

# CMC\_v2: Towards More Accurate COVID-19 Detection with Discriminative Video Priors

Junlin Hou<sup>1</sup>, Jilan Xu<sup>1,3</sup>, Nan Zhang<sup>2</sup>, Yi Wang<sup>3</sup>, Yuejie Zhang<sup>1\*</sup>,  
Xiaobo Zhang<sup>4\*</sup>, and Rui Feng<sup>1,2,4\*</sup>

<sup>1</sup> School of Computer Science, Shanghai Key Laboratory of Intelligent Information Processing, Fudan University, China

<sup>2</sup> Academy for Engineering and Technology, Fudan University, China  
{jlhou18, jilanxu18, 20210860062, yjzhang, fengrui}@fudan.edu.cn

<sup>3</sup> Shanghai AI Laboratory, China  
wygamle@gmail.com

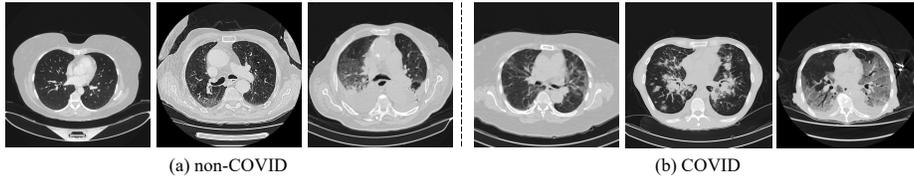
<sup>4</sup> Children's Hospital of Fudan University, National Children's Medical Center, Shanghai, China  
zhangxiaobo0307@163.com

**Abstract.** This paper presents our solution for the 2nd COVID-19 Competition, occurring in the framework of the AIMIA Workshop at the European Conference on Computer Vision (ECCV 2022). In our approach, we employ the winning solution last year which uses a strong 3D Contrastive Mixup Classification network (CMC\_v1) as the baseline method, composed of contrastive representation learning and mixup classification. In this paper, we propose CMC\_v2 by introducing natural video priors to COVID-19 diagnosis. Specifically, we adapt a pre-trained (on video dataset) video transformer backbone to COVID-19 detection. Moreover, advanced training strategies, including hybrid mixup and cutmix, slice-level augmentation, and small resolution training are also utilized to boost the robustness and the generalization ability of the model. Among 14 participating teams, CMC\_v2 ranked 1st in the 2nd COVID-19 Competition with an average Macro F1 Score of 89.11%.

**Keywords:** COVID-19 detection, Hybrid CNN-transformer, Contrastive learning, Hybrid mixup and cutmix

## 1 Introduction

The Coronavirus Disease 2019 SARS-CoV-2 (COVID-19), identified at the end of 2019, is a highly infectious disease, leading to an everlasting worldwide pandemic and collateral economic damage [29]. Early detection of COVID-19 is crucial to the timely treatment of patients, and beneficial to slowdown or even break viral transmission. COVID-19 detection aims to identify COVID from non-COVID cases. Among several COVID-19 detection means, chest computed tomography (CT) has been recognized as a key component in the diagnostic procedure for COVID-19. In CT, we resort to typical radiological findings to confirm COVID-19, including ground glass opacities, opacities with rounded morphology,



**Fig. 1.** Some examples of (a) non-COVID and (b) COVID cases from the COV19-CT-DB dataset. The non-COVID category includes no pneumonia and other pneumonia cases. The COVID category contains COVID-19 cases of different severity levels.

crazy-paving pattern, and consolidations [3]. As a CT volume contains hundreds of slices, delivering a convincing diagnosis from these data demands a heavy workload on radiologists. Relying on manual analysis is barely scalable considering the surging increasing number of infection cases. Regarding this, there is an urgent need for accurate automated COVID-19 diagnosis approaches.

Recently, deep learning approaches have achieved promising performance in fighting against COVID-19. They have been widely applied to various medical practices, including the lung and infection region segmentation [27,18,2,10] as well as the clinical diagnosis and assessment [28,25,22,9]. Though a line of works [10,28,9] has been employed for COVID-19 detection via CT analysis and yielded effective results, it is still worth pushing its detection performance to a new level in a faster and more accurate manner for a better medical assistant experience. Improving this performance is non-trivial, since the inner variances between CT scans of COVID are huge and its differences with some non-COVID like pneumonia are easily overlooked. Specifically, CT scans vary greatly in imaging across different devices and hospitals (Fig. 1), and they share several similar visual manifestations with other types of pneumonia. Further, the scarcity of CT scans of COVID-19 due to regulations in the medical area makes these challenges even harder, as we cannot simply turn to a deep model to learn these mentioned characteristics with a big number of annotated scans from scratch.

