

Joint learning of 3D lesion segmentation and classification for explainable COVID-19 diagnosis

Xiaofei Wang, Lai Jiang, Liu Li, Mai Xu, Xin Deng, Lisong Dai, Xiangyang Xu, Tianyi Li, Yichen Guo, Zulin Wang, and Pier Luigi Dragotti

Abstract—Given the outbreak of COVID-19 pandemic and the shortage of medical resource, extensive deep learning models have been proposed for automatic COVID-19 diagnosis, based on 3D computed tomography (CT) scans. However, the existing models independently process the 3D lesion segmentation and disease classification, ignoring the inherent correlation between these two tasks. In this paper, we propose a joint deep learning model of 3D lesion segmentation and classification for diagnosing COVID-19, called DeepSC-COVID, as the first attempt in this direction. Specifically, we establish a large-scale CT database containing 1,805 3D CT scans with fine-grained lesion annotations, and reveal 4 findings about lesion difference between COVID-19 and community acquired pneumonia (CAP). Inspired by our findings, DeepSC-COVID is designed with 3 subnets: a cross-task feature subnet for feature extraction, a 3D lesion subnet for lesion segmentation, and a classification subnet for disease diagnosis. Besides, the task-aware loss is proposed for learning the task interaction across the 3D lesion and classification subnets. Different from all existing models for COVID-19 diagnosis, our model is interpretable with fine-grained 3D lesion distribution. Finally, extensive experimental results show that the joint learning framework in our model significantly improves the performance of 3D lesion segmentation and disease classification in both efficiency and efficacy.

Index Terms—Multi-task Learning, CT scans, COVID-19, Deep Neural Networks.

I. INTRODUCTION

After being first identified in December 2019, COVID-19 has emerged as a pandemic of global health concern, causing unprecedented social and economic disruption [?], [?]. According to the WHO report [?], as of March 29, 2021, there were a total of 126,890,643 infected patients, 2,778,619 of whom died. The worldwide outbreak of COVID-19 has placed enormous pressure on healthcare systems and led to an extreme shortage of medical resources [?]. A feasible way to control the COVID-19 pandemic is to identify and isolate the infected cases [?], which requires an effective screening method with high sensitivity to detect infected people and their close contacts. Real-time reverse transcription polymerase chain reaction (RT-PCR) [?] is

X. Wang, L. Jiang, M. Xu, X. Deng, T. Li, Y. Guo and Z. Wang are with the Beihang University, Beijing 100191 China; L. Li and P. L. Dragotti are with the Imperial College London, London SW7 2AZ UK; L. Dai and X. Xu are with the Liyuan Hospital, Huazhong University of Science and Technology, Wuhan 430077 China. X. Wang, L. Jiang, and L. Li contributes equally to this work. M. Xu, X. Deng and X. Xu are the corresponding authors of this paper (E-mail: Maixu@buaa.edu.cn; cindydeng@buaa.edu.cn; 1993ly0538@hust.edu.cn). This work was supported by the NSFC project under grants 61922009, 61876013 and 62050175, Fundamental Research Funds for the Central Universities under grant 2020kfyXGYJ097 and Beijing Natural Science Foundation under Grant JQ20020.

a common method for COVID-19 detection; however, it suffers from a high false negative rate [?], [?] in the early stages of the disease. Recently, the antigen test has been developed for the rapid diagnosis of COVID-19; however, it still suffers from relatively low specificity and sensitivity. As reported in [?], the sensitivity of the antigen test is only 30.2%. Chest computed tomography (CT) has been demonstrated to have better sensitivity for detecting COVID-19, especially in regions with severe epidemic situations [?], [?]. Unfortunately, it is a time-consuming process for doctors to interpret and make a diagnosis during COVID-19 outbreak from each CT scan with hundreds of slices. Even an experienced radiologist can only interpret 4–10 chest CT scans per hour [?], [?]. Therefore, an automatic CT interpretation model is highly desired for accurate, efficient and trustworthy COVID-19 diagnosis.

There are three great challenges in developing an automatic CT interpretation model for COVID-19 diagnosis. (1) Although the tasks of 3D lesion segmentation and disease classification are highly correlated with each other for COVID-19 diagnosis, they cannot be simultaneously learned in the existing deep learning (DL) models [?], [?], [?], [?], [?], [?]. Hence, it is challenging to develop a joint deep learning model of 3D lesion segmentation and disease classification. (2) Despite being a new disease, COVID-19 has similar imaging manifestations as other types of pneumonia, e.g., community acquired pneumonia (CAP) [?]. Thus, it is a challenging task for the model to produce a differential diagnosis between COVID-19 and other similar types of pneumonia. (3) Most automatic diagnosis models [?], [?], [?], [?] are based on “black box” deep neural networks (DNNs) [?], [?], [?], [?], which lack sufficient explainability to assist radiologists in making credible diagnoses. The explainability of DNNs is another challenge in the design of automatic CT interpretation models for COVID-19 diagnosis.

To tackle the above challenges, we establish a large-scale CT database, called 3DLSC-COVID, which is the first CT database with fine-grained 3D lesion segmentation and classification labels of COVID-19, CAP and non-pneumonia. Based on the lesion characteristics found from this database, we propose a joint DL model, namely DeepSC-COVID, for accurate 3D lesion segmentation and the diagnosis of COVID-19. Specifically, the DeepSC-COVID model consists of three subnets, i.e., cross-task feature, 3D lesion segmentation and disease classification subnets, and is able to simultaneously generate the 3D segmented lesion and the classification results of COVID-19, CAP or non-pneumonia. In the classification subnet, a new multi-layer visualization mechanism is developed to generate the evidence masks that contain small and indistinct lesions

for disease diagnosis. In this way, the process of COVID-19 diagnosis in our model is explainable. Besides, in the training phase, a novel task-aware loss is proposed on the basis of our visualization mechanism for efficient interaction between the tasks of segmentation and classification. With the guidance of the segmented lesions, the classification subnet is able to focus on the lesions, such that the diagnosis of COVID-19 can be significantly accelerated with higher classification accuracy. Note that, different from the single-scale attention constrained mechanism [?], our task-aware loss has multi-scale attention constraint to generate more fine-grained visualization maps. Additionally, the task-aware loss is used in our method to optimize both tasks of segmentation and classification, thus being able to interact the information between these two tasks and to boost the performance of both tasks. In conclusion, the developed DeepSC-COVID model can provide the rapid, accurate and explainable diagnosis of COVID-19, meanwhile visualizing the fine-grained lesions for doctors.

To the best of our knowledge, our method is one of the pioneering works in joint learning of 3D lesion segmentation and disease classification based on 3D CT scans, especially for the disease of COVID-19. The main contributions of this paper are as follows. (1) We establish a large-scale database of CT scans, with fine-grained lesion annotations, for the diagnosis of COVID-19 and CAP. (2) We thoroughly analyze the new database, and yield 4 important findings about the lesion differences between the diseases. (3) We propose an explainable deep multi-task learning model for both tasks of 3D lesion segmentation and disease classification of COVID-19.

