

RCoNet: Deformable Mutual Information Maximization and High-order Uncertainty-aware Learning for Robust COVID-19 Detection

Shunjie Dong, Qianqian Yang, Yu Fu, Mei Tian, Cheng Zhuo, *Senior Member, IEEE*

Abstract—The novel 2019 Coronavirus (COVID-19) infection has spread world widely and is currently a major healthcare challenge around the world. Chest Computed Tomography (CT) and X-ray images have been well recognized to be two effective techniques for clinical COVID-19 disease diagnoses. Due to faster imaging time and considerably lower cost than CT, detecting COVID-19 in chest X-ray (CXR) images is preferred for efficient diagnosis, assessment and treatment. However, considering the similarity between COVID-19 and pneumonia, CXR samples with deep features distributed near category boundaries are easily misclassified by the hyper-planes learned from limited training data. Moreover, most existing approaches for COVID-19 detection focus on the accuracy of prediction and overlook the uncertainty estimation, which is particularly important when dealing with noisy datasets. To alleviate these concerns, we propose a novel deep network named $RCoNet_s^k$ for robust COVID-19 detection which employs *Deformable Mutual Information Maximization (DeIM)*, *Mixed High-order Moment Feature (MHMF)* and *Multi-expert Uncertainty-aware Learning (MUL)*. With DeIM, the mutual information (MI) between input data and the corresponding latent representations can be well estimated and maximized to capture compact and disentangled representational characteristics. Meanwhile, MHMF can fully explore the benefits of using high-order statistics and extract discriminative features of complex distributions in medical imaging. Finally, MUL creates multiple parallel dropout networks for each CXR image to evaluate uncertainty and thus prevent performance degradation caused by the noise in the data. The experimental results show that $RCoNet_s^k$ achieves the state-of-the-art performance on an open source COVIDx dataset of 15134 original CXR images across several metrics. Crucially, our method is shown to be more effective than existing methods with the presence of noise in the data.

Index Terms—Chest X-rays, COVID-19, $RCoNet_s^k$, DeIM, MHMF, MUL, Noisy Data, Uncertainty

I. INTRODUCTION

CORONAVIRUS disease 2019 (COVID-19) causes an ongoing pandemic that significantly impacts everyone's life since it was first reported, with hundreds of thousands of deaths and millions of infections emerging in over 200 countries [1], [2]. As indicated by the World Health Organization (WHO), due to its highly contagious nature and lack of corresponding vaccines, the most effective method to control the spread of COVID-19 infection is to keep social distance and contact tracing. Hence, early and fast diagnosis of COVID-19 has become significantly essential to control further spreading, and such that the patients could be hospitalized and receive proper treatment in time.

Since the emerge of COVID-19, *reverse transcription polymerase chain reaction (RT-PCR)*, as a viral nucleic

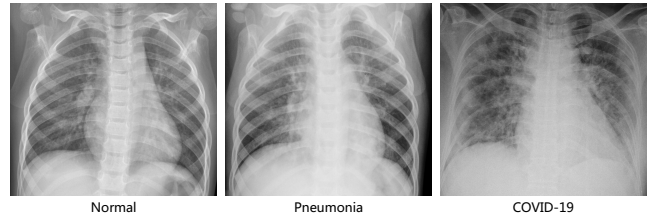


Fig. 1. Visual illustration of chest X-ray images, including normal, pneumonia and COVID-19.

acid detection method by gene sequencing, is the accepted standard for COVID-19 detection [3]. However, because of the low accuracy of RT-PCR and limited medical test kits in many hyper-endemic regions or countries, it is challenging to detect every individual affected by COVID-19 rapidly [4], [5]. Therefore, alternative testing methods, which are faster and more reliable than RT-PCR, are urgently needed to combat the disease.

Since most COVID-19 positive patients were diagnosed with pneumonia, radiological examinations could help detect and assess the disease. Recently, chest *computed tomography (CT)* has been shown to be efficient and reliable to achieve a real-time clinical diagnosis of COVID-19, outperforming over RT-PCR in terms of accuracy. Moreover, some deep learning based methods have been proposed for COVID-19 detection using chest CT images [6], [7], [8], [9]. For example, an adaptive feature selection approach was proposed in [10] for COVID-19 detection based on a trained deep forest model. In [11], an uncertainty vertex-weighted hypergraph learning method was designed to identify COVID-19 from *community acquired pneumonia (CAP)* using CT images. However, the routine use of CT, which is conducted via expensive equipments, takes considerably more time than X-ray imaging and brings a massive burden on radiology departments. Compared to CT, X-rays could significantly speed up disease screening, and hence become a preferred method for disease diagnosis.

Accordingly, deep learning based methods for detecting COVID-19 with *chest X-ray (CXR)* have been developed and shown to be able to achieve accurate and speedy detection [12], [13]. For instance, a tailored convolution neural network platform trained on open source dataset called COVIDNet in [14] was proposed for the detection of COVID-19 cases from CXR. Oh *et al.* [15] proposed a novel probabilistic gradient-weighted class activation map to enable infection segmentation and detection of COVID-19 on CXR images.

Fig. 1 shows three samples from the *COVID_x* dataset [14] which contains three different classes: normal, pneumonia and COVID-19. However, due to the similar pathological information between pneumonia and COVID-19 in the early stage, the CXR samples may have latent features distributed near the category boundaries, which can be easily misclassified by the hyper-plane learned from the limited training data. Moreover, to the best of our knowledge, most of the existing methods for COVID-19 detection are designed to extract the lower-dimension latent representations which may not be able to fully capture statistical characteristic of complex distributions (i.e., non-Gaussian distribution). Furthermore, quantifying uncertainty in COVID-19 detection is still a major yet challenging task for doctors, especially with the presence of noise in the training samples (i.e., *label* noise and *image* noise).

To address the above problems, we propose a novel deep network architecture, referred to as *RCoNet_s^k*, for robust COVID-19 detection which, in particular, contains the following three modules, i.e., *Deformable mutual Information Maximization* (DeIM), *Mixed High-order Moment Feature* (MHMF) and *Multi-expert Uncertainty-aware Learning* (MUL):

- The Deformable mutual Information Maximization (DeIM) module estimates and maximizes the mutual information (MI) between input data and learned high-level representations, which pushes the model to learn the discriminative and compact features. We employ deformable convolution layers in this module which are able to explore disentangled spatial features and mitigate the negative effect of similar samples across different categories.
- The Mixed High-order Moment Feature (MHMF) module, inspired by [16], fully explores the benefits of using a mix of high-order moment statistics to better characterize the feature distributions in medical imaging.
- The Multi-expert Uncertainty-aware Learning (MUL) creates multiple parallel dropout networks, each can be treated as an *expert*, to derive multiple experts based diagnosis similar to clinical practices, which improves the prediction accuracy. MUL also quantifies the prediction accuracy by obtaining the variance in prediction across different experts.
- The experimental results show that our proposal achieves the state-of-the-art performance in terms of most metrics both on open source COVID_x dataset of 15134 original CXR images and that of noisy setting.

The remaining of this paper is organized as follows: In Section II, we review related works on mutual information estimation and uncertainty learning as well. In Section III, after an overview of our proposed approach, we discuss the main components of *RCoNet_s^k*. In Section IV, we compare our proposed architecture with the existing deep learning based methods evaluated on a public available dataset of CXR images and also the same dataset but under noisy conditions. And we also conduct extensive experiments to demonstrate the benefits of DeIM, MHMF and MUL on the performance of the system. Finally, we conclude this paper in Section V.

II. BACKGROUND AND RELATED WORKS

In this section, we introduce related works on mutual information estimation and uncertainty learning that lay the foundation of this paper.

