

# SPATIO-TEMPORAL HYBRID FUSION OF CAE AND SWIN TRANSFORMERS FOR LUNG CANCER MALIGNANCY PREDICTION

Sadaf Khademi<sup>†</sup>, Shahin Heidarian<sup>‡</sup>, Parnian Afshar<sup>†</sup>, Farnoosh Naderkhani<sup>†</sup>, Anastasia Oikonomou<sup>††</sup>, Konstantinos N. Plataniotis<sup>‡‡</sup>, and Arash Mohammadi<sup>†</sup>

<sup>†</sup>Concordia Institute for Information Systems Engineering, Concordia University, Montreal, Canada

<sup>‡</sup>Department of Electrical and Computer Engineering, Concordia University, Montreal, Canada

<sup>††</sup>Department of Medical Imaging, Sunnybrook Health Sciences Centre, Toronto, Canada

<sup>‡‡</sup>Department of Electrical and Computer Engineering, University of Toronto, Toronto, Canada

## ABSTRACT

The paper proposes a novel hybrid discovery Radiomics framework that simultaneously integrates temporal and spatial features extracted from non-thin chest Computed Tomography (CT) slices to predict Lung Adenocarcinoma (LUAC) malignancy with minimum expert involvement. Lung cancer is the leading cause of mortality from cancer worldwide and has various histologic types, among which LUAC has recently been the most prevalent. LUACs are classified as pre-invasive, minimally invasive, and invasive adenocarcinomas. Timely and accurate knowledge of the lung nodules malignancy leads to a proper treatment plan and reduces the risk of unnecessary or late surgeries. Currently, chest CT scan is the primary imaging modality to assess and predict the invasiveness of LUACs. However, the radiologists' analysis based on CT images is subjective and suffers from a low accuracy compared to the ground truth pathological reviews provided after surgical resections. The proposed hybrid framework, referred to as the CAET-SWin, consists of two parallel paths: (i) The Convolutional Auto-Encoder (CAE) Transformer path that extracts and captures informative features related to inter-slice relations via a modified Transformer architecture, and; (ii) The Shifted Window (SWin) Transformer path, which is a hierarchical vision transformer that extracts nodules' related spatial features from a volumetric CT scan. Extracted temporal (from the CAET-path) and spatial (from the Swin path) are then fused through a fusion path to classify LUACs. Experimental results on our in-house dataset of 114 pathologically proven Sub-Solid Nodules (SSNs) demonstrate that the CAET-SWin significantly improves reliability of the invasiveness prediction task while achieving an accuracy of 82.65%, sensitivity of 83.66%, and specificity of 81.66% using 10-fold cross-validation.

**Index Terms**— Lung Adenocarcinoma, Lung Nodule Invasiveness, Transformer, Subsolid Nodule, Self-Attention.

## 1. INTRODUCTION

Lung Cancer (LC) is the deadliest and least-funded cancer worldwide [1, 2]. Non-small-cell LC is the major type of LC, and Lung Adenocarcinoma (LUAC) is the most prevalent histologic subtype [3]. Lung nodules manifesting as Ground Glass (GG) or Sub-solid Nodules (SSNs) on Computed Tomography (CT) scans have a higher risk of malignancy than other incidentally detected small solid nodules. SSNs are often diagnosed as adenocarcinoma and are generally classified into pure GG nodules and part-solid nodules according to their appearance on the lung window settings [4, 5]. A timely and accurate attempt to differentiate the LUACs is of utmost importance to guide a proper treatment plan, as in some

