




# CLIP-Lung: Textual Knowledge-Guided Lung Nodule Malignancy Prediction

Yiming Lei<sup>1</sup>, Zilong Li<sup>1</sup>, Yan Shen<sup>2</sup>, Junping Zhang<sup>1</sup>,  
and Hongming Shan<sup>3,4,5</sup> 

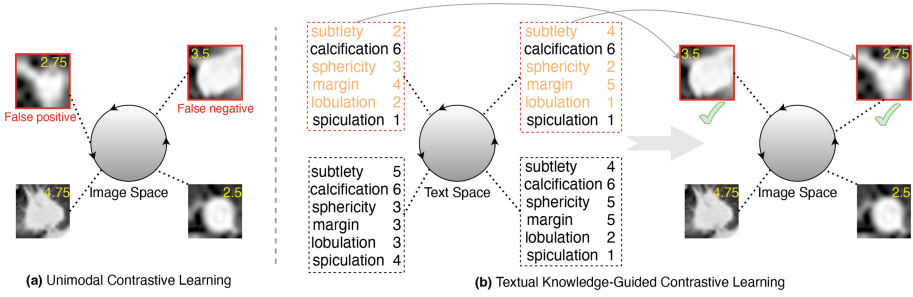
- <sup>1</sup> Shanghai Key Lab of Intelligent Information Processing, School of Computer Science, Fudan University, Shanghai, China  
<sup>2</sup> School of Pharmacy, China Pharmaceutical University, Nanjing, China  
<sup>3</sup> Institute of Science and Technology for Brain-Inspired Intelligence and MOE Frontiers Center for Brain Science, Fudan University, Shanghai, China  
[hmsan@fudan.edu.cn](mailto:hmsan@fudan.edu.cn)  
<sup>4</sup> Key Laboratory of Computational Neuroscience and Brain-Inspired Intelligence (Fudan University), Ministry of Education, Shanghai, China  
<sup>5</sup> Shanghai Center for Brain Science and Brain-Inspired Technology, Shanghai, China

**Abstract.** Lung nodule malignancy prediction has been enhanced by advanced deep-learning techniques and effective tricks. Nevertheless, current methods are mainly trained with cross-entropy loss using one-hot categorical labels, which results in difficulty in distinguishing those nodules with closer progression labels. Interestingly, we observe that clinical text information annotated by radiologists provides us with discriminative knowledge to identify challenging samples. Drawing on the capability of the contrastive language-image pre-training (CLIP) model to learn generalized visual representations from text annotations, in this paper, we propose CLIP-Lung, a textual knowledge-guided framework for lung nodule malignancy prediction. First, CLIP-Lung introduces both class and attribute annotations into the training of the lung nodule classifier without any additional overheads in inference. Second, we design a channel-wise conditional prompt (CCP) module to establish consistent relationships between learnable context prompts and specific feature maps. Third, we align image features with both class and attribute features via contrastive learning, rectifying false positives and false negatives in latent space. Experimental results on the benchmark LIDC-IDRI dataset demonstrate the superiority of CLIP-Lung, in both classification performance and interpretability of attention maps. Source code is available at <https://github.com/ymLeiFDU/CLIP-Lung>.

**Keywords:** Lung nodule classification · vision-language model · prompt learning

## 1 Introduction

Lung cancer is one of the most fatal diseases worldwide, and early diagnosis of the pulmonary nodule has been identified as an effective measure to prevent lung cancer. Deep learning-based methods for lung nodule classification

