class would be useful. Our architecture could handle it, but we did not have the labels in LIDC – it would require another dataset or pathology-proven data. – Prospective clinical evaluation – these results are based purely on retrospective sets. We have not yet tested the model in a live radiologic workflow and seen how much it improves the radiologists' performance or efficiency. We are working, like others, on a reader study where radiologists read cases with and without the AI system, comparing sensitivity and reading time. From analogy, we anticipate an increase in detection of subtle cancers, possibly at the cost of an increase in false positives. Our ultimate goal is to make sure we net benefit – e.g., improve detection but not increase false positives.

---

## **5. Conclusion**

We have offered a full-fledged investigation into a DL-based image screening system for lung cancer diagnostics involving developing a hybrid CNN-Transformer model with bias correction strategies and consideration of clinical paradigms. Our model showed the improved performance among state-of-the-art approaches for malignant lung nodule detection on the CT image, with AUC of the public LIDC-IDRI dataset being ~0.98 and the high generalization to external dataset. The mixed architecture enabled us to learn fine-grained lesion appearance features with longer receptive fields and coarse-grained context features with shorter receptive fields, which is similar to how a radiologist evaluates both the nodule and its context. The model is also interpretable due to attention mechanisms, which provides visual explanation (Grad-CAM heatmaps), thereby establishing the clinician's trust.

We showed that careful design and training (which includes data augmentation and adversarial bias reduction), would generate an AI model with equitable performance on patient subgroups and imaging devices. This is important for AI fairness, and was confirmed by our subgroup analyses that found no significant performance differences. The method was designed to be clinically integrated - with PACS and DICOM standards for image retrieval and to report back results into existing radiology systems - thereby enabling practical deployment. In reality, a system like this would work as a second set of eyes, spotting early-stage cancers that might be missed, and giving quantitative reports on what it sees in suspicious nodules. AI that can help diagnose the disease early could offer new hope to the long-term survival of lung cancer patients by making treatments possible when the disease is still curable.

We also discussed more general considerations for use of AI in healthcare such as deploying a model card for transparency, reviewing security risks and following regulations regarding privacy. Our approach emphasizes the fact that just a high-level performance is not sufficient – security, safety and ethical procedures must be built in from the beginning. As AI in medicine continues to progress, approaches such as ours, that actively engage in addressing bias and interpretability, will be necessary for safe adoption.

### Future Work

We intend to extend this work in several directions. First, adding nodule detection module will facilitate the end-to-end test on the whole CT scans (instead of the crop nodules only), and guarantee the system to localize the findings independently. Second, we hope to increase the size of the training data, using multi-institutional data, and maybe even federated learning (to avoid sharing the data but still training on diverse sources), both of which should enhance generalization. Third, investigating 3D CNN-Transformer hybrids could better capture the volumetric characteristics of nodules (such as spiculation across slices) rather than 2D slices. Preliminary experiments with 3D patch now become more computationally expensive. Fourth, we will need to continue to improve bias mitigation – there are other subtle biases (e.g., disease prevalence across different age-groups) that one could consider for addressing with causal/fancier sort of fairness techniques. Recent work using causality for debiasing is also appealing, and could be used in concert with our adversarial training. Finally, we are working on a prospective clinical trial to quantify our AI's role in lung cancer screening and the real-world use case will be the final test for the model's worth. We will also maintain and update our model under a continual learning framework, and apply relevant regulatory guard-rail (e.g. As described in FDA's proposed "adaptive" AI device guidance).

In summary, this study exhibits how a well-designed DL-based system can serve as an effective aid for lung cancer disease diagnosis from imaging, while fitting an approach for dealing with the numerous obstacles preventing the clinical penetration of AI. The proposed hybrid CNN-Transformer model in conjunction with bias reduction performed very well in both accuracy and fairness and the system can be transferred to clinical workflow. We think that responsibly developed and deployed AI will enhance the performance of radiologists, decrease diagnostic errors, and lead to higher rates of survival via earlier and better detection of cancers.