cases, a pre-invasive or minimally invasive SSN can be monitored with regular follow-up CT scans, whereas invasive lesions should undergo immediate surgical resection if they are deemed eligible. Most often, the SSNs type is diagnosed based on the pathological findings performed after surgical resections, which is not desired for prior treatment planning. Currently, radiologists use chest CT scans to assess the invasiveness of the SSNs based on their imaging findings and patterns prior to making decisions regarding the appropriate treatment. Such visual approaches, however, are time-consuming, subjective, and error-prone. So far, many studies have used high-resolution and thin-slice ( $< 1.5mm$ ) CT images for the SSN classification, which require longer analysis times, as well as more reconstruction time [6, 7]. However, lung nodules are mostly identified from CT scans performed for varied clinical purposes acquired using routine standard or low-dose scanning protocols with non-thin slice thicknesses (up to  $5mm$ ) [8]. In addition, recent lung cancer screening recommendation, suggests using low-dose CT scans with thicker slice-thicknesses (up to  $2.5mm$ ) [9, 10]. Capitalizing on the above discussion, the necessity of developing an automated invasiveness assessment framework that performs well regardless of technical settings has recently arisen among the research community and healthcare professionals.

**Related Works:** Generally speaking, existing works on the SSN invasiveness assessment can be categorized into two main classes: (i) Radiomics-based, and; (ii) Deep Learning-based frameworks, also referred to as Discovery Radiomics [11]. In the former class, data-characterization algorithms extract quantitative features from nodule masks and the original CT images, which are then analyzed using statistical or conventional Machine Learning (ML) models [12, 13]. As an example of such frameworks, a histogram-based model is developed in [8] to predict the invasiveness of primary adenocarcinoma SSNs from non-thin CT scans of 109 pathologically labeled SSNs. In this study, a set of histogram-based and morphological features along with additional features extracted via the functional Principal Component Analysis (PCA) is fed to a linear logistic regression. Discovery Radiomics approaches, on the other hand, extract informative and discriminative features in an automated fashion. Existing deep models working with volumetric CT scans can be classified in two categories: (i) 3D-based solutions [14], where the whole 3D volume of CT images are fed to the model. Processing a large 3D CT scan at once, however, results in extensive computational complexity requiring more computational resources, and enormous training datasets, and; (ii) 2D-based solutions [15–17], where individual 2D CT slices are first analyzed, which are then fused via an aggregation

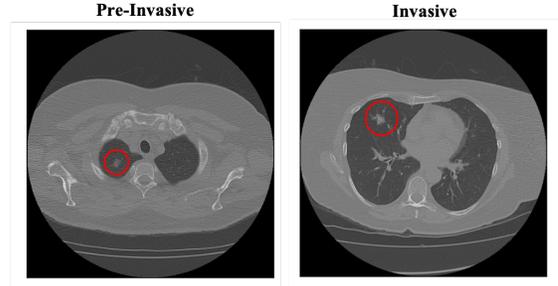
mechanism to represent characteristics of the whole volume.

Due to the time-series nature of the volumetric CT scans, which utilize a sequence of 2D images to provide a detailed representation of the organ, there has been recently a surge of interest in Category (ii), especially, in the application of sequential deep models for diagnostic/prognostic tasks based on 2D-CT scans. In 2017, a new sequential deep model, commonly known as ‘‘Transformer’’ [18], has been proposed illustrating superior performance in the tasks with sequential input data. Transformer models benefit from a novel self-attention mechanism, which is capable of capturing global context and dependencies between instances of the sequential data while requiring far less computational resources compared to conventional recurrent-based architectures such as Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM). Transformers are also superior to their counterparts in terms of parallelization and dynamic attention. Although the transformer model was initially designed for natural language processing, in 2020 Vision Transformer (ViT) was introduced to adopt the self-attention mechanism for image processing applications [19]. ViT-based models apply the self-attention mechanism to the sequence of patches extracted from a 2D image. It is worth mentioning that development of transformer-based models in processing sequential medical images is progressing rapidly especially models [20, 21] for COVID-19 identification and segmentation have shown promising results and potentials.