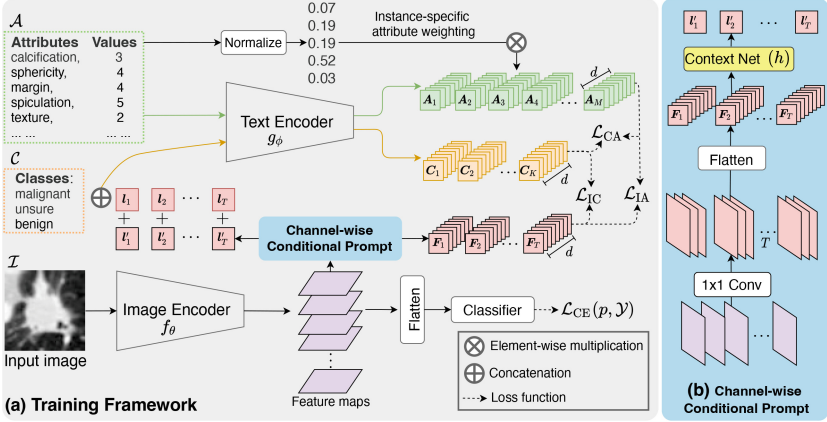


**Fig. 1.** Motivation of CLIP-Lung. (a) Unimodal contrastive learning. (b) Proposed textual knowledge-guided contrastive learning. Yellow values are the annotated malignancy scores. Dashed boxes contain pairs of textual attributes and annotated values. (Color figure online)

have been widely studied in recent years [9, 12]. Usually, the malignancy prediction is often formulated as benign-malignant binary classification [9, 10, 19], and the higher classification performance and explainable attention maps are impressive. Most previous works employ a learning paradigm that utilizes cross-entropy loss between predicted probability distributions and ground-truth one-hot labels. Furthermore, inspired by ordered labels of nodule progression, researchers have turned their attention to ordinal regression methods to evaluate the benign-unsure-malignant classification task [2, 11, 13, 18, 21], where the training set additionally includes nodules with uncertain labels. Indeed, the ordinal regression-based methods are able to learn ordered manifolds and to further enhance the prediction accuracy.

However, the aforementioned methods still face challenges in distinguishing visually similar samples with adjacent rank labels. For example, in Fig. 1(a), since we conduct unimodal contrastive learning and map the samples onto a spherical space, the false positive nodule with a malignancy score of 2.75 has a closer distance to that with a score of 4.75, and the false negative one should not be closer to that of score 2.5. To address this issue, we found that the text attributes, such as “subtlety”, “sphericity”, “margin”, and “lobulation”, annotated by radiologists, can exhibit the differences between these hard samples. Therefore, we propose leveraging text annotations to guide the learning of visual features. In practice, this also aligns with the fact that the annotated text information represents the direct justification for identifying lesion regions in the clinic. As shown in Fig. 1, this text information is beneficial for distinguishing visually similar pairs, while we conduct this behavior by applying contrastive learning that pulls semantic-closer samples and pushes away semantic-farther ones.

To integrate text annotations into the image-domain learning process, an effective text encoder providing accurate textual features is required. Fortunately, recent advances in vision-language models, such as contrastive language-image pre-training (CLIP) [16], provide us with a powerful text encoder pre-trained with text-based supervisions and have shown impressive results in downstream vision tasks. Nevertheless, it is ineffective to directly transfer CLIP to medical tasks due to the data covariate shift. Therefore, in this paper, we pro-



**Fig. 2.** Illustration of the proposed CLIP-Lung.

pose CLIP-Lung, a framework to classify lung nodules using image-text pairs. Specifically, CLIP-Lung constructs learnable text descriptions for each nodule from both class and attribute perspectives. Inspired by CoCoOp [20], we propose a channel-wise conditional prompt (CCP) module to allow nodule descriptions to guide the generation of informative feature maps. Different from CoCoOp, CCP constructs specific learnable prompts conditioned on grouped feature maps and triggers more explainable attention maps such as Grad-CAM [17], whereas CoCoOp provides only the common condition for all the prompt tokens. Then, we design a textual knowledge-guided contrastive learning based on obtained image features and textual features involving classes and attributes. Experimental results on LIDC-IDRI [1] dataset demonstrate the effectiveness of learning with textual knowledge for improving lung nodule malignancy prediction.

The contributions of this paper are summarized as follows.

- 1) We propose CLIP-Lung for lung nodule malignancy prediction, which leverages clinical textual knowledge to enhance the image encoder and classifier.
- 2) We design a channel-wise conditional prompt module to establish consistent relationships among the correlated text tokens and feature maps.
- 3) We simultaneously align the image features with class and attribute features through contrastive learning while generating more explainable attention maps.