A. Mutual Information Estimation

Mutual information (MI), as a fundamental concept in information theory, is widely applied to unsupervised feature learning for quantifying the correlation between random variables. MI has been exploited in a wide range of domains and tasks, including biomedical sciences [17], blind source separation (BSS, e.g., independent component analysis [18]), feature selection [19], [20] and causal inference [21]. For example, the object tracking task considered in [22] was treated as a problem of optimizing the mutual information between features extracted from a video with most color information removed and those from the original full-color video. Closely related work presented in [23] considered learning representations to predict cross-modal correspondence by maximizing MI between features from the multi-view encoders and the content of the held-out view. Moreover, Mutual Information Neural Estimation (MINE) proposed by [24] was designed to learn a general-purpose estimator of the MI between continuous variables based on dual representations of the KL-divergence, which are scalable, flexible and, most crucially, trainable via back-propagation. Based on MINE, our proposal estimates and maximizes the CXR image inputs and the corresponding latent representations to improve diagnosis performance.

B. Uncertainty in Deep Learning

Aiming at combating the significant negative effects of uncertainty in deep neural networks, uncertainty learning has been getting lots of research attention, which facilitates the reliability assessment and solves risk-based decision-making problems [25], [26], [27]. In recent years, various frameworks have been proposed to characterize the uncertainty in the model parameters of deep neural networks, referred to as *model uncertainty*, due to the limited size of training data [28], [29], which can be reduced by collecting more training data [26], [30], [31]. Meanwhile, another kind of uncertainty in deep learning, referred to as *data uncertainty*, measures the noise inherent in given training data, and hence cannot be eliminated by having more training data [32]. To combat these two kinds of uncertainty, lots of works on various computer vision tasks, i.e., face recognition [25], semantic segmentation [33], object detection [34], person re-identification [35], etc., have introduced deep uncertainty learning to improve the robustness of deep learning model and interpretability of discriminant. For face recognition task in [26], an uncertainty-aware probabilistic face embedding (PFE) was proposed to represent face images as distributions by utilizing data uncertainty. Exploiting the advantage of Bayesian deep neural networks, one recent study [36] leveraged the model uncertainty for analysis and learning of face representations. To our knowledge, our proposal is the first work that utilizes the high-order moment statistics and multiple expert networks to estimate uncertainty for COVID-19 detection using CXR images.

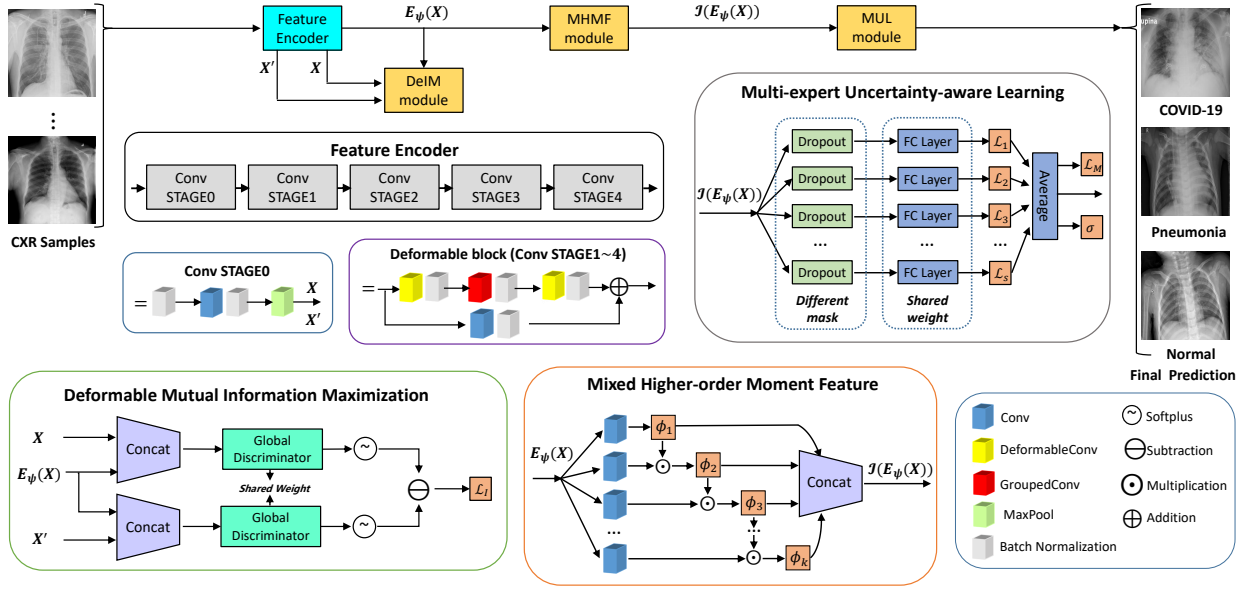


Fig. 2. The architecture of RCoNet_s^k for COVID-19 detection.

III. METHOD

In this section, we introduce the novel RCoNet_s^k for robust COVID-19 detection, which incorporates *Deformable mutual Information Maximization* (DeIM), *Mixed High-order Moment Feature* (MHMF) and *Multi-expert Uncertainty-aware Learning* (MUL), as illustrated in Fig. 2. k is the number of levels of moment features that are combined in MHMF, and s is the number of the expert network in MUL, which will be further clarified in the sequel. The CXR images are first processed by DeIM which consists of a stack of deformable convolution layers, extracting discriminative features. The compact features are then fed into MHMF module to generate high-order moment latent features, reducing negative effects caused by similar images. The proposed MUL utilizes the learned high-order features to generate final diagnoses.

A. Deformable Mutual Information Estimation and Maximization

Due to the similarity between COVID-19 and pneumonia in the latent space, we propose Deformable mutual Information Maximization (DeIM) to extract discriminative and informative features, reducing the negative influence caused by the lack of distinctiveness in the deep features. In particular, we train the model by maximizing the mutual information between the input and corresponding latent representation.

We use a stack of five convolutional stages, as shown in Fig. 2, to encode inputs into latent representations, which is denoted by a differentiable parametric function E_ψ :

$$E_\psi : \mathcal{X} \rightarrow \mathcal{Z}, \quad (1)$$

where ψ denotes the set of all the trainable parameters in these layers, and \mathcal{X} and \mathcal{Z} denote the input and output spaces, respectively.

The detailed architecture of each convolutional stage is presented in Fig. 2, which consists of several convolutional

layers each followed by a batch normalization layer. Note that we employ deformable convolutional layers which can better extract spatial information of the irregular infected area compared to conventional convolutional layers. More specifically, regular convolution operates on pre-defined rectangular grid from an input image or a set of input feature maps, while the deformable convolution operates on deformable grids that each grid point is moved by a learnable offset. For example, the receptive grid \mathcal{P} of a regular convolution with kernel size 3×3 is fixed and can be given by:

$$\mathcal{P} = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\}, \quad (2)$$

while, for deformable convolution, the receptive grid is moved by the learned offsets $\Delta p_n \in \mathbb{R}^2$ and the output is given as follows:

$$b(p_0) = \sum_{p_n \in \mathcal{P}} w(p_n) \cdot a(p_0 + p_n + \Delta p_n). \quad (3)$$

where $b(p_0)$ denotes the value at location p_0 on the output feature map b , p_n enumerates the locations in \mathcal{P} , $w(p_n)$ represents the weight at location p_n of the kernel, and $a(\cdot)$ is value at given location on the input feature map. We can see that with the introduction of offsets Δp_n , the receptive grid is no longer fixed to be a rectangle, and instead is deformable.