**Contributions:** The paper proposes a novel hybrid malignancy predictive framework, referred to as the CAET-SWin, which combines spatial and temporal features extracted by two parallel self-attention mechanisms (the CAET and SWin paths). Intuitively speaking, the CAET-SWin is designed to take advantage of the 3D nature of non-thin CT sequences by jointly modeling temporal (inter-slice) and spatial (within-slice) variations of CT slices with reduced computational complexity. Each of the constituent paths concentrates on a specific domain to capture local and global evidence of invasive nodules. The first path is built based on a Convolutional Auto-Encoder (CAE) algorithm [22] to form a sequential feature map of the whole lung provided in a set of CT slices. The obtained sequential feature maps are then directly fed to the encoder part of a transformer model containing multiple Multi-head Self Attention (MSA) layers to find temporal dependencies between slices. This path of the CAET-SWin has an intuitively pleasing modified structure, i.e., the patch flattening and linear projection components of the ViT architecture are removed and the output of the CAE is directly fed to the Position Embedding (PE) module. At the same time, nodule patches with each slice are analyzed in the second parallel path, which is a ViT with a hierarchical structure using the shifted windowing scheme of the SWin transformer. This path finds spatial connections among local windows formed for each image that significantly enhances spatial modeling power [23]. Finally, temporal and spatial features extracted in these parallel paths are fused via a stack of fully connected (FC) layers to provide the final predictions. An important aspect considered in the design of the proposed hybrid CAET-SWin framework is that it does not require a detailed annotation of the nodules needed in most existing models, which is a challenging and time-consuming task even for expert radiologists. The only required information from the radiologists/experts is the set of slices with evidence of a nodule and an approximate region containing a suspicious nodule.

## 2. MATERIALS AND METHODS

In this section, first, we briefly present the dataset utilized to develop and test the proposed CAET-SWin framework. Afterwards, we describe the lung segmentation module used as a pre-processing step to form the required input of the two parallel paths of the CAET-SWin



**Fig. 1.** Sample pre-invasive and invasive adenocarcinomas.

framework. Finally, we briefly introduce the self-attention mechanism, which is the main building block of both parallel paths of our hybrid framework.

### 2.1. Dataset

Most of recent studies [24] were developed and evaluated based on the public LIDC-IDRI [25] dataset, which does not have pathologically proven labels and focuses more on nodule detection than classification. In this study, we have used the dataset initially introduced in [8] and added five additional cases acquired from the same institution to further balance the dataset. This dataset contains non-thin volumetric CT scans of 114 pathologically proven SSNs (with technical parameters of 100–135 kVp and 80–120 mAs) segmented and reviewed by 2 experienced thoracic radiologists. All SSN labels are provided after surgical resections. SSNs are initially classified according to their histology into three categories of pre-invasive lesions including atypical adenomatous hyperplasia and adenocarcinoma in situ, minimally invasive, and invasive pulmonary adenocarcinoma [5]. Following the original study [8], we have grouped the first two categories to represent the pre-invasive and minimally invasive class with 58 cases, and kept the invasive nodules as the other class including 56 cases. Fig. 1 shows two sample lung adenocarcinomas from the dataset.

### 2.2. Lung Segmentation

As the pre-processing step, we have utilized a well-trained U-Net-based lung segmentation model, introduced in [26], to extract the lung parenchyma from the CT scans. This approach has been proven beneficial to enhance the learning process and final results of deep learning-based models in previous CT-related studies [15, 27, 28], by removing distracting components from the CT images. The extracted lung areas are then passed into two parallel paths (i.e., the CAET and the SWin paths) of the CAET-SWin framework as follows: (i) *Input of the CAET Path:* Extracted lung areas after the pre-processing step are down-sampled from (512, 512) to (256, 256) pixels to reduce the complexity and memory allocation without significant loss of information. (ii) *Input of the SWin Path:* Nodule patches within the segmented area resulting from the pre-processing step based on SSN annotation provided by radiologists and then zero-padded to (224, 224) pixels.