## 2 Methodology

### 2.1 Overview

**Problem Formulation.** In this paper, we arrange the lung nodule classification dataset as  $\{\mathcal{I}, \mathcal{Y}, \mathcal{C}, \mathcal{A}\}$ , where  $\mathcal{I} = \{\mathcal{I}_i\}_{i=1}^N$  is an image set containing  $N$  lung nodule images.  $\mathcal{Y} = \{y_i\}_{i=1}^N$  is the corresponding class label set and  $y_i \in \{1, 2, \dots, K\}$ ,

and  $K$  is the number of classes.  $\mathcal{C} = \{\mathbf{c}_k\}_{k=1}^K$  is a set of text embeddings of classes. Finally,  $\mathcal{A} = \{\mathbf{a}_m\}_{m=1}^M$  is the set of attribute embeddings, where each element  $\mathbf{a}_m \in \mathbb{R}^{d \times 1}$  is a vector representing the embedding of an attribute word such as “spiculation”. Then, for a given sample  $\{\mathbf{I}_i, y_i\}$ , our aim is to learn a mapping  $f_\theta : \mathbf{I}_i \mapsto y_i$ , where  $f$  is a deep neural network parameterized by  $\theta$ .

**CLIP-Lung.** In Fig. 2(a), the training framework contains an image encoder  $f_\theta$  and a text encoder  $g_\phi$ . First, the input image  $\mathbf{I}_i$  is fed into  $f_\theta$  and then generates the feature maps. According to Fig. 2(b), the feature maps are converted to channel-wise feature vectors  $f_\theta(\mathbf{I}_i) = \mathbf{F}_{t,:}$  and then to learnable tokens  $\mathbf{l}'_t$ . Second, we initialize the context tokens  $\mathbf{l}_t$  and add them with  $\mathbf{l}'_t$  to construct the learnable prompts, where  $T$  is the number of context words. Next, the concatenation of the class token and  $\mathbf{l}_t + \mathbf{l}'_t$  is used as input of text encoder yielding the class features  $g_\phi(\mathbf{c}_k) = \mathbf{C}_{k,:}$ , note that  $\mathbf{C}_{k,:}$  is conditioned on channel-wise feature vectors  $\mathbf{F}_{t,:}$ . Finally, the attribute tokens  $\mathbf{a}_m$  are also fed into the text encoder to yield corresponding attribute features  $g_\phi(\mathbf{a}_m) = \mathbf{A}_{m,:}$ . Note that the vectors  $\mathbf{F}_{t,:}$ ,  $\mathbf{l}_{t,:}$ ,  $\mathbf{l}'_{t,:}$ , and  $\mathbf{C}_{k,:}$  are with the same dimension  $d = 512$  in this paper. Consequently, we have image feature  $\mathbf{F} \in \mathbb{R}^{T \times d}$ , class feature  $\mathbf{C} \in \mathbb{R}^{K \times d}$ , and attribute feature  $\mathbf{A} \in \mathbb{R}^{M \times d}$  to conduct the textual knowledge-guided contrastive learning.

## 2.2 Instance-Specific Attribute Weighting

For the attribute annotations, all the lung nodules in the LIDC-IDRI dataset are annotated with the *same* eight attributes: “subtlety”, “internal structure”, “calcification”, “sphericity”, “margin”, “lobulation”, “spiculation”, and “texture” [4, 8], and the annotated value for each attribute ranges from 1 to 5 except for “calcification” that is ranged from 1 to 6. In this paper, we fix the parameters of a pre-trained text encoder so that the generated eight text feature vectors are the same for all the nodules. Therefore, we propose an instance-specific attribute weighting scheme to distinguish different nodules. For the  $i$ -th sample, the weight for each  $\mathbf{a}_m$  is calculated through normalizing the annotated values:

$$w_m = \frac{\exp(v_m)}{\sum_{m=1}^M \exp(v_m)}, \quad (1)$$

where  $v_m$  denotes the annotated value for  $\mathbf{a}_m$ . Then the weight vectors of the  $i$ -th sample is represented as  $\mathbf{w}_i = [w_1, w_2, \dots, w_M]^\top \in \mathbb{R}^{M \times 1}$ . Hence, the element-wise multiplication  $\mathbf{w}_i \cdot \mathbf{A}_i$  is unique to  $\mathbf{I}_i$ .