II. RELATED WORK

Imaging-based COVID-19 databases. Although many people infected by COVID-19, it is still not easy to build a large-scale imaging-based COVID-19 databases, due to the privacy of the patients and hospitals. Table I summarizes the representative CT/X-rays based COVID-19 databases that are public online. As can be seen from this table, most existing public databases lack fine-grained lesion annotation, and only a few of them have small scale of lesion segmentation labels. This is probably because of the lack of the experts with rich experience in diagnosing COVID-19. In contrast, this paper establishes a large-scale database of 1,805 CT scans with 458,730 slices, in which 157,696 slices are annotated with lesions. It is worth mentioning that the lesion-annotated slices of our database are around 17 times more than those of the largest 3D lesion segmentation database [?].

Automatic COVID-19 diagnosis on CT scans. In the past few months, many DL-based methods were developed for COVID-19 diagnosis on CT scans [?], [?], [?], [?], [?], [?]. They mainly focus on two tasks: disease classification and lesion segmentation. In order to automatically diagnose COVID-19, Li *et al.* [?] proposed a COVID-19 detection neural network (COVNet) using ResNet-50 [?] as the backbone. With a series of CT slices as inputs, COVNet generates a classification result for each CT scan. Similarly, ouyang [?] *et al.* designed a dual-sampling attention network for classifying COVID-19 and CAP. Specifically, they proposed an online attention module with a 3D convolutional network to focus on the infection

TABLE I
SUMMARY OF THE EXISTING COVID-19 DATABASES.

Database	Type	# Slices	#Cases	Lesion Annotation*
[?]	X-rays	761	412	-
[?]	X-rays	2,905	-	-
[?]	X-rays	16,756	13,645	-
[?]	X-rays+CT	5,381	1,311	-
[?]	CT	349	216	-
[?]	CT	1,103	64(videos)	-
[?]	CT	1,110	-	-
[?]	CT	2,482	120	-
[?]	CT	453	99	-
[?]	CT	361,221	2,246	2D (4,695 slices)
[?]	CT	144,167	750	2D (3,855 slices)
[?]	CT	3,520	20	3D (1,844 slices)
[?]	CT	76,250	558	3D (9,015 slices)
Ours	CT	458,730	1,805	3D (157,696 slices)

* “2D” means that only part of the slices of one CT scan are annotated, while “3D” denotes that all the slices with lesions of a CT scan are annotated.

regions in lungs for the diagnosis. Different from the disease classification methods, other works [?], [?] focused on COVID-19 lesion segmentation. Specifically, Zhou *et al.* [?] proposed a fully automatic machine-agnostic method that can segment and quantify the infection regions on CT scans from different sources. Wang *et al.* [?] designed a noise-robust framework for automatic segmentation of COVID-19 pneumonia lesions from CT scans. Unfortunately, all above methods neglect the correlation between disease classification and lesion segmentation. In fact, the lesion segmentation results act as explainable diagnostic evidence for disease classification; meanwhile, the classification results are able to further improve the accuracy of lesion segmentation.

Only a few DL-based methods [?], [?], [?] have been developed to perform both tasks of lesion segmentation and disease classification for COVID-19. Specifically, Mahmud *et al.* [?] proposed a hybrid attention based network for lesion segmentation, diagnosis, and severity prediction of COVID-19. In their training stage, the lesion segmentation network is optimized firstly and is then integrated into the training of diagnosis and severity prediction. Similarly, Jin *et al.* [?] proposed a sequential optimization pipeline, in which they first train the lesion segmentation network alone, and then use the segmentation results to train the classification network. However, all these methods cannot be seen as multi-task learning according to the definition of [?], since they separately learn the two tasks, ignoring the information sharing between two tasks. In contrast, our DeepSC-COVID method is a multi-task deep learning work, as it can jointly learn the two tasks of 3D lesion segmentation and classification for COVID-19, achieving task-aware information sharing through the proposed cross-task feature subnet and the novel task-aware loss. This way, the tasks of lesion segmentation and disease classification can boost each other to achieve better performance.

III. DATABASE AND ANALYSIS

A. Database Establishment

This retrospective study was performed in accordance with the Declaration of Helsinki of the World Medical Association and was approved by the medical ethics committee of Liyuan Hospital, Tongji Medical College, Huazhong University of Science and Technology. Besides, all data were anonymized.