### 2.3. Multi-Head Self-Attention Mechanism

Transformer architecture constitutes the main component of the two parallel paths of the CAET-SWin framework, which uses a self-attention mechanism to capture dependencies among various instances of the input sequence. Self-attention is developed based on a scaled dot-product attention function, mapping a query and a set of key-value pairs to an output. The query ( $Q$ ), keys ( $K$ ), and values ( $V$ ) are learnable representative vectors for the instances in the input sequence with dimensions  $d_k$ ,  $d_k$ , and  $d_v$ , respectively. The output of the self-attention module is computed as a weighted average of the values, where the weight assigned to each value is computed by a similarity function of the query and the correspond-

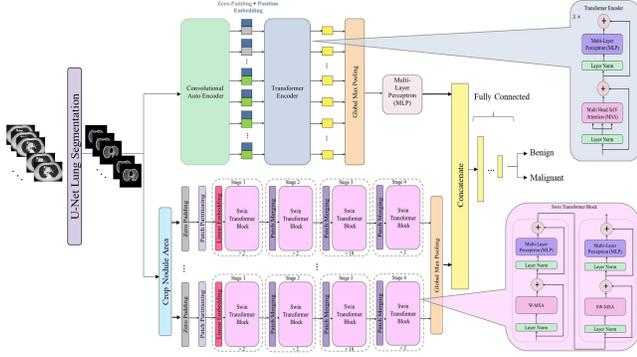


Fig. 2. Pipeline of the CAET-SWin Transformer.

ing key after applying a softmax function [18]. More specifically, the attention values on a set of queries are computed simultaneously, packed together into matrix  $Q$ . The keys and values are similarly represented by matrices  $K$  and  $V$ . The output of the attention function is computed as follows

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (1)$$

where superscript  $T$  denotes transpose of a given matrix. It is also beneficial to linearly project the queries, keys, and values  $h$  times with various learnable linear projections to vectors with  $d_k$ ,  $d_k$  and  $d_v$  dimensions, respectively, before applying the attention function. On each of the projected versions of queries, keys, and values, the attention function is performed in parallel, resulting in  $d_v - \text{dimensional}$  output values. These values are then concatenated and once again linearly projected via a FC layer. This process is called Multi-head Self Attention (MSA), which helps the model to jointly attend to information from different representation subspaces at different positions [18]. The output of the MSA module is

$$\text{MSA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)\mathbf{W}^O, \quad (2)$$

where  $\text{head}_i = \text{Attention}(Q\mathbf{W}_i^Q, K\mathbf{W}_i^K, V\mathbf{W}_i^V)$ , (2)

where the projections are achieved by parameter matrices  $\mathbf{W}_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $\mathbf{W}_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $\mathbf{W}_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ , and  $\mathbf{W}^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$ . This completes an overview of the preliminary material, next, we develop the CAET-SWin framework.

### 3. THE PROPOSED CAET-SWin FRAMEWORK

The proposed hybrid CAET-SWin framework combines spatial and temporal features extracted by two parallel self-attention mechanisms (the so called CAET and Swin paths) to perform malignancy prediction based on CT images. Through its hybrid architecture, the CAET-SWin learns from the 3D structure of non-thin CT scans by simultaneous extraction of inter-slice (through its CAET path) and within-slice (through its Swin path) features. Extracted features are then fused through a FC fusion path (implemented in series to the CAET and Swin parallel paths) to form the output. In what follows, we present each of the three paths of the CAET-SWin framework.

#### 3.1. The CAE-Transformer Path

To extract temporal (inter-slice) features, the first feature extractor module of the proposed hybrid CAET-SWin is named CAE-Transformer, which integrates CAE and a modified version of the transformer encoder [18, 19]. Fig. 2 illustrates the pipeline of the CAE-Transformer framework, along with the architecture of a transformer encoder. To develop the CAE-Transformer module, we initially pre-trained a CAE model based on CT images with the evidence of a nodule in public LIDC-IDRI dataset. This is performed to represent CT images with compressed and informative feature maps.

Table 1. Results obtained by the CAET-SWin and its counterparts.