## 2.3 Channel-Wise Conditional Prompt

CoCoOp [20] firstly proposed to learn language contexts for vision-language models conditioned on visual features. However, it is inferior to align context words with partial regions of the lesion. Therefore, we propose a channel-wise conditional prompt (CCP) module, in Fig. 2(b), to split latent feature maps into  $T$  groups and then flatten them into vectors  $\mathbf{F}_{t,:}$ . Next, we denote  $h(\cdot)$  as a context

network that is composed of a multi-layer perceptron (MLP) with one hidden layer, and each learnable context token is now obtained by  $\mathbf{l}'_t = h(\mathbf{F}_{t,:})$ . Hence, the conditional prompt for the  $t$ -th token is  $\mathbf{l}_t + \mathbf{l}'_t$ . In addition, CCP also outputs the  $\mathbf{F}_{t,:}$  for image-class and image-attribute contrastive learning.

## 2.4 Textual Knowledge-Guided Contrastive Learning

Recall that our aim is to enable the visual features to be similar to the textual features of the annotated classes or attributes and be dissimilar to those of irrelevant text annotations. Consequently, we accomplish this goal through contrastive learning [3, 5, 7]. In this paper, we conduct such image-text contrastive learning by utilizing pre-trained CLIP text encoder [16]. In Fig. 2, we align  $\mathbf{F} \in \mathbb{R}^{T \times d}$  with  $\mathbf{C} \in \mathbb{R}^{K \times d}$  and  $\mathbf{A} \in \mathbb{R}^{M \times d}$ , *i.e.*, using class and attribute knowledge to regularize the feature maps.

**Image-Class Alignment.** First, the same to CLIP, we align the image and class information by minimizing the cross-entropy (CE) loss for the sample  $\{\mathbf{I}_i, y_i\}$ :

$$\mathcal{L}_{\text{IC}} = - \sum_{t=1}^T \sum_{k=1}^K y_i \log \frac{\exp(\sigma(\mathbf{F}_{t,:}, \mathbf{C}_{k,:})/\tau)}{\sum_{k'=1}^K \exp(\sigma(\mathbf{F}_{t,:}, \mathbf{C}_{k',:})/\tau)}, \quad (2)$$

where  $\mathbf{C}_{k,:} = g_\phi(\mathbf{c}_k \oplus (\mathbf{l}_1 + \mathbf{l}'_1, \mathbf{l}_2 + \mathbf{l}'_2, \dots, \mathbf{l}_T + \mathbf{l}'_T)) \in \mathbb{R}^{d \times 1}$  and “ $\oplus$ ” denotes concatenation, *i.e.*,  $\mathbf{C}_{k,:}$  is conditioned on learnable prompts  $\mathbf{l}_t + \mathbf{l}'_t$ .  $\sigma(\cdot, \cdot)$  calculates the cosine similarity and  $\tau$  is the temperature term. Therefore,  $\mathcal{L}_{\text{IC}}$  implements the contrastive learning between channel-wise features and corresponding class features, *i.e.*, the ensemble of grouped image-class alignment results.

**Image-Attribute Alignment.** In addition to image-class alignment, we further expect the image features to correlate with specific attributes. So we conduct image-attribute alignment by minimizing the InfoNCE loss [5, 16]:

$$\mathcal{L}_{\text{IA}} = - \sum_{t=1}^T \sum_{m=1}^M \log \frac{\exp(\sigma(\mathbf{F}_{t,:}, \mathbf{w}_{m,:} \cdot \mathbf{A}_{m,:})/\tau)}{\sum_{m'=1}^M \exp(\sigma(\mathbf{F}_{t,:}, \mathbf{w}_{m',:} \cdot \mathbf{A}_{m',:})/\tau)}. \quad (3)$$

Hence,  $\mathcal{L}_{\text{IA}}$  indicates which attribute the  $\mathbf{F}_{t,:}$  is closest to since each vector  $\mathbf{F}_{t,:}$  is mapped from the  $t$ -th group of feature maps through the context network  $h(\cdot)$ . Therefore, certain feature maps can be guided by specific annotated attributes.