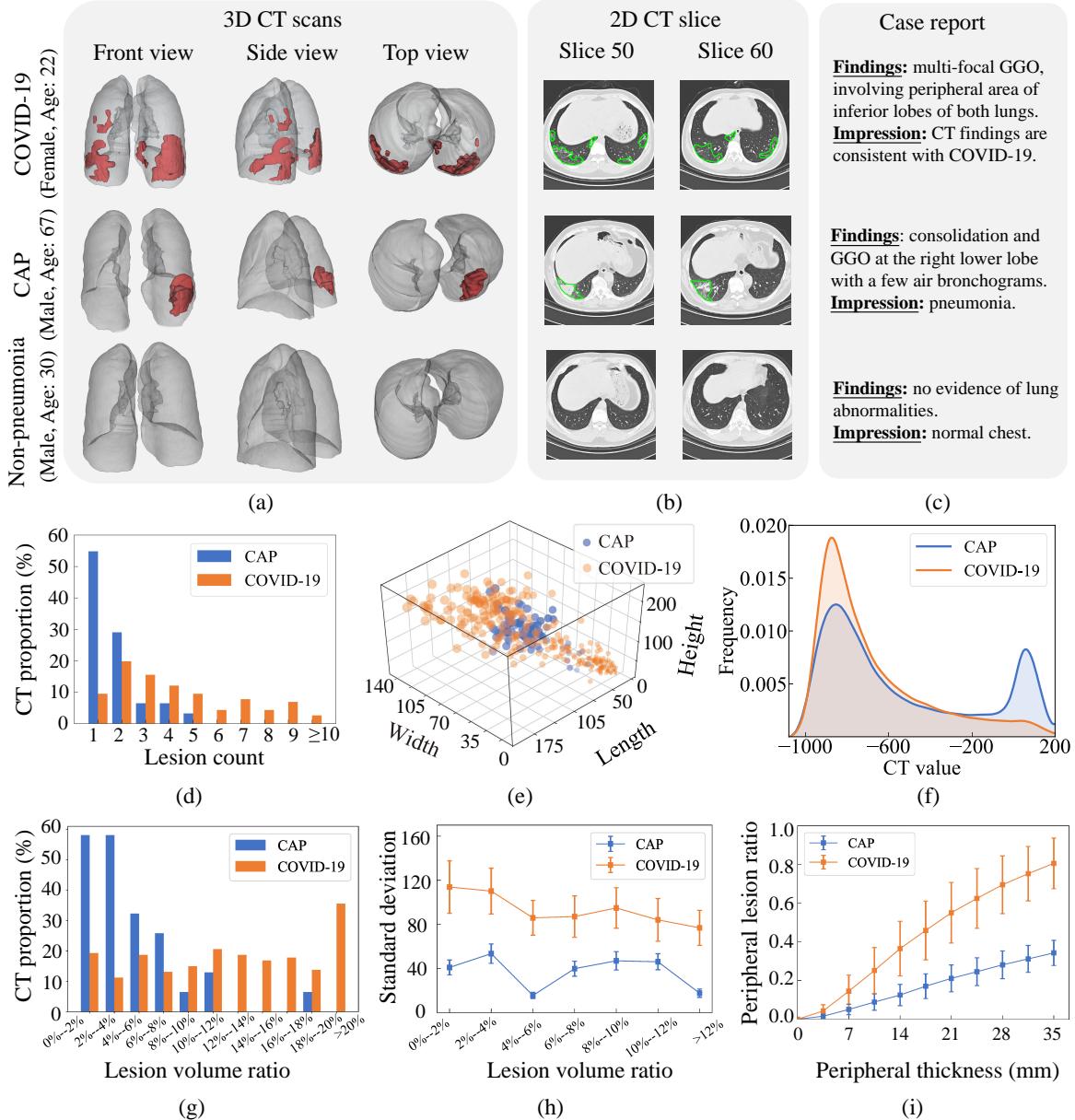


Fig. 1. Illustration and statistical analysis of our database. (a) Front, side and top views of 3D chest CT scans with 3D lesion segmentation for one COVID-19 (upper row), one CAP (middle row) and one non-pneumonia (lower row) individuals. Note that the lesions in lungs are marked in red. (b), Two selected 2D CT slices, corresponding to 3D CT scan in the same row. The lesions in each slice are encircled by green lines. (c) Case reports of the three individuals. (d) Histogram of lesion counts in the CT scans for all COVID-19 and CAP individuals in our 3DLSC-COVID database. (e) Width, length and height of each lesion in 3D CT scans for COVID-19 and CAP, respectively. For better visualization, only 372 lesions are randomly selected from our database. (f) Distribution curves of CT values in the lesions of CAP and COVID-19 CT scans, respectively. (g) Histograms of lesion count in the CT scans for CAP and COVID-19, respectively. (h) Standard deviations of lesion distribution for CAP and COVID-19 with varied lesion volume ratios. (i) Changes of peripheral lesion ratio with peripheral thickness varying from 0 to 35 mm for CAP and COVID-19, respectively. Note that the results of these charts (d, f-i) are obtained upon all CAP and COVID-19 CT scans in our 3DLSC-COVID database.

For establishing our 3DLSC-COVID database ¹, a total of 1,805 3D chest CT scans with more than 570,000 CT slices were collected from 2 standard CT scanners of Liyuan Hospital, i.e., UIH uCT 510 and GE Optima CT600. Among all CT scans, there were 794 positive cases of COVID-19, which were further confirmed by clinical symptoms and RT-PCR from January 16 to April 16, 2020. Additionally, 540 positive cases of CAP and 471 non-pneumonia cases were randomly selected from the same

hospital between November 5, 2016 and April 28, 2020. In contrast to existing CT-based COVID-19 databases [?], [?], [?], our 3DLSC-COVID database is the first CT database with both fine-grained 3D lesion segmentation and disease classification labels for the COVID-19 and CAP diagnosis. More details about patients and CT scans of the 3DLSC-COVID database are summarized in the supplementary material.

For lesion segmentation, we recruited 2 resident radiologists with over 2 years of experience to annotate the areas and boundaries of the lesions in each 2D CT slice. Then, for

¹The 3DLSC-COVID database is available at IEEE Dataport <https://dx.doi.org/10.21227/mxb3-7j48>.

each CT scan, the 2D annotated lesions of all CT slices were merged to obtain the 3D lesions. Subsequently, the 2 resident radiologists were asked to further refine the segmented lesions in 3D viewing mode. At last, both 2D and 3D lesions were reviewed and corrected by a senior radiologist with over 10 years of experience in thoracic radiology. Some examples of the annotated CT scans for COVID-19, CAP and non-pneumonia individuals are shown in Fig. 1.

B. Analysis of characteristics of 3D lesion.

We characterize the 3D lesions of COVID-19 and CAP via thoroughly analyzing the lesion annotations in our 3DLSC-COVID database. Four important findings are obtained in terms of the count, size, CT value and spatial distribution of 3D lesions, which are briefly introduced as follows.

Finding 1: The number of lesions for COVID-19 is considerably more than that for CAP.

Analysis: As shown in Fig. 1 (d), the number of lesions for COVID-19 is around 2.6 times than that for CAP, i.e., averagely 4.4 lesions per CT scan for COVID-19 versus 1.7 for CAP. To be more specific, 54.8% CT scans of CAP in our database contain only one lesion, while around 55.2% CT scans for COVID-19 have 4 or more lesions. Fig. 1 (a) visualizes the segmented lesions of COVID-19 and CAP, which also indicate the obvious difference of lesion counts between COVID-19 and CAP.

Finding 2: The overall lesion volume in COVID-19 CT scans is significantly larger than that for CAP. Additionally, compared with CAP, the lesions of COVID-19 vary significantly in size.

Analysis: Fig. 1 (g) shows the histogram of the lesion volume ratio (LVR) for all COVID-19 and CAP individuals in the 3DLSC-COVID database. Here, LVR indicates the proportion of lesions to the whole lung. As shown in this figure, the average LVR for COVID-19 is around 3.3 times higher than that for CAP per CT scan, i.e., the average LVR is 14.3% for COVID-19 versus 4.4% for CAP. The CAP cases with LVR larger than 12% accounts for only 5%, while the COVID-19 cases with LVR larger than 12% accounts for above 50%.

In addition to the lesion volume, we compare the 3D size of lesions for COVID-19 and CAP. The 3D size is measured by drawing a bounding box for each lesion, which is defined as the minimum cuboid to wrap the lesion. Fig. 1 (e) shows the 3-D scatter diagram with axes of width, length and height, drawn on 372 randomly selected lesions from our database. As can be seen in this figure, the 3D size of lesions for CAP is concentrated. Specifically, the width, length and height of over 90 % CAP lesions are densely distributed in the range of [35 mm, 105 mm] (span = 70 mm), [42 mm, 105 mm] (span = 63 mm) and [85 mm, 160 mm] (span = 75 mm), respectively. In contrast, the 3D lesion size of COVID-19 is with a larger span, i.e., [15 mm, 140 mm] (span = 125 mm) in width, [15 mm, 170 mm] (span = 155 mm) in length and [30 mm, 240 mm] (span = 210 mm) in height, respectively. This verifies that the lesions of COVID-19 vary significantly in size compared to those of CAP.