We optimize E_ψ by maximizing the mutual information between the input and the output, i.e., $I(X; Z)$, where $Z \triangleq E_\psi(X)$. The precise mutual information requires knowledge probability density functions (PDFs) of X and Z , which is intractable to obtain in practice. To overcome this issue, Mutual Information Neural Estimation (MINE) proposed in [24] estimates mutual information by using a lower-bound on the Donsker-Varadhan representation [37] of the KL-divergence:

$$\begin{aligned} I(X; Z) &:= D_{KL}(\mathbb{J} || \mathbb{M}) \geq \hat{T}_\theta^{(DV)}(X; Z) : \\ &= \mathbb{E}_{\mathbb{J}}[T_\theta(x, z)] - \log \mathbb{E}_{\mathbb{M}}[e^{T_\theta(x, z)}], \end{aligned} \quad (4)$$

where \mathbb{J} represents the joint probability of X and Z , i.e., $\mathbb{J} \triangleq P(X, Z)$, and \mathbb{M} denotes the product of marginal probabilities of X and Z , $\mathbb{M} \triangleq P(X)P(Z)$. $T_\theta : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$ denotes a *global discriminator* modeled by a neural network with parameters θ , which is trained to maximize $\hat{I}_\theta^{(DV)}(X; Z)$ to approximate the actual mutual information. Hence, we can simultaneously estimate and maximize $I(X; E_\psi(X))$ by maximizing $\hat{I}_\theta^{(DV)}(X; Z)$:

$$(\hat{\theta}, \hat{\psi}) = \underset{\theta, \psi}{\operatorname{argmax}} \hat{I}_\theta^{(DV)}(X; E_\psi(X)). \quad (5)$$

Since the encoder E_ψ and the mutual information estimator T_θ are optimized simultaneously with the same objective function, we can share some layers between them, and replace the T_θ with $T_{\theta, \psi}$ to account for this fact.

Since we are primarily interested in maximizing the mutual information rather than estimating the precise value, we can alternatively use a Jensen-Shannon MI estimator (JSD) [38], which offers more interpretable trade-off:

$$\hat{I}_{\theta, \psi}^{(DeJSD)}(X; E_\psi(X)) := \mathbb{E}_{\mathbb{P}} \left[-\log \left(1 + e^{-T_{\theta, \psi}(x, E_\psi(x))} \right) \right] - \mathbb{E}_{\mathbb{P} \times \tilde{\mathbb{P}}} \left[\log \left(1 + e^{T_{\theta, \psi}(x', E_\psi(x'))} \right) \right], \quad (6)$$

where x is an input sample of an empirical probability distribution \mathbb{P} , x' denotes a fake sample from distribution $\tilde{\mathbb{P}}$, where $\tilde{\mathbb{P}} = \mathbb{P}$. This estimator is illustrated by the DeIM block shown in Fig. 2, which has the latent representation $E_\psi(x)$, the input sample x and the fake sample x' as input, and the difference between the outputs of the two softplus operations as the estimation of MI.

Another alternative MI estimator is called Noise-Contrastive Estimator (NCE) [39], which is defined as:

$$\hat{I}_{\theta, \psi}^{(DeNCE)}(X; E_\psi(X)) := \mathbb{E}_{\mathbb{P}} \left[T_{\theta, \psi}(x, E_\psi(x)) - \mathbb{E}_{\tilde{\mathbb{P}}} \left[\log \sum_{x'} e^{T_{\theta, \psi}(x', E_\psi(x'))} \right] \right]. \quad (7)$$

The experiments have found that using the NCE estimator outperforms the JSD estimator in some cases, but appears to be quite similar most of the time.

The existing works [40] that implement these estimators use some latent representation of x , which is then merged with some randomly generated features to obtain “fake” samples that satisfy $\mathbb{P} = \tilde{\mathbb{P}}$. In contrast, we use the samples from other categories as the “fake” samples, i.e., x' , instead. For example, if the input is a pneumonia sample, then the fake sample is either a normal or COVID sample. We note that this can push the learned encoder to derive more distinguishable features for samples from different categories.

B. Mixed High-order Moment Feature

The presence of the image noise and label noise in CXR datasets may cause image latent representations generated by deep neural networks to be scattered in the entire feature space. To deal with this issue, [25], [26], [35] represent each image as a Gaussian distribution, that is defined by a mean (a standard feature vector) and a variance. However, the deep features of

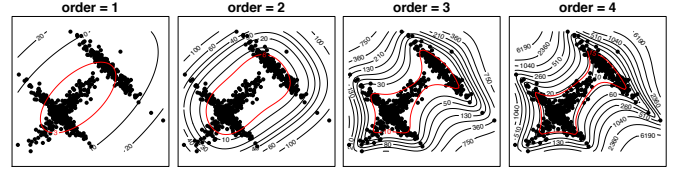


Fig. 3. Data points from three Gaussian distributions and the corresponding moment feature of order 1 to 4

CXR samples we considered in this paper typically follow a complex, non-Gaussian distribution [41], [42], which cannot be fully captured by its first-order (mean) or second-order statistics (variance).

We seek a better combination of different orders of statistics to more precisely characterize the latent representation of the CXR images. We illustrate the moment features of different orders [16] in Fig. 3, where we plot 350 data points in \mathbb{R}^2 sampled from a distribution that combines three different Gaussian distributions. We can observe that the high-order moment features are more expressive of statistical characteristic compared to low-order one. More specifically, it captures the shape of the cloud of samples more accurately. Therefore, we include the Mixed High-order Moment Feature (MHMF) module in the proposed model, as shown in Fig. 2, which outputs a combination of high-order moment features with the latent representation $E_\psi(X)$ as input. This will potentially solve the scattering problem, and, more importantly, capture the subtle differences between CXR images of similar categories, i.e., pneumonia and COVID-19 in our case.

We show how to obtain the complicated high-order moment feature in the following. Define r -th order moment feature as $\phi_r(\mathbf{a})$, where $\mathbf{a} \in \mathbb{R}^{H \times W \times C}$ denotes a latent feature map of dimension $H \times W \times C$. Lots of recent works adopt the Kronecker product to compute high-order moment feature [42]. However, calculating Kronecker product of high dimensional feature maps is significantly computational intensive, and hence infeasible for real-world applications. Inspired by [43], [44], [45], we approximate $\phi_r(\mathbf{a})$ by exploiting r random projectors which relies on certain factorization schemes, such as Random Maclaurin [46]. We use 1×1 convolution kernels as the random projectors to estimate the expectations of high-order moment features. That is,

$$\phi_r(\mathbf{a}) \approx \mathcal{K}_1(\mathbf{a}) \odot \mathcal{K}_2(\mathbf{a}) \odot \cdots \odot \mathcal{K}_r(\mathbf{a}) \in \mathbb{R}^{H \times W \times C}, \quad (8)$$

where \odot represents the Hadamard (element-wise) product, and $\mathcal{K}_1, \mathcal{K}_2, \dots, \mathcal{K}_r$ are 1×1 convolution kernels with random weights.

Note that Random Maclaurin produces an estimator that is independent of the input distribution, which causes the estimated high-order moments to contain non-informative high-order moment components. We eliminate these components by learning the weights of the projectors, i.e., the 1×1 convolution kernels, from the data. Also note that the Hadamard product of a number of random projectors may end up with the estimated high-order moment features to be similar to low-order ones. To solve this problem, we use a recursive way to estimate the

high-order moments instead,

$$\phi_r(\mathbf{a}) = \phi_{r-1}(\mathbf{a}) \odot \mathcal{K}_r(\mathbf{a}). \quad (9)$$

Since different order moments capture different informative statistics, we design the MHMF module to keep the estimated moments of different levels of order, as shown in Fig. 2, the output of which is given as:

$$\mathcal{J}(\mathbf{a}) = [\phi_1(\mathbf{a}); \phi_2(\mathbf{a}); \dots; \phi_r(\mathbf{a})] \in \mathbb{R}^{H \times W \times rC}. \quad (10)$$

Hence, $\mathcal{J}(\mathbf{a})$ is rich enough to capture the complicated statistics, and produce discriminative features for the input of different categories.