**Class-Attribute Alignment.** Although the image features have been aligned with classes and attributes, the class embeddings obtained by the pre-trained CLIP encoder may shift in the latent space, which may result in inconsistent class space and attribute space, *i.e.*, annotated attributes do not match the corresponding classes, which is contradictory to the actual clinical diagnosis. To avoid this weakness, we further align the class and attribute features:

$$\mathcal{L}_{\text{CA}} = - \sum_{k=1}^K \sum_{m=1}^M \log \frac{\exp(\sigma(\mathbf{C}_{k,:}, \mathbf{w}_{m,:} \cdot \mathbf{A}_{m,:})/\tau)}{\sum_{m'=1}^M \exp(\sigma(\mathbf{C}_{k,:}, \mathbf{w}_{m',:} \cdot \mathbf{A}_{m',:})/\tau)}, \quad (4)$$

and this loss implies semantic consistency between classes and attributes.

Finally, the total loss function is defined as follows:

$$\mathcal{L} = \mathbb{E}_{I_i \in \mathcal{I}} [\mathcal{L}_{\text{CE}} + \mathcal{L}_{\text{IC}} + \alpha \cdot \mathcal{L}_{\text{IA}} + \beta \cdot \mathcal{L}_{\text{CA}}], \quad (5)$$

where  $\alpha$  and  $\beta$  are hyperparameters for adjusting the losses and are set as 1 and 0.5, respectively.  $\mathcal{L}_{\text{CE}}$  denotes the cross-entropy loss between predicted probabilities obtained by the classifier and the ground-truth labels. Note that during the inference phase, test images are only fed into the trained image encoder and classifier. As a result, CLIP-Lung does not introduce any additional computational overhead in inference.

### 3 Experiments

#### 3.1 Dataset and Implementation Details

**Dataset.** LIDC-IDRI [1] is a dataset for pulmonary nodule classification or detection based on low-dose CT, which involves 1,010 patients. According to the annotations, we extracted 2,026 nodules, and all of them were labeled with scores from 1 to 5, indicating the malignancy progression. We cropped all the nodules with a square shape of a doubled equivalent diameter at the annotated center, then resized them to the volume of  $32 \times 32 \times 32$ . Following [9, 11], we modified the first layer of the image encoder to be with 32 channels. According to existing works [11, 18], we regard a nodule with an average score between 2.5 and 3.5 as *unsure* nodules, *benign* and *malignant* categories are those with scores lower than 2.5 and larger than 3.5, respectively. In this paper, we construct three sub-datasets: LIDC-A contains three classes of nodules both in training and test sets; according to [11], we construct the LIDC-B, which contains three classes of nodules *only* in the training set, and the test set contains benign and malignant nodules; LIDC-C includes benign and malignant nodules both in training and test sets.

**Experimental Settings.** In this paper, we apply the CLIP pre-trained ViT-B/16 as the text encoder for CLIP-Lung, and the image encoder we used is ResNet-18 [6] due to the relatively smaller scale of training data. The image encoder is initialized randomly. Note that for the text branch, we froze the parameters of the text encoder and updated the learnable tokens  $\mathbf{l}$  and  $\mathbf{l}'$  during training. The learning rate is 0.001 following the cosine decay, while the optimizer is stochastic gradient descent with momentum 0.9 and weight decay 0.00005. The temperature  $\tau$  is initialized as 0.07 and updated during training. All of our experiments are implemented with PyTorch [15] and trained with NVIDIA A100 GPUs. The experimental results are reported with average values through five-fold cross-validation. We report the recall and F1-score values for different classes and use “ $\pm$ ” to indicate standard deviation.