Finding 3: Compared to the CAP lesions which can either display low or high density in CT images, the COVID-19 lesions tend to mainly display low density (darker).

Analysis: The densities of CAP and COVID-19 lesions are investigated in terms of CT values. Fig. 1 (f) shows the distribution curves of CT values in the lesions of CAP and COVID-19, respectively. Note that smaller CT values indicate lower density. As can be seen, for COVID-19 lesions, the distribution curve only has one peak, with more than 70% of the CT values concentrated between -960 Hounsfield unit (HU) and -600 HU. In contrast, for CAP lesions, the CT value distribution has two primary peaks, i.e., over 75% of the CT values are distributed in [-970 HU, -580 HU] and [-70 HU, 140 HU]. As such, the COVID-19 lesions tend to mainly display low density, while the CAP lesions can either display low or high density. A possible medical explanation for this finding lies in the lesion types. Specifically, the COVID-19 lesions are mainly ground-glass opacity (GGO) [?], [?], which is a pattern of hazy increased lung opacity that shows low contrast with surrounding regions. In addition to GGO, CAP has another type of lesion called consolidation. The consolidation is a typical pneumonia lesion that has the homogeneous increase in lung parenchymal attenuation of CT scans, which is in highly contrast with surrounding regions. Some examples of the segmented lesions in 2D CT slices for COVID-19 and CAP can be seen in Fig. 1 (b).

Finding 4: The COVID-19 lesions are mostly scattered in the peripheral area of lungs. In contrast, the CAP lesions are more concentrated, which are mainly distributed in the central area of lungs.

Analysis: The spatial distribution of the lesions is evaluated for CAP and COVID-19, by measuring the standard deviation of the lesion centers in all CT scans. Here, the lesion center is the central point of the lesion bounding box, and a larger standard deviation indicates more scattered lesion distribution. For specific analysis, we divide all CT scans into different groups upon their lesion volume ratios. Fig. 1 (h) shows the standard deviations of lesion distribution in different groups of CAP and COVID-19. As can be seen in this figure, the standard deviation of lesions in COVID-19 is significantly larger than that in CAP for all CT groups. In particular, for the CT group with lesion volume ratio from 4% to 6%, the standard deviation is 86.0 on average for COVID-19, compared with only 15.6 for CAP. This demonstrates that the spatial distribution of COVID-19 lesions is more scattered than that of CAP.

Next, the lesion distribution areas in CT scans are analyzed for CAP and COVID-19, by calculating the proportion of lesions within the peripheral lung areas to the overall lesions, denoted as peripheral lesion ratio (PLR). To calculate PLR, given a CT scan, we first generate the 3D binary masks of the lung areas by a state-of-the-art lung segmentation algorithm [?]. For the CT scan with slice S , width W and height H , the lung and lesion masks are denoted as $\mathbf{U} \in \mathbb{R}^{S \times W \times H}$ and $\mathbf{L} \in \mathbb{R}^{S \times W \times H}$, respectively. Then, PLR is defined as follows:

$$\text{PLR} = \frac{\sum_{s=1}^S \sum_{i=1}^W \sum_{j=1}^H [\mathbf{U}_s - E(\mathbf{U}_s, \sigma)]_{i,j} \mathbf{L}_{s,i,j}}{\sum_{s=1}^S \sum_{i=1}^W \sum_{j=1}^H \mathbf{L}_{s,i,j}}, \quad (1)$$

where \mathbf{U}_s is the s -th slice of the lung mask \mathbf{U} , and $E(\mathbf{U}_s, \sigma)$ is the erosion operation with the erosion kernel of σ in radius.

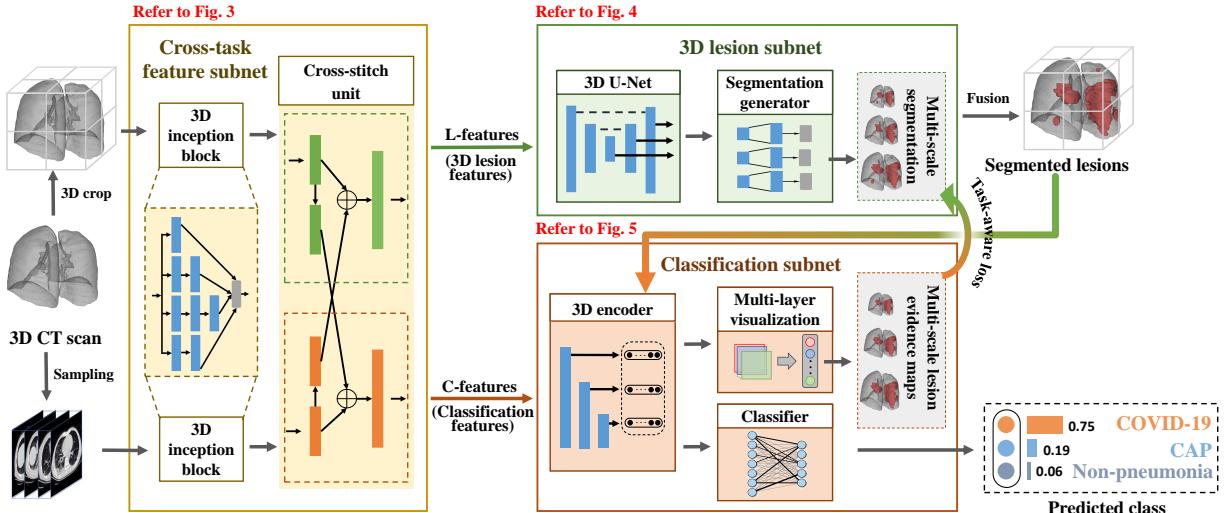


Fig. 2. Framework of the proposed DeepSC-COVID model.

Note that the difference between the lung mask and its erosion result [$\mathbf{U}_s - E(\mathbf{U}_s, \sigma)$] can be regarded as the peripheral lung areas, which is controlled by the hyper-parameter of σ denoted as peripheral thickness in the following. Fig. 1 (i) shows the PLR with different peripheral thickness for the CT scans of COVID-19 and CAP in the 3DLSC-COVID database. As shown, the COVID-19 lesions are more possibly distributed in the peripheral area of the lung, e.g., PLR = 62.4% for COVID-19 lesions *versus* 24.5% for CAP lesions. This indicates the significant difference of lesion distribution between CAP and COVID-19 in CT scans.

The above findings reveal the typical characteristics of lesions for COVID-19, and are used as guidance to design our DeepSC-COVID model for automatic CT interpretation in COVID-19 diagnosis.