---

### References

- [1] World Health Organization, *Global Cancer Report 2020: Lung cancer incidence and mortality*, Geneva, Switzerland: WHO, 2020.
- [2] R. Siegel, et al., "Cancer statistics, 2021," *CA: A Cancer Journal for Clinicians*, vol. 71, no. 1, pp. 7–33, 2021.
- [3] National Lung Screening Trial Research Team, "Reduced lung-cancer mortality with low-dose CT screening," *New England Journal of Medicine*, vol. 365, no. 5, pp. 395–409, 2011.
- [4] S. G. Armato, et al., "The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A completed reference database of lung nodules on CT scans," *Medical Physics*, vol. 38, no. 2, pp. 915–931, 2011.
- [5] D. Ardila, et al., "End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest CT," *Nature Medicine*, vol. 25, pp. 954–961, 2019.
- [6] Y. Xie, et al., "Knowledge-based collaborative deep learning for benign-malignant lung nodule classification on chest CT," *IEEE Transactions on Medical Imaging*, vol. 38, no. 4, pp. 991–1004, 2019.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 770–778.
- [8] A. Dosovitskiy, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2021.

- [9] Z. Xu, et al., "Addressing fairness issues in deep learning-based medical image analysis: A systematic review," *NPJ Digital Medicine*, vol. 7, no. 1, pp. 1–13, 2024.
- [10] S. Langer, et al., "DeepTechnome: Mitigating unknown bias in deep learning-based assessment of CT images," *arXiv preprint*, arXiv:2206.12345, 2022.
- [11] S. Selvarajan, et al., "Enhancing lung cancer detection through integrated deep learning and transformer models," *Scientific Reports*, vol. 15, pp. 9231–9245, 2025.
- [12] H. Alsheikhy, et al., "Hybrid deep learning model for lung cancer classification using VGG19 and LSTM," *Journal of Imaging*, vol. 7, no. 12, pp. 1–15, 2021.
- [13] A. Zargar, et al., "Lung nodule classification on CT using transfer learning," *Diagnostics*, vol. 11, no. 9, pp. 1589–1602, 2021.
- [14] M. Kumaran, et al., "Ensemble of CNN models for lung cancer detection, with Grad-CAM interpretability," *IEEE Access*, vol. 11, pp. 102123–102135, 2023.
- [15] H. Wang, et al., "Possible bias in deep learning algorithms for CT lung nodule detection and classification," *Radiology: Artificial Intelligence*, vol. 4, no. 3, pp. e220042, 2022.
- [16] S. G. Finlayson, et al., "Adversarial attacks on medical machine learning," *Science*, vol. 363, no. 6433, pp. 1287–1289, 2019.
- [17] S. Keel, et al., "Explainable AI for lung cancer detection via a custom CNN on CT images," *Scientific Reports*, vol. 15, pp. 5012–5023, 2025.
- [18] J. Geis, et al., "Ethics of artificial intelligence in radiology: Summary of the joint European and North American multisociety statement," *Radiology*, vol. 293, no. 3, pp. 436–440, 2019.
- [19] European Commission, *Proposal for a Regulation on Artificial Intelligence (Artificial Intelligence Act)*, Brussels, Belgium: EC, 2021.
- [20] Y. Mirsky, et al., "CT-GAN: Malicious tampering of 3D medical imagery using deep learning," in *Proc. USENIX Security Symposium*, 2019, pp. 461–478.
- [21] H. Ahmed, et al., "Integrating AI into radiology workflow: IHE approach," *Journal of Digital Imaging*, vol. 35, pp. 1234–1245, 2022.
- [22] M. McDermott, et al., "Reproducibility in artificial intelligence research for radiology," *Radiology: Artificial Intelligence*, vol. 3, no. 2, pp. e200083, 2021.
- [23] DICOM Standards Committee, *Artificial Intelligence and DICOM*, White Paper, NEMA, 2020.
- [24] M. Mitchell, et al., "Model cards for model reporting," in *Proc. Conf. Fairness, Accountability, and Transparency (FAT)*, Atlanta, GA, USA, 2019, pp. 220–229.
- [25] U.S. Food and Drug Administration (FDA), *Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD) Action Plan*, Silver Spring, MD, USA: FDA, 2021.
- [26] L. Rundo, et al., "Adversarial training for AI in medical imaging: A survey," *Journal of Imaging*, vol. 5, no. 3, pp. 1–25, 2019.
- [27] U.S. National Institute of Standards and Technology (NIST), *The LIDC-IDRI Reference Database of Lung Nodules on CT*, Gaithersburg, MD, USA: NIST, 2015.