C. Multi-expert Uncertainty-aware Learning

The MHMF module, as described in section III-B, generates mixed high-order moment features of each sample in the latent space, which we aim to further exploit to derive compact and disentangled information for COVID-19 detection. Meanwhile, quantifying uncertainty in disease detection is undoubtedly significant to understand the confidence level of computer-based diagnoses. Motivated by the clinical practices, we present a novel neural network in this section, referred to as Multi-expert Uncertainty-aware Learning (MUL), which takes in the mixed high-order moment features and outputs the prediction and the quantification of the diagnostic uncertainty caused by the noise in the data.

The structure of Multi-expert Uncertainty-aware Learning module is shown in Fig. 2, which consists of multiple dropout layers that process the output from MHMF in parallel, each of which together with the following several fully connected layers can be regarded as an *expert* for COVID-19 detection. We note that each dropout layer uses different masks which results in different subsets of latent information to be kept, while the following fully connected layers share the same weights across different experts. The masks for the dropout layers are generated randomly at each iteration during training, but fixed during the inference time. We denote the input-output function of each expert by $C_e^j(\cdot)$, $j = 1, \dots, N$, where N is the total number of experts. Hence, we have the classification loss \mathcal{L}_e^j of j -th expert given as follows:

$$\mathcal{L}_e^j = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_w(C_e^j(\mathcal{J}(E_\psi(x_i))), y_i), \quad (11)$$

where n represents the total number of labeled CXR samples, and y_i denotes the one-hot representation of the class label, $i = 1, \dots, n$, and we recall that $\mathcal{J}(\cdot)$ denotes the MHMF operation given in Eq. (10) and $E_\psi(\cdot)$ is the preprocessing step on the CXR samples. Note that, the total number of COVID-19 cases is much smaller than non-COVID cases, i.e., normal and pneumonia cases. This imbalance in the dataset leads to a high ratio of false-negative classification. To mitigate this negative effect, we employ a weighted cross-entropy $\mathcal{L}_w(\cdot)$ given as follows:

$$\mathcal{L}_w(\hat{y}_i, y_i) = -\frac{1}{C} \sum_{c=1}^C \lambda_c \cdot y_{i,c} \log \hat{y}_{i,c}, \quad (12)$$

where C is the total number of classes, $y_{i,c}$ is the c -th element of y_i , and $\hat{y}_{i,c}$ denotes the corresponding prediction. λ_c represents the weight that controls how much the error on class c contributes to the loss, $c = 1, \dots, C$. Finally, the loss \mathcal{L}_M of the whole MUL module is derived by averaging the loss values of all the experts:

$$\mathcal{L}_M = \frac{1}{N} \sum_{j=1}^N \mathcal{L}_e^j. \quad (13)$$

We use the variance of classification loss \mathcal{L}_e^j with regards to the average loss \mathcal{L}_M to quantify the uncertainty, denoted by σ , which is given as:

$$\sigma = \frac{1}{N} \sum_{j=1}^N (\mathcal{L}_M - \mathcal{L}_e^j)^2. \quad (14)$$

The proposed MUL module improves the diagnostic accuracy as the final prediction combines the results from multiple experts, and also mitigates the negative effects caused by the noise in the data by introducing the dropout layers. Moreover, the experiments have revealed that the more experts in MUL module the faster the system converges during training.

D. Training

The whole architecture of RCoNet_s^k is presented in Fig. 2, where the CXR images are first processed by a stack of deformable convolution layers, then transformed to high-order moment latent features by the MHMF module, which are then fed to the MUL module to generate final diagnoses. The loss used to optimize RCoNet_s^k is given as follows

$$\mathcal{L}_{total} = \mathcal{L}_M - \alpha \mathcal{L}_I, \quad (15)$$

where \mathcal{L}_M is the prediction loss given by Eq. (13), and \mathcal{L}_I denotes the mutual information between the input X and the latent representation $E_\psi(X)$ estimated by either Eq. (6) or Eq. (7). α is a positive hyper-parameter that governs how much \mathcal{L}_M and \mathcal{L}_I contribute to the total loss. During training, the trainable parameters of the whole systems are updated iteratively to minimize \mathcal{L}_{total} , which is to jointly minimize the prediction loss \mathcal{L}_M thus to improve the accuracy, and maximize the mutual information \mathcal{L}_I .

IV. EXPERIMENTS AND RESULTS

A. Dataset

We use a public chest X-ray dataset, referred to as COVIDx, to evaluate the proposed model, which is published by the authors of COVID-Net [14]. This dataset contains a total of 13975 CXR images from 13870 patients of 3 classes: (a) normal (no infections); (b) pneumonia (non-COVID-19 pneumonia); (c) COVID-19. It contains samples from five open source available data repositories <https://github.com/lindawangg/COVID-Net/blob/master/docs/COVIDx.md>. Three random CXR samples of these three classes are shown in Fig. 1. To reduce the negative effect caused by extremely unbalanced training samples, i.e., very limited number of COVID-19 positive cases compared to the other two categories, we further include other open-source

TABLE I
DETAILS OF PATIENT DATA USED FOR TRAINING AND TESTING

Data	Number of Patients Per Class			Total Patients
	Normal	Pneumonia	COVID-19	
Train	7966	5451	207	13624
Test	885	594	31	1510

CXR datasets from <https://www.kaggle.com/c/rsna-pneumonia-detection-challenge/data>. Following [14], [47], the dataset is finally divided into 13624 training and 1510 test samples. The numbers of samples from different categories used for training and testing are summarized in Table I. Moreover, we also adopted various data augmentation techniques to generate more COVID-19 training samples, such as flipping, translation, rotation using random five different angles, to tackle the data imbalance issue such that the proposed model can learn an effective mechanism of detecting COVID-19.

B. Evaluation Metrics

In our experiments, we use the following six metrics to evaluate the COVID-19 detection performance of different approaches:

- Accuracy (ACC): ACC calculates the proportion of images that are correctly identified. $ACC = \frac{TP+TN}{TP+TN+FP+FN}$.
- Sensitivity (SEN): SEN is the ratio of the positive cases that have been correctly detected to all the positive cases. $SEN = \frac{TP}{TP+FN}$.
- Specificity (SPE): SPE is the ratio of the negative cases that have been correctly classified to all the negative cases. $SPE = \frac{TN}{TN+FP}$.
- Balance (BAC): BAC is the mean value of SEN and SPE . $BAC = \frac{SEN+SPE}{2}$.
- Positive Predictive Value (PPV): PPV is the ratio of correctly detected positive cases to all cases that are detected to be positive. $PPV = \frac{TP}{TP+FP}$.
- F1-score ($F1$): $F1$ uses a combination of accuracy and sensitivity to calculate a balanced average result. $F1 = \frac{2 \times ACC \times SEN}{ACC+SEN}$.

TN , TP , FN and FP represent the total number of true negatives, true positives, false negatives, and false positives, respectively.

C. Compared Methods

We compare the proposed RCoNet $_s^k$ with the following five existing deep learning methods for COVID-19 detection:

- PbcNN [15]: A patch-based convolutional neural network with a relatively small number of trainable parameters.
- COVID-Net [14]: A tailored deep convolutional neural network that uses a projection-expansion-projection design pattern.
- DenseNet-121 [48]: A densely connected convolutional network that connects each layer to every other layer in a feed-forward fashion.

TABLE II
DETAILS OF 10% NOISY PATIENT DATA USED FOR TRAINING.