IV. METHODOLOGY

A. Framework of DeepSC-COVID

As illustrated in Fig. 2, the proposed DeepSC-COVID model² consists of 3 subnets: cross-task feature, 3D lesion and classification subnets. For 3D lesion segmentation, due to the limited GPU memory, it is difficult to input the full-sized CT scans. As such, the original CT scan is cropped into smaller non-overlapping 3D patches. For classification, the 3D CT scan is preprocessed by slice sampling at an average interval to remove the redundancy between adjacent slices, in order to improve the classification efficiency.

After preprocessing, both the cropped 3D CT patch and sampled 2D CT slices are fed into the cross-task feature subnet with 3D inception blocks and cross-stitch unit. Specifically, based on the classic 2D inception block [?], the 3D inception block is designed to extract the multi-scale 3D features from the cropped 3D CT patch and sampled 2D CT slices, respectively. Then, the cross-stitch unit is developed to mix the features to generate 3D lesion features (L-features) and classification

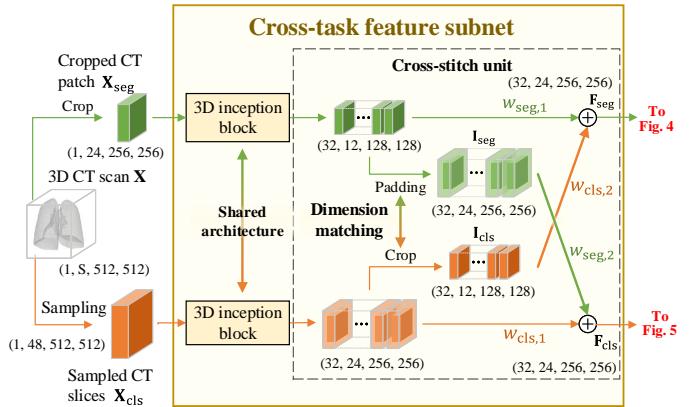


Fig. 3. Structure of the cross-task feature subnet in the proposed DeepSC-COVID model.

features (C-features). These two features are fed into the 3D lesion and classification subnets, respectively. In the 3D lesion subnet, a 3D U-Net and a segmentation generator are designed to segment the multi-scale 3D lesions of COVID-19 or CAP. In the classification subnet, a 3D encoder and a classifier are developed to predict the probability scores for COVID-19, CAP and non-pneumonia. Besides, the task-aware loss is proposed for learning the task interaction across the 3D lesion and classification subnets. To obtain the evidence masks of the classification subnet, we propose a multi-layer visualization method for extracting the pathological regions for disease diagnosis. Finally, according to the predicted probabilities, the input 3D CT scan can be classified as COVID-19, CAP or non-pneumonia.

B. Cross-task feature subnet.

Let \mathbf{X}_{seg} and \mathbf{X}_{cls} denote the cropped 3D patch and the sampled CT slices after preprocessing, the details of which is introduced in Section V-A. Given \mathbf{X}_{seg} and \mathbf{X}_{cls} , the cross-task feature subnet is designed to jointly extract the 3D features for the subsequent 3D lesion segmentation and classification

²The source codes of our DeepSC-COVID model are available at Github (<https://github.com/XiaofeiWang2018/DeepSC-COVID>)

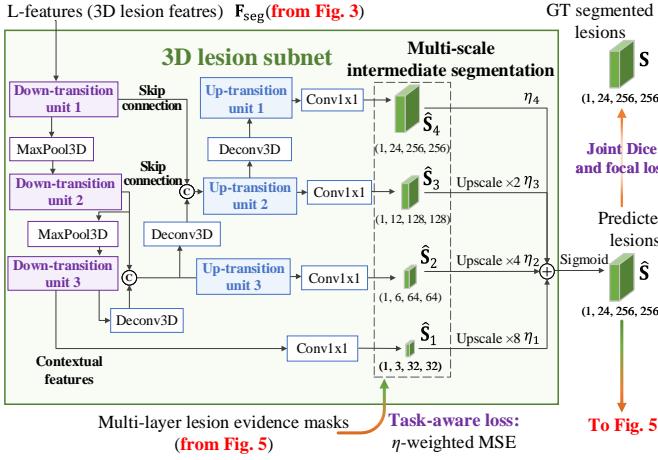


Fig. 4. Structure of the 3D lesion subnet in the proposed DeepSC-COVID model.

subnets. The structure of the cross-task feature subnet is shown in Fig. 3, which consists of 2 cascaded components, i.e., the 3D inception block and the cross-stitch unit. The structure details about these 2 components are described in the following paragraphs.

(1) *3D inception block*. Based on the classic 2D inception block [?], the 3D inception block is developed to extract the multi-scale 3D features. The 3D inception block has 4 branches with cascaded 3D convolutional layers. Benefiting from the multiple receptive fields of different branches, the multi-scale 3D features are extracted, followed by the group normalization [?] and rectified linear unit (ReLU) activation. The specific kernel size, stride and output channel for each 3D convolutional layer are shown in supplementary material. The X_{seg} and X_{cls} are input to two different 3D inception blocks, to extract 3D features for the segmentation and classification tasks, respectively. These two 3D inception blocks do not share parameters, allowing for inception blocks to extract more efficient features for each single task.

(2) *Cross-stitch unit*. Next, the cross-stitch unit is proposed to enhance the information interaction between the segmentation and classification tasks. Specifically, given the extracted 3D features from the 3D inception blocks, dimension matching is first conducted to unify the receptive field of the extracted features via zero-padding and 3D cropping. Let I_{seg} and I_{cls} denote the dimension-matched features for segmentation and classification, respectively. Then, for enhancing the information interaction, I_{seg} and I_{cls} are linearly combined to generate the final cross-task 3D lesion features (L-features) F_{seg} and classification features (C-features) F_{cls} via the following formulation:

$$\begin{bmatrix} F_{seg} \\ F_{cls} \end{bmatrix} = \begin{bmatrix} w_{seg,1} & w_{cls,2} \\ w_{seg,2} & w_{cls,1} \end{bmatrix} \begin{bmatrix} I_{seg} \\ I_{cls} \end{bmatrix}, \quad (2)$$

where $w_{cls,1}$, $w_{cls,2}$, $w_{seg,1}$ and $w_{seg,2}$ are learnable weights. Then, F_{seg} and F_{cls} are fed into the subsequent 3D lesion and classification subnets for further processing.