#### 3.2 Experimental Results and Analysis

**Performance Comparisons.** In Table 1, we compare the classification performances on the LIDC-A dataset, where we regard the benign-unsure-malignant

**Table 1.** Classification results on the test set of LIDC-A.

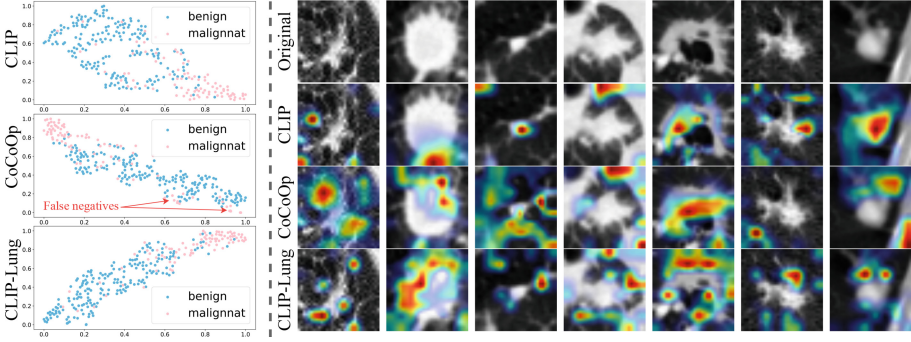
Method	Accuracy	Benign		Malignant		Unsure	
		Recall	F1	Recall	F1	Recall	F1
CE Loss	$54.2 \pm 0.6$	72.2	62.0	64.4	61.3	29.0	36.6
Poisson [2]	$52.7 \pm 0.7$	60.5	56.8	58.4	58.7	41.0	44.1
NSB [13]	$53.4 \pm 0.7$	<b>80.7</b>	63.0	<b>67.3</b>	63.8	16.0	24.2
UDM [18]	$54.6 \pm 0.4$	76.7	64.3	49.5	53.5	32.5	39.5
CORF [21]	$56.8 \pm 0.4$	71.3	63.3	61.3	62.3	38.5	44.3
CLIP [16]	$56.6 \pm 0.3$	59.5	59.2	55.2	60.0	53.9	52.2
CoCoOp [20]	$56.8 \pm 0.6$	59.0	59.2	55.2	60.0	<b>55.1</b>	52.8
<b>CLIP-Lung</b>	<b><math>60.9 \pm 0.4</math></b>	67.5	<b>64.4</b>	60.9	<b>66.3</b>	53.4	<b>54.1</b>

**Table 2.** Classification results on test sets of LIDC-B and LIDC-C.

Method	LIDC-B					LIDC-C				
	Accuracy	Benign		Malignant		Accuracy	Benign		Malignant	
		Recall	F1	Recall	F1		Recall	F1	Recall	F1
CE Loss	$83.3 \pm 0.6$	92.4	88.4	63.4	70.3	$85.5 \pm 0.5$	91.5	89.7	72.3	75.6
Poisson [2]	$81.8 \pm 0.4$	94.2	87.7	54.5	65.1	$84.0 \pm 0.3$	87.9	88.3	75.2	74.5
NSB [13]	$78.1 \pm 0.5$	90.6	85.8	50.5	60.7	$84.9 \pm 0.7$	91.0	89.2	71.3	74.6
UDM [18]	$79.3 \pm 0.4$	87.0	86.2	62.4	67.7	$84.6 \pm 0.5$	88.8	88.8	75.2	75.2
CORF [21]	$81.5 \pm 0.3$	<b>95.9</b>	87.8	49.5	62.8	$83.0 \pm 0.2$	87.9	87.7	72.3	72.6
CLIP [16]	$83.6 \pm 0.6$	92.0	88.7	64.4	70.4	$87.5 \pm 0.3$	92.0	91.0	77.0	78.8
CoCoOp [20]	$86.8 \pm 0.7$	94.5	90.9	69.0	75.9	$88.2 \pm 0.6$	<b>95.0</b>	91.8	72.4	78.8
<b>CLIP-Lung</b>	<b><math>87.5 \pm 0.3</math></b>	94.5	<b>91.7</b>	<b>72.3</b>	<b>79.0</b>	<b><math>89.5 \pm 0.4</math></b>	94.0	<b>92.8</b>	<b>80.5</b>	<b>82.8</b>