Model	Accuracy [95% CI]	Sensitivity %	Specificity %
Ref. [8]	81.00 [58.1 94.6]	80.00	<b>81.80</b>
Ref. [20]	61.66 [49.9 73.5]	61.33	62.33
CAET (GAP)	56.21 [46.7 65.7]	42.00	69.66
CAET (Flatten)	64.84 [57.9 71.7]	54.66	74.66
CAET (GMP)	69.46 [58.8 80.2]	64.33	74.66
SWin	78.10 [70.0 86.2]	76.66	79.66
CAET-SWin	<b>82.65 [75.6 89.8]</b>	<b>83.66</b>	81.66

The CAE model consists of an encoder and a decoder part, where the former is responsible for generating a compact representation of the input image through a stack of five convolution and five max-pooling layers (kernel size of  $2 \times 2$ ) followed by a FC layer with the size of 256. The decoder component aims to reconstruct the original image using the compressed feature representation generated by the encoder. By minimizing the Mean Squared Error (MSE) between the original and the reconstructed image, the CAE learns to produce highly informative feature representations for the input images. The pre-trained model on the LIDC-IDRI dataset is then fine-tuned on our in-house dataset (presented in Section 2).

The CAE component plays the role of the embedding layer in basic transformer architecture. In other words, instead of patch flattening and linear projection, which is commonly used in the transformer architecture, the output of the CAE is directly fed to the PE module. As the number of slices with the evidence of a nodule varies between different subjects (from 2 to 25 slices per nodule), we have taken the maximum number of slices in our dataset (i.e., 25 slices) and zero-padded the input sequences based on this number such that all sequences provided by CAE would have the same dimension of (25, 256). The PE layer is then incorporated into the model to add information about the position of slices in the input sequence. More specifically, a transformer encoder is initialized by applying the MSA on the normalized CAE-generated feature maps corresponding to the input instances, followed by a residual connection, which adds low-level features of the input to the output of the MSA module. A Layer Normalization (LN) is then applied to the results. The normalized values are then passed to the next module, which contains a Multi-Layer Perceptron (MLP), followed by another residual connection as shown in Fig. 2.

The CAE-Transformer path is constructed by stacking 3 transformer encoder blocks on top of each other with a projection dimension of 256, key and query dimensions of 128, and 5 heads in each MSA module. Finally, the outputs obtained by stack of transformer encoders from all input instances (slices) are passed to a Global Max Pooling (GMP) layer followed by an MLP layer with 32 neurons and a Relu activation function to generate the feature vectors.

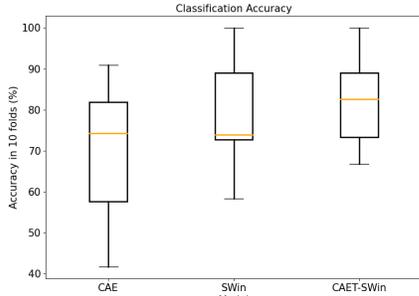
#### 3.2. The Swin-Transformer Path

The second path of the CAET-SWin framework consists of a Swin transformer (SWin-B architecture [29]) followed by a GMP layer that concentrates on local variations of nodules related to each subject. The Swin transformer is a hierarchical transformer that uses shifting-window MSA and includes four stages as shown in Fig. 2. First, the input image is divided into  $4 \times 4$  patches. In Stage 1, the feature dimension of each patch is projected into 128 by a linear embedding layer and then Swin transformer blocks are applied on the patches. The hierarchical structure is built based on a patch merging layer at the beginning of Stages 2-4 reducing the number of patches. This layer concatenates the features of  $2 \times 2$  neighboring patches and applies a linear layer. Afterward, several Swin transformer blocks are employed for feature transformation at each stage.

All layers of a Swin transformer are similar to the original transformer except for the MSA module, which is replaced based on a

**Table 2.** 10 fold cross-validation accuracy of CAET-SWin framework and its constituent parts.

Model	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Fold 6	Fold 7	Fold 8	Fold 9	Fold 10
CAET	81.82	54.54	45.45	90.91	41.67	66.67	<b>83.33</b>	66.67	81.82	81.82
SWin	90.91	72.73	63.64	100	<b>91.67</b>	75.00	58.33	83.33	72.73	72.73
CAET-SWin	<b>90.91</b>	<b>72.73</b>	<b>72.73</b>	<b>100</b>	83.33	<b>75.00</b>	66.67	<b>83.33</b>	<b>100</b>	<b>81.82</b>