C. 3D lesion subnet

Given L-features F_{seg} extracted from the cross-task feature subnet, the 3D lesion subnet is developed to segment the 3D

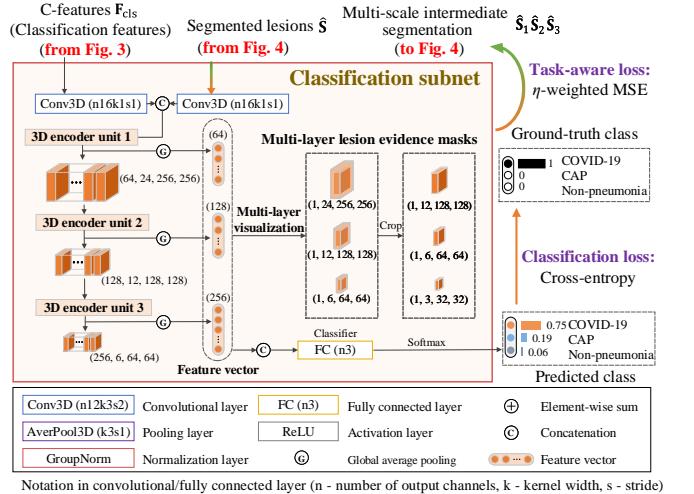


Fig. 5. Structure of the classification subnet in the proposed DeepSC-COVID model.

lesions of CT scans. The structure of the 3D lesion subnet is shown in Fig. 4. In the 3D lesion subnet, a U-shaped 3D structure, which is composed of three down-transition units and three up-transition units, is designed to extract the features for precisely localizing 3D lesions. Specifically, the input L-features F_{seg} are progressively contracted and down-sampled through three down-transition units followed by 3D max pooling layers with stride of 2. In this way, the contextual information of 3D CT scans can be captured in the outputs of the last down-transition units, namely contextual features. Subsequently, the contextual features are progressively expanded and up-sampled through three up-transition units followed by deconvolutional layers [?] with stride of 2. Note that the skip connections are adopted between the up-transition unit and its corresponding down-transition unit, in order to provide boundary information during the up-sampling process. Detailed structures of down-transition and up-transition units can be found in supplementary material.

Then, the outputs of the last down-transition unit and each up-transition unit are further processed by the 3D convolution layers to generate the multi-scale intermediate segmentation. Assuming that \hat{S}_i is the segmentation result at the i -th scale, the final segmentation lesion \hat{S} is calculated as follows:

$$\hat{S} = \text{sigmoid}\left(\sum_{i=1}^4 \eta_i \cdot \text{UP}(\hat{S}_i, 2^{4-i})\right). \quad (3)$$

In the above equation, $\{\eta_i\}_{i=1}^4$ are the hyper-parameters to balance the intermediate segmentation at different scales, and $\text{UP}(\cdot, t)$ is the t -time upscale operation. During the training stage, segmented lesion \hat{S} is supervised by its corresponding ground-truth lesion. Furthermore, the intermediate segmentation result is also supervised by the multi-scale lesion evidence masks from the classification subnet, with minimization on the task-aware loss. The details of the evidence masks and the task-aware loss are introduced in the following sections.

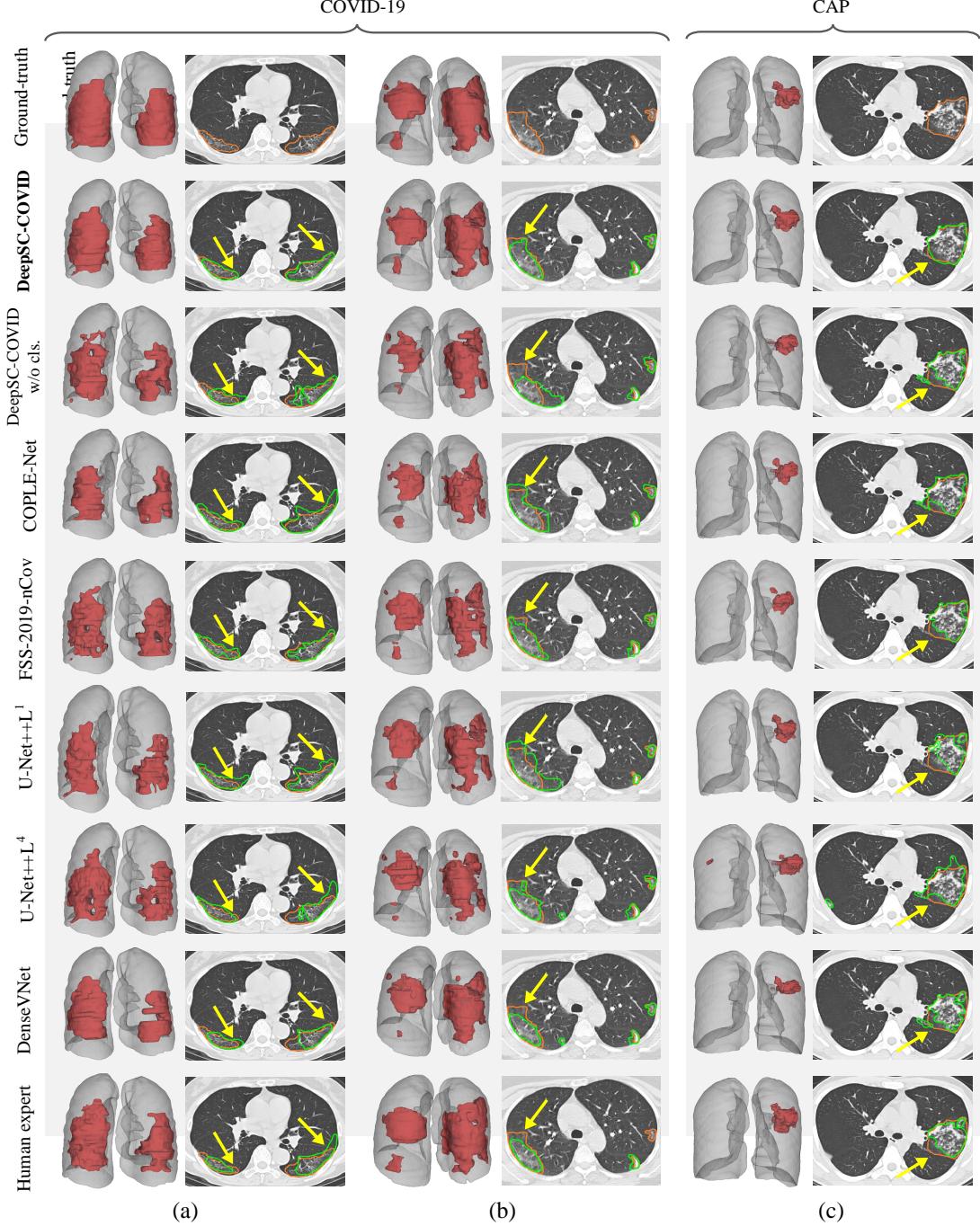


Fig. 6. Visual comparison of 3D and 2D lesion segmentation results. (a-b) Segmentation results of two COVID-19 samples. (c) Segmentation results of CAP sample. For 3D visualizations, the lesions and lungs are shown in red and grey for better view. For 2D visualizations, orange and green curves indicate the ground-truth segmentation results and the results generated by different methods.