as an ordinal relationship. Compared with ordinal classification methods such as Poisson, NSB, UDM, and CORF, CLIP-Lung achieves the highest accuracy and F1-scores for the three classes, demonstrating the effectiveness of textual knowledge-guided learning. CLIP and CoCoOp also outperform ordinal classification methods and show the superiority of large-scale pre-trained text encoders. Furthermore, CLIP-Lung obtained higher recalls than CLIP and CoCoOp *w.r.t.* benign and malignant classes, however, the recall of unsure is lower than theirs. We argue that this is due to the indistinguishable textual annotations, such as similar attributes of different nodules. In addition, we verify the effect of textual branch of CLIP-Lung using MV-DAR [12] on LIDC-A dataset. The obtained accuracy values with and without the textual branch are 58.9% and 57.3%, respectively, demonstrating the effectiveness of integrating textual knowledge.

Table 2 presents a performance comparison of CLIP-Lung on the LIDC-B and LIDC-C datasets. Notably, CLIP-Lung obtains higher evaluation values other than recalls of benign class. This disparity is likely attributed to the similarity in appearances and subtle variations in text attributes among the benign nodules. Consequently, aligning these distinct feature types becomes challenging, resulting in a bias towards the text features associated with malignant nodules.



**Fig. 3.** The t-SNE (Left) and Grad-CAM (Right) results.

**Table 3.** Ablation study on different losses. We report classification accuracies.

$\mathcal{L}_{IC}$	$\mathcal{L}_{IA}$	$\mathcal{L}_{CA}$	LIDC-A	LIDC-B	LIDC-C
✓			$56.8 \pm 0.6$	$86.8 \pm 0.7$	$88.2 \pm 0.6$
✓	✓		$59.4 \pm 0.4$	$86.8 \pm 0.6$	$86.7 \pm 0.4$
	✓	✓	$58.1 \pm 0.2$	$85.7 \pm 0.6$	$87.5 \pm 0.5$
✓		✓	$56.9 \pm 0.3$	$84.7 \pm 0.4$	$84.0 \pm 0.7$
✓	✓	✓	<b><math>60.9 \pm 0.4</math></b>	<b><math>87.5 \pm 0.5</math></b>	<b><math>89.5 \pm 0.4</math></b>

**Visual Features and Attention Maps.** To illustrate the influence of incorporating class and attribute knowledge, we provide the t-SNE [14] and Grad-CAM [17] results obtained by CLIP, CoCoOp, and CLIP-Lung. In Fig. 3, we can see that CLIP yields a non-compact latent space for two kinds of nodules. CoCoOp and CLIP-Lung alleviate this phenomenon, which demonstrates that the learnable prompts guided by nodule classes are more effective than fixed prompt engineering. Unlike CLIP-Lung, CoCoOp does not incorporate attribute information in prompt learning, leading to increased false negatives in the latent space. From the attention maps, we can observe that CLIP cannot precisely capture spiculation and lobulation regions that are highly correlated with malignancy. Simultaneously, our CLIP-Lung performs better than CoCoOp, which demonstrates the guidance from textual descriptions such as “spiculation”.

**Ablation Studies.** In Table 3, we verify the effectiveness of different loss components on the three constructed datasets. Based on  $\mathcal{L}_{IC}$ ,  $\mathcal{L}_{IA}$  and  $\mathcal{L}_{CA}$  improve the performances on LIDC-A, indicating the effectiveness of capturing fine-grained features of ordinal ranks using class and attribute texts. However, they perform relatively worse on LIDC-B and LIDC-C, especially the  $\mathcal{L}_{IC} + \mathcal{L}_{CA}$ . That is to say,  $\mathcal{L}_{IA}$  is more important in latent space rectification, *i.e.*, image-attribute consistency. In addition, we observe that  $\mathcal{L}_{IC} + \mathcal{L}_{IA}$  performs better than  $\mathcal{L}_{IA} + \mathcal{L}_{CA}$ , which is attributed to that  $\mathcal{L}_{CA}$  regularizes the image features indirectly.