Training Date	Clean	Noise	Total
Normal	7170	796 (Peumonia+COVID-19)	7966
Pneumonia	4906	545 (COVID-19+Normal)	5451
COVID-19	187	20 (Peumonia+Normal)	207

- CoroNet [49]: A deep convolutional neural network model based on Xception architecture pre-trained on ImageNet dataset.
- ReCoNet [47]: A residual image-based COVID-19 detection network that exploits a CNN-based multi-level preprocessing filter block and a multi-task learning loss.

D. Implementation

We implement our RCoNet $_s^k$ using the PyTorch library and apply ResNeXt [50] as the backbone network. We train the model with the Adam optimizer with an initial learning rate of 2×10^{-4} and a weight decay factor of 1×10^{-4} . All the experiments are run on an NVIDIA GeForce GTX 1080Ti GPU. We set the batch size to be 8, and resize all images to 224×224 pixels. The hyperparameter α in the loss function given in Eq. (15) is set to be within the range of $[0, 0.4]$. The drops rate of each dropout layer in the MUL module is randomly chosen from $\{0.1, 0.3, 0.5\}$. The loss weight λ_c for each category, which is used to calculate the weighted sum of the loss as given in Eq. (12), is set to be 1, 1, and 20 for the normal, pneumonia, COVID-19 samples, respectively, corresponding to the number of training samples in each. We adopt 5-fold cross-validation training that we randomly divide the training sets into five equal-size subsets and train the model five times that using different four subsets for training, and the remaining one for validation each time. We also evaluate our proposed model with different number of order moments for the MHMF module k , and different number of experts s .

To evaluate the performance of the proposed model with the presence of label noise, we derive a noisy dataset from the given dataset in the following way: we randomly select a given percentage of training samples in each category, and assign wrong labels to these sample. In particular, to ensure that the fake COVID-19 samples are less than the real ones, we assign the COVID-19 labels to selected normal and pneumonia samples in a way the the number of normal and pneumonia samples assigned with COVID-19 label equals to the number of COVID-19 samples assigned with either normal and pneumonia label. We show a realization of the derived noisy dataset when the percentage of fake samples is set to be 10% in Table II.

E. Results and Discussions

Performance on Clean Data: The numerical results on the clean dataset without any artificial noise added are shown in Table III. The results are presented in the form of $a \pm b$, where a and b denote the average and variance values of each metric on five independent experiments, respectively. We can see that RCoNet $_4^5$, i.e., the proposed model with $k = 4$ levels

TABLE III
PERFORMANCE COMPARISON OF DIFFERENT APPROACHES FOR COVID-19 DETECTION ON THE COVIDX DATASET

Method	ACC (%)	SEN (%)	SPE (%)	BAC (%)	PPV (%)	F1 (%)	Param (M)	FLOPs (G)
PbCNN [15]	88.90±1.63	85.90±1.69	96.40±2.10	91.15±1.31	88.65±1.52	87.37±2.14	11.60	-
COVID-Net [14]	95.10±1.34	91.37±1.37	95.76±2.04	93.57±0.89	94.73±0.97	93.20±0.85	117.4	15.10
DenseNet-121 [48]	97.40±1.67	96.08±0.88	97.23±1.01	96.66±1.21	96.05±1.00	96.74±1.04	7.61	bf5.59
CoroNet [49]	95.00±1.58	96.90±1.57	97.50±1.93	97.20±1.07	95.00±1.03	95.60±0.95	33.00	-
ReCoNet [47]	97.48±1.05	97.39±1.67	97.53±1.28	97.46±0.87	97.17±0.76	97.43±0.59	2.52	7.68
RCoNet _{s=3} ^{k=4}	96.12±0.33	95.71±0.41	96.38±0.29	96.05±0.20	95.86±0.62	95.91±0.56	.73	7.61
RCoNet _{s=4} ^{k=4}	96.78±0.57	96.48±0.69	96.91±0.74	96.70±0.34	96.94±0.53	96.63±0.58	6.74	7.70
RCoNet _{s=3} ^{k=4}	97.46±0.43	97.25±0.79	97.62±0.40	97.44±0.82	97.59±0.91	97.35±0.38	6.75	7.79
RCoNet _{s=4} ^{k=4}	97.89±0.53	97.33±0.45	98.24±0.39	97.79±0.62	97.93±0.74	97.61±0.48	6.77	7.91
RCoNet _{s=4} ^{k=3}	97.50±0.62	97.76±0.87	97.18±0.63	97.47±0.73	97.10±0.91	97.63±0.71	6.77	8.00

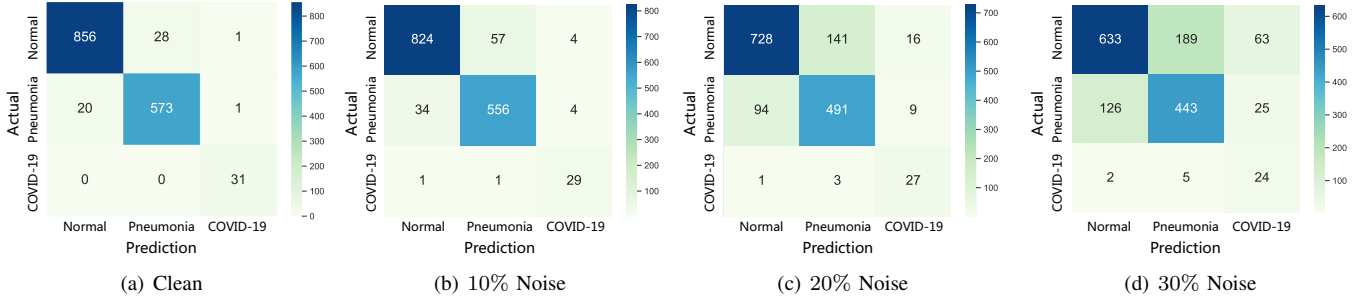


Fig. 4. Confusion matrices of the proposed RCoNet_s^k trained on noisy dataset with different percentages of noisy samples.

TABLE IV
PERFORMANCE COMPARISON OF DIFFERENT APPROACHES ON COVIDX DATASET WITH NOISY SAMPLES

Noise	Method	ACC (%)	SEN (%)	SPE (%)
10%	PbCNN [15]	83.22	81.98	89.01
	COVID-Net [14]	91.03	87.94	90.62
	DenseNet-121 [48]	91.97	87.94	92.17
	CoroNet [49]	89.45	88.74	90.06
	ReCoNet [47]	91.63	90.82	91.16
	RCoNet _{s=3} ^{k=4}	92.78	92.21	93.51
	RCoNet _{s=4} ^{k=4}	92.98	93.39	93.12
20%	PbCNN [15]	78.42	75.90	80.29
	COVID-Net [14]	82.51	82.77	81.95
	DenseNet-121 [48]	82.16	81.01	82.21
	CoroNet [49]	82.33	81.10	81.89
	ReCoNet [47]	83.26	82.72	83.17
	RCoNet _{s=3} ^{k=4}	84.18	84.56	85.79
	RCoNet _{s=4} ^{k=4}	84.30	84.01	85.99
30%	PbCNN [15]	67.76	66.47	70.61
	COVID-Net [14]	71.98	70.13	71.55
	DenseNet-121 [48]	72.74	72.36	72.96
	CoroNet [49]	71.87	72.02	71.54
	ReCoNet [47]	73.26	72.53	73.11
	RCoNet _{s=3} ^{k=4}	74.56	74.20	75.54
	RCoNet _{s=4} ^{k=4}	74.69	74.51	76.94
RCoNet _{s=4} ^{k=3}	74.88	74.37	75.21	

of mixed moment features and $s = 4$ experts, achieves notable performance improvement over the comparison methods in terms of most metrics considered, including ACC, SPE, BAC, PPV and F1 score. We note the performance of RCoNet_s^k can be further improved with a different set of k and s . For instance, RCoNet_{s=3}^{k=4} achieves better SEN and F1 score than RCoNet_{s=4}^{k=4}. The higher ACC and F1 score validate that RCoNet_s^k is able to obtain latent features, i.e., the mixed moment features of

different levels of order, that maintains inter-class separability and intra-class compactness better than other models. Note that RCoNet_{s=3}^{k=4} leads to a higher SEN than all other methods, which is particularly important to COVID-19 detection, since successfully detecting COVID-19 positive cases is the key to control the spread of this super contagious disease. Moreover, it can be observed that RCoNet_s^k has smaller variance compared to the others, which demonstrates the robustness and stability of our model.