**Fig. 3.** Boxplots for CAET-SWin and its constituent parts.

shifting window mechanism to calculate self-attention within local windows of size  $M \times M$ . The advantage of this method is to find cross-window connections and reduce computation complexity. As shown in Fig. 2, the Swin transformer block of the CAET-SWin consists of two modules, W-MSA and SW-MSA. W-MSA module is a regular partitioning scheme that divides the image into non-overlapping windows and self-attention is computed in each window. In the next step, SW-MSA computes self-attention in new windows generated by shifting the W-MSA windows by  $(\lfloor \frac{M}{2} \rfloor, \lfloor \frac{M}{2} \rfloor)$  pixels. The output of Stage 4 is the feature vector related to each slice fed to a GMP layer to aggregate the volume-level features of each subject. Here, we fine-tune the weights of a pre-trained Swin-B transformer trained on ImageNet-21k dataset [30] with input image size of  $224 \times 224$  and window size of  $7 \times 7$ .

### 3.3. Fusion Path

The output of the CAET path consists of 32 features capturing temporal relations of CT slices associated with a patient. On the other hand, the output of the Swin path includes 1024 features modeling inter-slice (spatial) correlations in different CT slices of a given subject. In order to take advantage of slice information in both spatial and temporal domains, we concatenated the two output vectors to form the final feature vector, which is passed through 4 FC layers with 512, 128, 32, and 2 neurons. The last FC layer uses a Softmax activation function to produce probability scores for the two classes.

## 4. EXPERIMENTAL RESULTS

We evaluated the performance of the proposed CAET-SWin Transformer framework using the 10-fold cross-validation method. It should be noted that the models presented in parallel paths were trained independently. The CAE model was pre-trained using a batch size of 128, learning rate of  $1e-4$ , and 200 epochs. The best model on the randomly sampled 20% of the dataset was selected as the candidate model. The model was then fine-tuned on the in-house dataset using a lower learning rate of  $1e-6$  and 50 epochs. To fine-tune the final CAE, only the middle FC layer and its previous and next convolution layers were trained while the other layers have been kept unchanged. The CAE-generated features were then used to train the transformer encoder. The transformer was trained using the Adam optimizer with a learning rate of  $1e-3$ , label smoothing with the  $\alpha = 0.1$ , and 200 epochs. Simultaneously, the pre-trained Swin-B transformer was trained by AdamW optimizer with a learning rate of  $1e-5$ , weight-decay of 0.05, and 50 epochs with early stopping training strategy (patience = 10). At the last step, FC layers

were trained by means of the Adam optimizer with a learning rate of  $1e-2$  and 20 epochs. Also, dropout layers were incorporated to prevent the model from getting over-fitted. The loss function used in the whole process was cross entropy. The classification results of the CAET-SWin framework are presented in Table 1 in terms of accuracy, sensitivity, and specificity.

We have compared performance of the proposed CAET-SWin framework with the results obtained by the model proposed in [8,20] and stand-alone models in parallel paths of our hybrid model. We added a FC layer with 2 neurons right after the last layer of each model with a Softmax activation function to classify SSNs and evaluate effects of each feature set. The experimental results provided in Table 1 show that simultaneous attention to both time and space domains empowers the overall model in such a way that CAET-SWin achieved the best performance among its constituent parts in all three evaluation metrics. More details regarding the performance of each fold is provided in the Table 2 and the distribution of classification accuracy for these 3 models is shown in Fig. 3. Additionally, we implemented two modified versions of CAET by replacing the aggregation method with the Global Average Pooling (GAP) layer and Flattening layer. However, as presented in Table 1 CAET utilizing GMP layer outperformed mentioned models. To compare the ability of the CAE algorithm in extracting efficient features fed into the transformer, we employed a pre-trained basic architecture ViT in a voting scheme presented in [20] to classify SSNs and results demonstrated the superiority of CAET in using CAE instead of patch flattening and linear projection compared to ViT model.