B. 3D lesion segmentation results

We qualitatively and quantitatively evaluate the lesion segmentation performance of our DeepSC-COVID model. Table II reports the 3D lesion segmentation results of our DeepSC-COVID and other state-of-the-art segmentation models. As shown in this table, our model achieves high accuracy in 3D lesion segmentation, i.e., 73.3%, 80.2%, 95.6%, 71.8%, and 2.8 mm in terms of Dice similarity coefficient (DSC), sensitivity, specificity, normalized surface Dice (NSD) and root

mean square symmetric surface distance (RMSD), respectively. In contrast to our model, the accuracy of other segmentation models is relatively low, e.g., the DSC scores are only 61.2%, 65.3%, 63.7% and 67.2% for UNet++L¹ [?], UNet++L⁴ [?], DenseVNet [?] and COUPLE-Net [?], respectively. Note that UNet++L¹ and UNet++L⁴ are the lightest and heaviest versions in [?]. Similar results can be found for other metrics, including sensitivity, specificity, RMSD and NSD. Additionally, Fig. 6 visualizes the segmentation results of our and the comparison

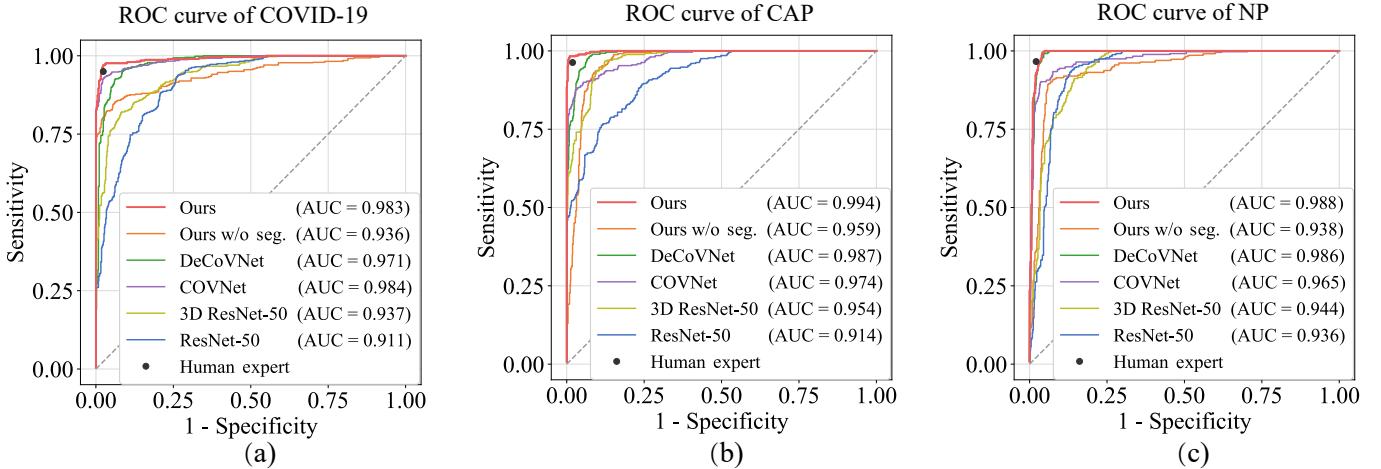


Fig. 7. The ROC curves of our DeepSC-COVID model, other models and human expert in the identification of COVID-19 (a), CAP (b) and NP (c).

models. As shown, our DeepSC-COVID model can locate both the COVID-19 and CAP lesions with higher accuracy than other models. In addition to the segmentation accuracy, we outperform most compared models in terms of segmentation efficiency, i.e., it takes 2.2 seconds for our model to segment a 3D CT scan, while other models require 1.0 to 10.8 seconds to process one 3D CT scan.

To further show the superiority of our DeepSC-COVID model, we compare the segmentation performance between our model and a human expert. Here, the human expert is a radiologist with 5 years of working experience. As shown in Table II, our model significantly outperforms the human expert in 3D lesion segmentation, with an improvement of 14%, 19%, 6.9%, 11.6%, and 10 mm in terms of DSC, sensitivity, specificity, NSD and RMSD, respectively. It is not surprising to see the low dice score of the human expert, since lesion segmentation in medical images is known to suffer from high inter-reader variability [?]. To conclude, our DeepSC-COVID model performs considerably better than the other segmentation models and the human expert for 3D lesion segmentation.

C. Disease classification results

Table II shows the classification results of our DeepSC-COVID model and 4 other state-of-the-art models for classifying COVID-19, CAP and non-pneumonia individuals. As can be seen, the classification accuracy (94.5%) of our model is considerably higher than those of the alternative methods, i.e., ResNet-50 [?] (77.7%) , 3D ResNet-50 [?] (82.5%) , COVNet [?] (87.2%) and DeCoVNet [?] (89.2%). Moreover, the sensitivity, specificity and area under the receiver operating characteristic (ROC) curve (AUC) of our model are the highest among all models. Table II also compares F₁-score between our and other models. Our model has an F₁-score of 94.2%, while ResNet-50, 3D ResNet-50, COVNet and DeCoVNet only yield values of 77.9%, 82.0%, 87.5% and 88.9%, respectively. In addition, Fig. 7 shows the ROC curves for each category, which visualize the tradeoff between sensitivity and specificity. Compared with other four models [?], [?], [?], [?], our ROC curve is closer to the upper-left corner, indicating that our model achieves better

classification results than do the 4 other models. To summarize, our DeepSC-COVID model is considerably better than 4 other models with respect to classifying COVID-19, CAP and non-pneumonia.

Compared to the human expert, the proposed DeepSC-COVID model offers a great advantage in diagnosis speed, i.e., the classification speed of our model is 2.2 seconds, which is significantly faster than the human expert (378.7 seconds). In addition, our model is comparable to the human expert in diagnosis accuracy, i.e., the average sensitivity, specificity and F₁-score of our model are only around 1.0% lower than those of the human expert. Fig. 7 plots the classification performance of the human expert on sensitivity-specificity plane. As can be seen, for both COVID-19 and CAP classification, the point of the human expert is located in the lower-right areas of our ROC curves, which indicates that given the same specificity of the human expert, our model can achieve higher sensitivity by adjusting the classification threshold. All of these results indicate that the DeepSC-COVID model offers high classification accuracy and speed, which offers capability for auxiliary medical diagnosis and large-scale COVID-19 screening.