We also evaluate the complexity of the proposed model in terms of numbers of parameters and computational cost, i.e., Float-point operations (FLOPs), which is presented in Table III. It can be observed that the proposed model has much fewer parameters than several existing methods, except ReCoNet. However, we note that the FLOPs of RCoNet_s^k is quite close to that of ReCoNet, which means it takes a similar amount of time to diagnose COVID-19 from CXR images by these two model. We can also observe that the increase of k and s , i.e., the number of mixed moment features and the number of experts in MUL, only causes a small, or even neglectable, amount of increase in the number of parameters and FLOPs as well, which suggests that we can improve the performance of the proposed model by optimizing k and s , without the concern on the significant increase of the complexity.

Performance on Noisy Data: We further compare the proposed model to the existing ones when there is noise present in the training dataset. We generate three noisy training datasets in the aforementioned way from the clean dataset with 10%, 20% and 30% samples with wrong labels, respectively. The results, which we take the averages from five independent experiments, are presented in Table IV. It can be easily seen that the more fake samples we add the more it degrades

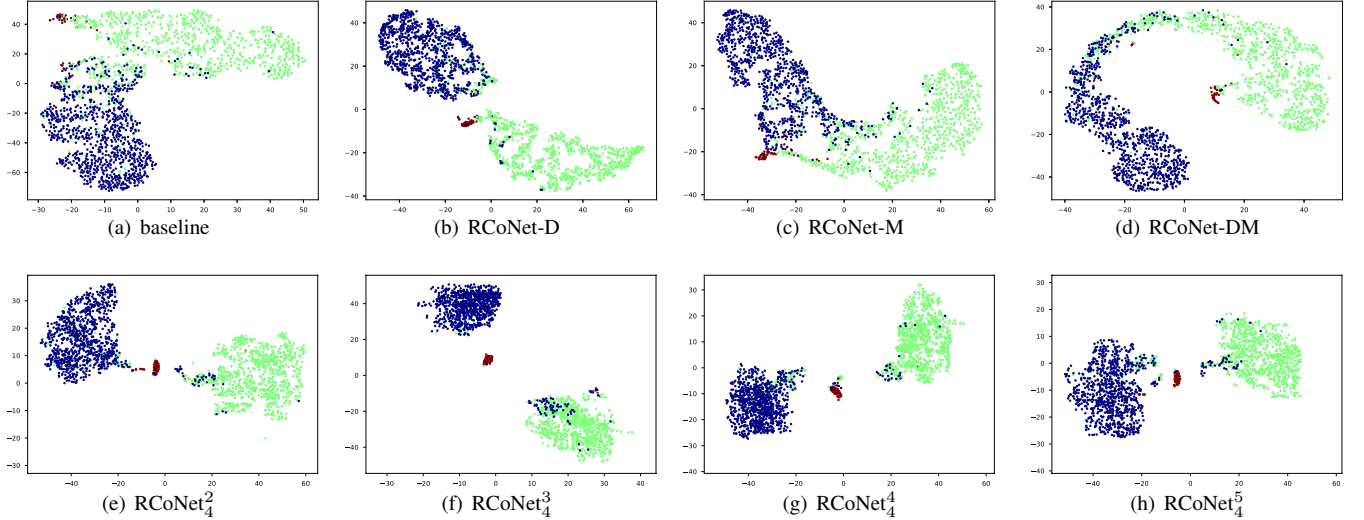


Fig. 5. The t-SNE visualization of the latent features generated by different methods. Blue, green and red dots represent normal, pneumonia and COVID-19 samples, respectively.

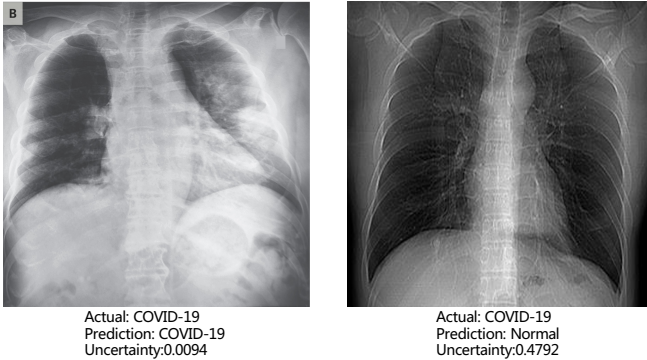


Fig. 6. Example CXR samples with their predictions and the corresponding uncertainty levels by RCoNet_4^k .

the performance of all the methods. Note that the proposed RCoNet_4^k still gets the state-of-the-art results in all considered cases with different percentages of noisy samples in the training dataset. Moreover, the performance gain over the existing methods slightly increases with the ratio of noisy samples, verifying that our model is more robust to the noise. Note that the extreme case of 30% noisy samples leads to great performance degradation of all the models. In practice, the percentage of label noise is usually around 10% to 20%. We present the confusion matrices in Fig. 4 to summarize the prediction accuracy of different categories. We can observe that, although with very limited number of COVID-19, our model still maintains high accuracy of detecting COVID-19 cases, even with the presence of noisy samples.

Uncertainty Estimation: One remarkable advantage of our model is the ability to quantify the uncertainty in the final prediction, which is significantly crucial for COVID-19 detection. This is done by obtaining the variance in the output of different experts in MUL as described in Section III-C. The larger the variance is, the more different experts disagree with each other, and, hence, the more uncertain the model is about

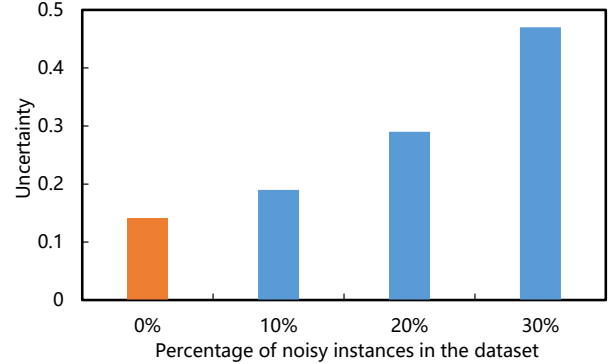


Fig. 7. Comparison on uncertainty level of the predictions by RCoNet_4^k .

the final prediction. We present two CXR samples in Fig. 6, including the predictions and the corresponding uncertainty level by RCoNet_s^k . We can see that the correctly classified CXR image has a low uncertainty level about its prediction, i.e., 0.0094, and the misclassified CXR sample with a high uncertainty level, i.e., 0.4792, suggests that an alternative way of diagnosis should be sought to correct this prediction. This greatly improves the reliability of the prediction by RCoNet_s^k , and reduces the chance of misdiagnosis. We also show in Fig. 7 the average uncertainty levels of RCoNet_s^k trained on clean and noisy datasets with different ratios of noisy samples. It can be observed that the uncertainty level increases almost linearly with the percentage of noisy samples in the dataset, which highlights the negative impact of noise on model training.

F. Analysis

We further numerically analyse the benefits of the three key modules of RCoNet_s^k , i.e., the DeIM, MHMF and MUL modules in this section.

Effectiveness of DeIM: We utilize t-SNE method [51] to visualize the latent features, presented in Fig. 5, which are

TABLE V
IMPACT OF THE MHMF AND MUL ON THE MODEL PERFORMANCE.