As stated previously, correctly identifying the invasiveness level of nodules could have a great impact on the treatment plan and its success. Therefore, the correct diagnosis of malignant nodules is relatively more important than detection of early stages of nodule transmutation. From this point of view, sensitivity would be a more capable evaluation metric than other criteria. The CAET-SWin achieved a sensitivity of 83.66% which is about 4% higher than the reference study [8] while specificity is kept at the same value. In other words, in our hybrid model fewer cases of malignancy are missed. Hence, potentially it could be presented as a more reliable algorithm for recognizing malignant nodules. Furthermore, the confidence interval (CI) which describes the uncertainty level of a model is narrower for CAET-SWin illustrating that the proposed hybrid model is more precise/reliable than its radiomics counterpart.

## 5. CONCLUSION

The paper proposed a hybrid transformer-based framework, referred to as the CAET-SWin, to accurately and reliably predict the invasiveness of lung adenocarcinoma subsolid nodules from non-thin 3D CT scans. The proposed CAET-SWin model achieves this objective by combining spatial (within-slice) and temporal (inter-slice) features extracted by its two constituent parallel paths (the CAET and Swin paths) designed based on the self-attention mechanism. The CAET-SWin significantly improved reliability of the invasiveness prediction task compared to its radiomics-based counterpart while increasing the accuracy by 1.65% and sensitivity by 3.66%. Investigating effects of embedding radiomics and morphological features in the CAET-SWin framework is a fruitful direction for future research.