D. Multi-task gain

To evaluate the gain of multi-task learning, additional experiments are conducted with single tasks of segmentation and classification. Specifically, we first remove the classification subnet from our model for the single segmentation task, and then remove the segmentation subnet for single classification task. Table II reports the results of single-task learning. As reported, the accuracy of single-task learning is lower than that of multi-task learning for both tasks. Specifically, for segmentation, the multi-task gain values are 7.1%, 7.5%, 2.9%, 7.5% and 3.4 mm in terms of DSC, sensitivity, specificity, NSD and RMSD, respectively. For classification, the multi-task gain achieves values of 9.3%, 9.4%, 5.4%, 9.3% and 4.4% in terms of accuracy, sensitivity, specificity, F₁-score and AUC, respectively. Additionally, the results of the single segmentation task are visualized in Fig. 6. The ROC curve of single classification task are shown in Fig. 7.

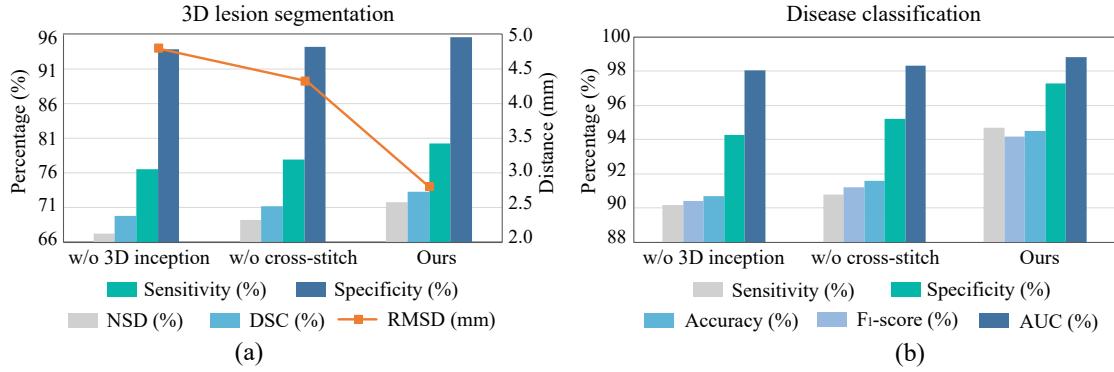


Fig. 8. The impacts of 3D inception block and cross-stitch unit on the performance of segmentation and classification. (a) Results of 3D lesion segmentation, in terms of sensitivity, specificity, NSD, DSC and RMSD. (b) Results of disease classification, in terms of sensitivity, specificity, accuracy, F_1 -score and AUC.

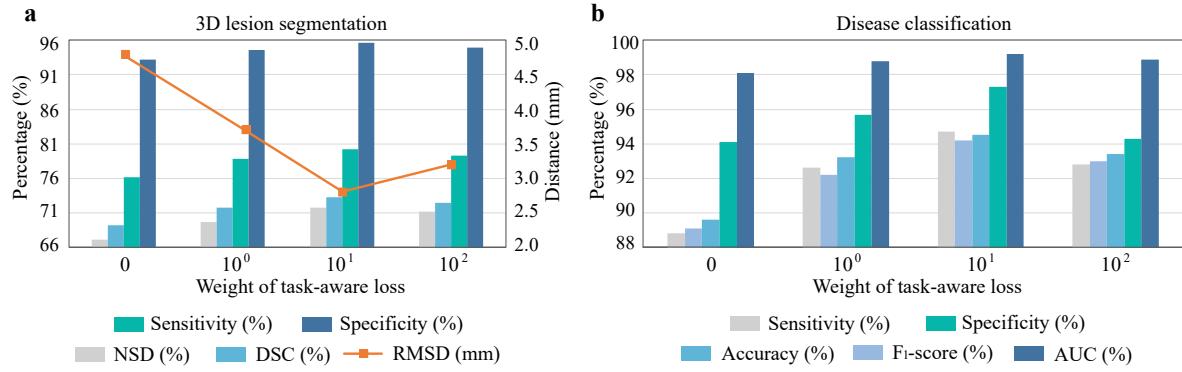


Fig. 9. Segmentation and classification results with different weights on the task-aware loss. (a) Results of 3D lesion segmentation, in terms of sensitivity, specificity, NSD, DSC and RMSD. (b) Results of disease classification, in terms of sensitivity, specificity, accuracy, F_1 -score and AUC.

E. Ablation study

Here, we analyze the effectiveness of different components in the proposed DeepSC-COVID model on the tasks of 3D lesion segmentation and disease classification through ablation study.

Effectiveness of 3D inception block. We first analyze the impact of 3D inception block on 3D lesion segmentation and disease classification. Specifically, we replace the 3D inception block by conventional 3D convolutional layer, in which the kernel size, stride and output channel are the same as the 3D inception block. Fig. 8 shows the segmentation and classification results with and without the 3D inception block. As shown, the performance of both segmentation and classification tasks significantly degrades after replacing the 3D inception block. This indicates the effectiveness of our 3D inception block in extracting effective multi-scale 3D features for both tasks.

Effectiveness of cross-stitch unit. We further conduct the ablation experiment to evaluate the impact of the cross-stitch unit on segmentation and classification performance, by removing it from the cross-task feature subnet in the proposed DeepSC-COVID model. Fig. 8 shows the segmentation and classification results with and without the cross-stitch unit. We can see from this figure that the performance of both the segmentation and classification degrades, when the cross-stitch unit is removed. This validates the positive contribution of cross-stitch unit to our model.

Effectiveness of task-aware loss. Finally, we evaluate the im-

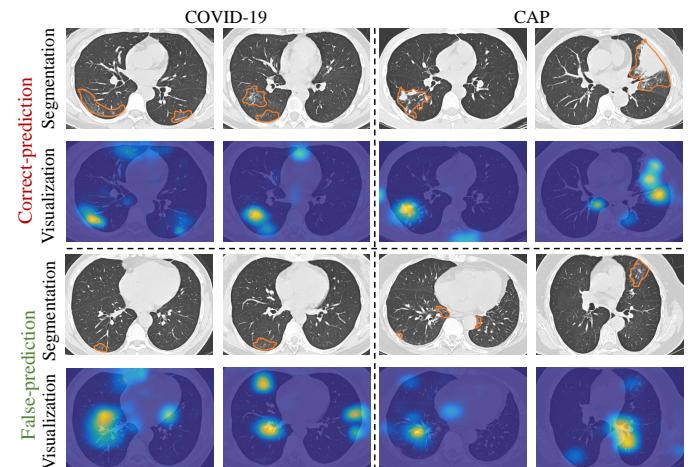


Fig. 10. Visualization maps and segmentation results of both correct-prediction and false-prediction cases.

pact of the proposed task-aware loss. To be specific, we train the DeepSC-COVID model with different weights λ_{ta} on the task-aware loss, i.e., $\lambda_{ta} = 0, 10^0, 10^1, 10^2$ in equation (11) of the main text. Note that $\lambda_{ta} = 0$ indicates that the task-aware loss is fully removed. Fig. 9 shows the segmentation and classification results with different λ_{ta} . As shown, the DeepSC-COVID model performs the worst, when the task-aware loss is fully removed