RCoNet_s^k	$s=1$	$s=2$	$s=3$	$s=4$	$s=5$	$s=6$	$s=7$
$k=1$	95.4	95.7	95.9	96.1	96.1	96.0	95.8
$k=2$	96.3	96.4	96.6	96.8	96.8	96.7	96.4
$k=3$	97.2	97.2	97.3	97.5	97.4	97.3	97.3
$k=4$	97.4	97.6	97.8	97.9	97.9	97.7	97.5
$k=5$	97.2	97.3	97.3	97.4	97.5	97.5	97.3
$k=6$	96.8	97.0	97.0	97.1	97.0	96.9	96.9

generated by the bottleneck layers of the baseline model, i.e., ResNeXt, RCoNet_s^k and three variants of RCoNet_s^k : (a) RCoNet-D: a model contains only DeIM; (b) RCoNet-M: a model contains only MUL; (c) RCoNet-DM: a model contains DeIM and MUL but not MHMF. Comparing the latent feature distribution by the baseline model shown in Fig. 5(a), and that by RCoNet-D presented in Fig. 5(b), we can tell that the introduction of DeIM leads to better class separation in the latent space.

Effectiveness of MHMF: We can observe in Fig. 5(a) - Fig. 5(d) that the latent features of the COVID-19 samples, generated by the models without MHMF, always distribute around the category boundary, and are not quite separable from those of some pneumonia samples. Meanwhile, the latent feature distributions presented in Fig. 5(e) - Fig. 5(h) derived by the models with MHMF show significant separability between different categories, which implies that MHMF can extract discriminative features. We also include numerical results of RCoNet_s^k , trained and tested on COVIDx dataset, with regards to different values of k , i.e., the number of levels of the moment features to be mixed, and s , i.e., the number of experts, in Table V in terms of accuracy. We can observe that, for a given value of s , the accuracy increases first with the value of k but decreases after k is larger than 4. It demonstrates that including more levels of moment feature could improve the model performance. However, the overly high-order moments may lead to performance degradation, which may be because these features are not useful for COVID detection.

Effectiveness of MUL: From Table V, we observe that, for a given value of k , accuracy increases first with the value of s but saturates around $s = 5$. This implies that having more experts in MUL can increase the prediction accuracy but it is not necessary to have too many.

Parameter Sensitivity and Convergence: We evaluate how sensitive the model performance in terms of accuracy to the value of α . We show the average accuracy of five independent experiments by RCoNet_4^4 trained on the dataset with different ratios of noisy samples in Fig. 8. As we can see, the larger α , which means the prediction loss, i.e., \mathcal{L}_M , contributes less to the total loss, not necessarily leads to degradation in the accuracy. This means maximizing the mutual information between the input and the latent features could keep useful information within the latent features, thus improving the prediction accuracy. We have also shown the learning curves of different models in Fig. 9, which shows that RCoNet_4^4 converges slightly faster than the others, including COVID-Net, ReCoNet and CoroNet.

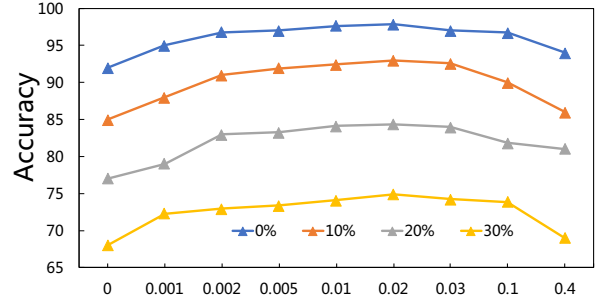


Fig. 8. The prediction accuracy by RCoNet_4^4 with regards to different values of α .

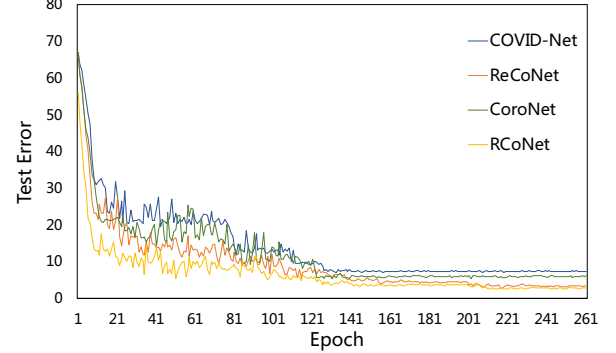


Fig. 9. Comparison on the learning trajectories of different models.

V. CONCLUSIONS

In this paper, we proposed a novel deep network model, named RCoNet_s^k , for robust COVID-19 detection, which contains three key components, i.e., *Deformable mutual Information Maximization* (DeIM), *Mixed High-order Moment Feature* (MHMF) and *Multi-expert Uncertainty-aware Learning* (MUL). DeIM estimates and maximizes the mutual information between input data and the latent representations simultaneously to obtain the category separability in the latent space. We proposed MHMF to overcome the limited expressive capability of low-order statistics, and instead use a combination of both low and high order moment features to extract more informative and discriminative features. MUL generates the final diagnosis and the uncertainty estimation, by combining the output of multiple parallel dropout networks, each as an expert. We numerically validated that the proposed RCoNet trained on either the public COVIDx dataset or the noisy version of it, outperforms the existing methods in terms of all the metrics considered. We note that these three modules can be easily implemented into other frameworks for different tasks.

REFERENCES

- [1] K. Zhang, X. Liu, J. Shen, Z. Li, Y. Sang, X. Wu, Y. Zha, W. Liang, C. Wang, K. Wang *et al.*, “Clinically applicable ai system for accurate diagnosis, quantitative measurements, and prognosis of covid-19 pneumonia using computed tomography,” *Cell*, 2020.
- [2] Z. Han, B. Wei, Y. Hong, T. Li, J. Cong, X. Zhu, H. Wei, and W. Zhang, “Accurate screening of covid-19 using attention based deep 3d multiple instance learning,” *IEEE Transactions on Medical Imaging*, 2020.
- [3] X. Mei, H.-C. Lee, K.-y. Diao, M. Huang, B. Lin, C. Liu, Z. Xie, Y. Ma, P. M. Robson, M. Chung *et al.*, “Artificial intelligence-enabled rapid diagnosis of patients with covid-19,” *Nature Medicine*, pp. 1–5, 2020.