## 6. REFERENCES

- [1] S.D. Kamath, S.M. Kircher, and A.B. Benson, "Comparison of Cancer Burden and Nonprofit Organization Funding Reveals Disparities in Funding Across Cancer Types," *Journal of the National Comprehensive Cancer Network*, vol. 17, no. 7, pp. 849–854, jul 2019.
- [2] F. Bray et al., "Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries," *CA: A Cancer Journal for Clinicians*, vol. 68, no. 6, pp. 394–424, nov 2018.
- [3] R.S. Herbst, D. Morgensztern, and C. Boshoff, "The biology and management of non-small cell lung cancer," *Nature*, vol. 553, no. 7689, pp. 446–454, jan 2018.
- [4] H.Y. Kim, Y.M. Shim, K.S. Lee, J. Han, C.A. Yi, and Y.K. Kim, "Persistent Pulmonary Nodular Ground-Glass Opacity at Thin-Section CT: Histopathologic Comparisons," *Radiology*, vol. 245, no. 1, pp. 267–275, oct 2007.
- [5] J. Lai, Q. Li, F. Fu, Y. Zhang, Y. Li, Q. Liu, and H. Chen, "Subsolid Lung Adenocarcinomas: Radiological, Clinical and Pathological Features and Outcomes," *Seminars in Thoracic and Cardiovascular Surgery*, jun 2021.
- [6] X. Cui et al., "A Subsolid Nodules Imaging Reporting System (SSN-IRS) for Classifying 3 Subtypes of Pulmonary Adenocarcinoma," *Clinical Lung Cancer*, vol. 21, no. 4, pp. 314–325.e4, jul 2020.
- [7] X. Shao, R. Niu, Z. Jiang, X. Shao, and Y. Wang, "Role of PET/CT in Management of Early Lung Adenocarcinoma," *American Journal of Roentgenology*, vol. 214, no. 2, pp. 437–445, feb 2020.
- [8] A. Oikonomou and otehrs, "Histogram-based models on non-thin section chest CT predict invasiveness of primary lung adenocarcinoma subsolid nodules," *Scientific Reports*, vol. 9, no. 1, pp. 6009, dec 2019.
- [9] E.A. Kazerooni and otehrs, "ACR–STR Practice Parameter for the Performance and Reporting of Lung Cancer Screening Thoracic Computed Tomography (CT)," *Journal of Thoracic Imaging*, vol. 29, no. 5, pp. 310–316, sep 2014.
- [10] K. Fujii, K. McMillan, M. Bostani, C. Cagnon, and M. McNitt-Gray, "Patient Size–Specific Analysis of Dose Indexes From CT Lung Cancer Screening," *American Journal of Roentgenology*, vol. 208, no. 1, pp. 144–149, jan 2017.
- [11] D. Gu, G. Liu, and Z. Xue, "On the performance of lung nodule detection, segmentation and classification," *Computerized Medical Imaging and Graphics*, vol. 89, pp. 101886, apr 2021.
- [12] C. Gao, P. Xiang, J. Ye, P. Pang, S. Wang, and M. Xu, "Can texture features improve the differentiation of infiltrative lung adenocarcinoma appearing as ground glass nodules in contrast-enhanced CT?," *European Journal of Radiology*, vol. 117, pp. 126–131, aug 2019.
- [13] J. Uthoff et al., "Machine learning approach for distinguishing malignant and benign lung nodules utilizing standardized perinodular parenchymal features from CT," *Medical Physics*, vol. 46, no. 7, pp. 3207–3216, jul 2019.
- [14] S. Liu, Y. Xie, A. Jirapatnakul, and A.P. Reeves, "Pulmonary nodule classification in lung cancer screening with three-dimensional convolutional neural networks," *Journal of Medical Imaging*, vol. 4, no. 04, pp. 1, nov 2017.
- [15] S. Heidarian et al., "COVID-FACT: A Fully-Automated Capsule Network-Based Framework for Identification of COVID-19 Cases from Chest CT Scans," *Frontiers in Artificial Intelligence*, vol. 4, may 2021.
- [16] S. Heidarian et al., "Ct-Caps: Feature Extraction-Based Automated Framework for Covid-19 Disease Identification From Chest Ct Scans Using Capsule Networks," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. jun 2021, pp. 1040–1044, IEEE.
- [17] M.M. Farhangi, N. Petrick, B. Sahiner, H. Frigui, A.A. Amini, and A. Pezeshk, "Recurrent attention network for false positive reduction in the detection of pulmonary nodules in thoracic CT scans," *Medical Physics*, vol. 47, no. 5, pp. 2150–2160, may 2020.
- [18] A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. 2017, vol. 30, Curran Associates, Inc.
- [19] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," oct 2020.
- [20] X. Gao, Y. Qian, and A. Gao, "COVID-VIT: Classification of COVID-19 from CT chest images based on vision transformer models," jul 2021.
- [21] A.A.E. Ambita, E.N.V. Boquio, and P.C. Naval, "COViT-GAN: Vision Transformer for COVID-19 Detection in CT Scan Images with Self-Attention GAN for Data Augmentation," pp. 587–598. 2021.
- [22] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, "Stacked Convolutional Auto-Encoders for Hierarchical Feature Extraction," pp. 52–59. 2011.
- [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 10012–10022.
- [24] P. Afshar, et al., "3D-MCN: A 3D Multi-scale Capsule Network for Lung Nodule Malignancy Prediction," *Nature Scientific Reports*, vol. 10, 7948, 2020.
- [25] 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, jan 2011.
- [26] J. Hofmanninger, F. Prayer, J. Pan, S. Röhrich, H. Prosch, and G. Langs, "Automatic lung segmentation in routine imaging is primarily a data diversity problem, not a methodology problem," *European Radiology Experimental*, vol. 4, no. 1, pp. 50, dec 2020.
- [27] P. Afshar et al., "Human-level COVID-19 Diagnosis from Low-dose CT Scans Using a Two-stage Time-distributed Capsule Network," may 2021.
- [28] A. Mohammadi et al., "Diagnosis/Prognosis of COVID-19 Chest Images via Machine Learning and Hypersignal Processing: Challenges, opportunities, and applications," *IEEE Signal Processing Magazine*, vol. 38, no. 5, pp. 37–66, sep 2021.
- [29] Lei Zhang and Yan Wen, "A transformer-based framework for automatic covid19 diagnosis in chest cts," in *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2021, pp. 513–518.
- [30] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.