## 4 Conclusion

In this paper, we proposed a textual knowledge-guided framework for pulmonary nodule classification, named CLIP-Lung. We explored the utilization of clinical textual annotations based on large-scale pre-trained text encoders. CLIP-Lung aligned the different modalities of features generated from nodule classes, attributes, and images through contrastive learning. Most importantly, CLIP-Lung establishes correlations between learnable prompt tokens and feature maps using the proposed CCP module, and this guarantees explainable attention maps localizing fine-grained clinical features. Finally, CLIP-Lung outperforms compared methods, including CLIP on LIDC-IDRI benchmark. Future work will focus on extending CLIP-Lung with more diverse textual knowledge.

**Acknowledgements.** This work was supported in part by Natural Science Foundation of Shanghai (No. 21ZR1403600), National Natural Science Foundation of China (Nos. 62101136 and 62176059), China Postdoctoral Science Foundation (No. 2022TQ0069), Shanghai Municipal of Science and Technology Project (No. 20JC1419500), and Shanghai Center for Brain Science and Brain-inspired Technology.

## References

1. Armato, S.G., III., McLennan, G., Bidaut, L., et al.: The lung image database consortium (LIDC) and image database resource initiative (IDRI): a completed reference database of lung nodules on CT scans. *Med. Phys.* **38**(2), 915–931 (2011)
2. Beckham, C., Pal, C.: Unimodal probability distributions for deep ordinal classification. In: *Proceedings of the International Conference on Machine Learning*, pp. 411–419 (2017)
3. Chen, X., He, K.: Exploring simple Siamese representation learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15750–15758 (2021)
4. Hancock, M.C., Magnan, J.F.: Predictive capabilities of statistical learning methods for lung nodule malignancy classification using diagnostic image features: an investigation using the lung image database consortium dataset. In: *Medical Imaging 2017: Computer-Aided Diagnosis*, vol. 10134, pp. 558–569 (2017)
5. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738 (2020)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778 (2016)
7. Henaff, O.: Data-efficient image recognition with contrastive predictive coding. In: *International Conference on Machine Learning*, pp. 4182–4192 (2020)
8. Joshi, A., Sivaswamy, J., Joshi, G.D.: Lung nodule malignancy classification with weakly supervised explanation generation. *J. Med. Imaging* **8**(4), 044502–044502 (2021)
9. Lei, Y., Tian, Y., Shan, H., Zhang, J., Wang, G., Kalra, M.K.: Shape and margin-aware lung nodule classification in low-dose CT images via soft activation mapping. *Med. Image Anal.* **60**, 101628 (2020)

10. Lei, Y., Zhang, J., Shan, H.: Strided self-supervised low-dose CT denoising for lung nodule classification. *Phenomics* **1**, 257–268 (2021)
11. Lei, Y., Zhu, H., Zhang, J., Shan, H.: Meta ordinal regression forest for medical image classification with ordinal labels. *IEEE/CAA J. Autom. Sin.* **9**(7), 1233–1247 (2022)
12. Liao, Z., Xie, Y., Hu, S., Xia, Y.: Learning from ambiguous labels for lung nodule malignancy prediction. *IEEE Trans. Med. Imaging* **41**(7), 1874–1884 (2022)
13. Liu, X., Zou, Y., Song, Y., Yang, C., You, J., K Vijaya Kumar, B.: Ordinal regression with neuron stick-breaking for medical diagnosis. In: *Proceedings of the European Conference on Computer Vision*, pp. 335–344 (2018)
14. Van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**(11) (2008)
15. Paszke, A., et al.: Automatic differentiation in PyTorch (2017)
16. Radford, A., et al.: Learning transferable visual models from natural language supervision. In: *International Conference on Machine Learning*, pp. 8748–8763 (2021)
17. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-CAM: visual explanations from deep networks via gradient-based localization. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 618–626 (2017)
18. Wu, B., Sun, X., Hu, L., Wang, Y.: Learning with unsure data for medical image diagnosis. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 10590–10599 (2019)
19. Xie, Y., Zhang, J., Xia, Y.: Semi-supervised adversarial model for benign-malignant lung nodule classification on chest CT. *Med. Image Anal.* **57**, 237–248 (2019)
20. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for vision-language models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16816–16825 (2022)
21. Zhu, H., et al.: Convolutional ordinal regression forest for image ordinal estimation. *IEEE Trans. Neural Netw. Learn. Syst.* **33**(8), 4084–4095 (2022)