- [4] W. Xie, C. Jacobs, J.-P. Charbonnier, and B. van Ginneken, "Relational modeling for robust and efficient pulmonary lobe segmentation in ct scans," *IEEE Transactions on Medical Imaging*, 2020.
- [5] X. Ouyang, J. Huo, L. Xia, F. Shan, J. Liu, Z. Mo, F. Yan, Z. Ding, Q. Yang, B. Song *et al.*, "Dual-sampling attention network for diagnosis of covid-19 from community acquired pneumonia," *IEEE Transactions on Medical Imaging*, 2020.
- [6] H. X. Bai, R. Wang, Z. Xiong, B. Hsieh, K. Chang, K. Halsey, T. M. L. Tran, J. W. Choi, D.-C. Wang, L.-B. Shi *et al.*, "Ai augmentation of radiologist performance in distinguishing covid-19 from pneumonia of other etiology on chest ct," *Radiology*, p. 201491, 2020.
- [7] A. A. Ardakani, A. R. Kanafi, U. R. Acharya, N. Khadem, and A. Mohammadi, "Application of deep learning technique to manage covid-19 in routine clinical practice using ct images: Results of 10 convolutional neural networks," *Computers in Biology and Medicine*, p. 103795, 2020.
- [8] H. Kang, L. Xia, F. Yan, Z. Wan, F. Shi, H. Yuan, H. Jiang, D. Wu, H. Sui, C. Zhang *et al.*, "Diagnosis of coronavirus disease 2019 (covid-19) with structured latent multi-view representation learning," *IEEE transactions on medical imaging*, 2020.
- [9] D.-P. Fan, T. Zhou, G.-P. Ji, Y. Zhou, G. Chen, H. Fu, J. Shen, and L. Shao, "Inf-net: Automatic covid-19 lung infection segmentation from ct images," *IEEE Transactions on Medical Imaging*, 2020.
- [10] L. Sun, Z. Mo, F. Yan, L. Xia, F. Shan, Z. Ding, W. Shao, F. Shi, H. Yuan, H. Jiang *et al.*, "Adaptive feature selection guided deep forest for covid-19 classification with chest ct," *arXiv preprint arXiv:2005.03264*, 2020.
- [11] D. Donglin, S. Feng, Y. Fuhua, X. Liming, M. Zhanhao, D. Zhongxiang, S. Fei, L. Shengrui, W. Ying, S. Ying, H. Miaofei, G. Yaozong, S. He, G. Yue, and S. Dinggang, "Hypergraph learning for identification of covid-19 with ct imaging," 2020.
- [12] Z. Y. Zu, M. D. Jiang, P. P. Xu, W. Chen, Q. Q. Ni, G. M. Lu, and L. J. Zhang, "Coronavirus disease 2019 (covid-19): a perspective from china," *Radiology*, p. 200490, 2020.
- [13] M. Siddhartha and A. Santra, "Covidlite: A depth-wise separable deep neural network with white balance and clahe for detection of covid-19," *arXiv preprint arXiv:2006.13873*, 2020.
- [14] L. Wang and A. Wong, "Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images," *arXiv preprint arXiv:2003.09871*, 2020.
- [15] Y. Oh, S. Park, and J. C. Ye, "Deep learning covid-19 features on cxr using limited training data sets," *IEEE Transactions on Medical Imaging*, 2020.
- [16] E. Pauwels and J. B. Lasserre, "Sorting out typicality with the inverse moment matrix sos polynomial," in *Advances in Neural Information Processing Systems*, 2016, pp. 190–198.
- [17] F. Maes, A. Collignon, D. Vandermeulen, G. Marchal, and P. Suetens, "Multimodality image registration by maximization of mutual information," *IEEE transactions on Medical Imaging*, vol. 16, no. 2, pp. 187–198, 1997.
- [18] A. Hyvärinen and E. Oja, "Independent component analysis: algorithms and applications," *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [19] N. Kwak and C.-H. Choi, "Input feature selection by mutual information based on parzen window," *IEEE transactions on pattern analysis and machine intelligence*, vol. 24, no. 12, pp. 1667–1671, 2002.
- [20] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 27, no. 8, pp. 1226–1238, 2005.
- [21] A. J. Butte and I. S. Kohane, "Mutual information relevance networks: functional genomic clustering using pairwise entropy measurements," in *Biocomputing 2000*. World Scientific, 1999, pp. 418–429.
- [22] C. Vondrick, A. Shrivastava, A. Fathi, S. Guadarrama, and K. Murphy, "Tracking emerges by colorizing videos," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 391–408.
- [23] R. Arandjelovic and A. Zisserman, "Look, listen and learn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 609–617.
- [24] M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm, "Mutual information neural estimation," in *International Conference on Machine Learning*, 2018, pp. 531–540.
- [25] J. Chang, Z. Lan, C. Cheng, and Y. Wei, "Data uncertainty learning in face recognition," *arXiv preprint arXiv:2003.11339*, 2020.
- [26] Y. Shi and A. K. Jain, "Probabilistic face embeddings," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6902–6911.
- [27] A. Kendall, V. Badrinarayanan, and R. Cipolla, "Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding," *arXiv preprint arXiv:1511.02680*, 2015.
- [28] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural networks," *arXiv preprint arXiv:1505.05424*, 2015.
- [29] Y. Gal, "Uncertainty in deep learning," *University of Cambridge*, vol. 1, p. 3, 2016.
- [30] D. J. MacKay, "A practical bayesian framework for backpropagation networks," *Neural computation*, vol. 4, no. 3, pp. 448–472, 1992.
- [31] R. M. Neal, *Bayesian learning for neural networks*. Springer Science & Business Media, 2012, vol. 118.
- [32] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *Advances in neural information processing systems*, 2017, pp. 5574–5584.
- [33] S. Isobe and S. Arai, "Deep convolutional encoder-decoder network with model uncertainty for semantic segmentation," in *2017 IEEE International Conference on Innovations in Intelligent Systems and Applications (INISTA)*. IEEE, 2017, pp. 365–370.
- [34] J. Choi, D. Chun, H. Kim, and H.-J. Lee, "Gaussian yolov3: An accurate and fast object detector using localization uncertainty for autonomous driving," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 502–511.
- [35] T. Yu, D. Li, Y. Yang, T. M. Hospedales, and T. Xiang, "Robust person re-identification by modelling feature uncertainty," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 552–561.
- [36] U. Zafar, M. Ghafoor, T. Zia, G. Ahmed, A. Latif, K. R. Malik, and A. M. Sharif, "Face recognition with bayesian convolutional networks for robust surveillance systems," *EURASIP Journal on Image and Video Processing*, vol. 2019, no. 1, p. 10, 2019.
- [37] M. D. Donsker and S. S. Varadhan, "Asymptotic evaluation of certain markov process expectations for large time, i," *Communications on Pure and Applied Mathematics*, vol. 28, no. 1, pp. 1–47, 1975.
- [38] S. Nowozin, B. Cseke, and R. Tomioka, "f-gan: Training generative neural samplers using variational divergence minimization," in *Advances in neural information processing systems*, 2016, pp. 271–279.
- [39] M. U. Gutmann and A. Hyvärinen, "Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics," *Journal of Machine Learning Research*, vol. 13, no. Feb, pp. 307–361, 2012.
- [40] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," in *Advances in Neural Information Processing Systems*, 2019, pp. 15 509–15 519.
- [41] J. Xu, P. Ye, Q. Li, H. Du, Y. Liu, and D. Doermann, "Blind image quality assessment based on high order statistics aggregation," *IEEE Transactions on Image Processing*, vol. 25, no. 9, pp. 4444–4457, 2016.
- [42] C. Chen, Z. Fu, Z. Chen, S. Jin, Z. Cheng, X. Jin, and X.-S. Hua, "Homm: Higher-order moment matching for unsupervised domain adaptation," *arXiv preprint arXiv:1912.11976*, 2019.
- [43] P. Jacob, D. Picard, A. Histace, and E. Klein, "Metric learning with horde: High-order regularizer for deep embeddings," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6539–6548.
- [44] H. Jégou and O. Chum, "Negative evidences and co-occurrences in image retrieval: The benefit of pca and whitening," in *European conference on computer vision*. Springer, 2012, pp. 774–787.
- [45] M. Opitiz, G. Waltner, H. Possegger, and H. Bischof, "Bier-boosting independent embeddings robustly," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5189–5198.
- [46] P. Kar and H. Karnick, "Random feature maps for dot product kernels," pp. 583–591, 2012.
- [47] S. Ahmed, M. H. Yap, M. Tan, and M. K. Hasan, "Reconet: Multi-level preprocessing of chest x-rays for covid-19 detection using convolutional neural networks," *medRxiv*, 2020.
- [48] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [49] A. I. Khan, J. L. Shah, and M. M. Bhat, "Coronet: A deep neural network for detection and diagnosis of covid-19 from chest x-ray images," *Computer Methods and Programs in Biomedicine*, p. 105581, 2020.
- [50] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," pp. 5987–5995, 2017.
- [51] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," in *International conference on machine learning*, 2014, pp. 647–655.