To tackle these challenges, we exploit video priors along with the given limited number of CT scans to learn an effective feature space for COVID-19 detection, along with contrastive training and some hybrid data augmentation means for further data-efficient learning. Specifically, we employ the advanced 3D contrastive mixup classification network (CMC-COV19D, abbr. CMC\_v1) [8], the winner in the ICCV 2021 COVID-19 Diagnosis Competition of AI-enabled Medical Image Analysis Workshop [13], as a baseline. CMC\_v1 introduces contrastive representation learning to discover discriminative representations of COVID-19 cases. Besides, a joint training loss is devised by combining the classification loss, mixup loss, and contrastive loss. In this work, we propose CMC\_v2 by introducing the following mechanisms customized for 3D models. (1) To capture the long-range lesion span across the slices in the CT scans, we adopt a hybrid CNN-transformer model, i.e. Uniformer [17] as the backbone network. The combination of convolution and self-attention reduces the network parameters

and computational costs. It relieves the potential overfitting when deploying 3D models to small-scale medical datasets. Besides, we empirically show that initializing the model with 3D weights pre-trained on video datasets is promising as modeling the relationship among slices is critical for COVID-19 detection. (2) We develop a hybrid mixup and cutmix augmentation strategy to enhance the models’ generalization ability. Due to the limited memory, a gather-and-dispatch mechanism is also customized for the modern Distributed DataParallel (DDP) scheme in Multi-GPU training. (3) We showcase both the 2D slice-level augmentation and the small resolution training bring improvements. By applying the intra-and-inter model ensemble [8], CMC\_v2 won the first prize in the 2nd COVID-19 detection challenge of the Workshop “AI-enabled Medical Image Analysis – Digital Pathology & Radiology/COVID19 (AIMIA)”. CMC\_v2 significantly outperforms the baseline model provided by the organizers by 16% Macro F1 Score.

The remainder of this paper is organized as follows. Section 2 reviews related works. In Section 3, we first recap the CMC\_v1 network, the basis of CMC\_v2, and then introduce the newly proposed modules in CMC\_v2. Section 4 describes the COV19-CT-DB dataset used in this paper. Section 5 provides the experimental settings and results. Section 6 concludes our work.

## 2 Related Work

### 2.1 COVID-19 detection

Numerous deep learning approaches have made great efforts to separate COVID patients from non-COVID subjects. Despite the binary classification, the task is challenging as the non-COVID cases include both common pneumonia subjects and non-pneumonia subjects.

The majority of deep learning approaches are based on Convolutional Neural Networks (CNN). [25] was a pioneering work that designed a CNN model to classify COVID-19 and typical viral pneumonia. Song et al. [22] proposed a deep learning-based CT diagnosis system (Deep Pneumonia) to detect patients with COVID-19 from patients with bacteria pneumonia and healthy people. Li et al. [18] developed a 3D COVNet based on ResNet50, aiming to extract both 2D local and 3D global features to classify COVID-19, CAP, and non-pneumonia. Xu et al. [31] introduced a location-attention model to categorize COVID-19, Influenza-A viral pneumonia, and healthy cases. It took the relative distance-from-edge of segmented lesion candidates as extra weight in a fully connected layer to offer distance information.

Recently, Vision Transformer (ViT) has demonstrated its potentials by achieving competitive results on a variety of computer vision tasks. Relevant studies have also been conducted on the COVID-19 diagnosis. Gao et al. [6] used a ViT based on the attention models to classify COVID and non-COVID CT images. To integrate the advantages of convolution and transformer for COVID-19 detection, Park et al. [20] presented a novel architecture that utilized CNN as a feature

extractor for low-level Chest X-ray feature corpus, upon which Transformer was trained for downstream diagnosis tasks with the self-attention mechanism.

## 2.2 Advanced network architecture

In our approach, we adopt two representative deep learning architectures as the backbones, namely ResNeSt-50 and Uniformer-S. Here, we briefly introduce the closely related ResNet and Transformer architectures and their variants.

In the family of ResNets, ResNet [7] introduced a deep residual learning framework to address the network degradation problem. ResNeXt [30] established a simple architecture by adopting group convolution in the ResNet bottleneck block. ResNeSt [33] presented a modular split-attention block within the individual network blocks to enable attention across feature-map groups.

Although CNN models have shown promising results, the limited receptive field makes it hard to capture global dependency. To solve this problem, Vision Transformer (ViT) [4] was applied to the sequences of image patches for an image classification task. Later on, Swin Transformer [19] proposed to use shifted windows between consecutive self-attention layers, which had the flexibility to model at various scales and had linear computational complexity with respect to the image size. Multi-scale Vision Transformer [5] connected the seminal idea of multi-scale feature hierarchies with transformer models for video and image recognition. Pyramid Vision Transformer [26] used a progressive shrinking pyramid to reduce the computations of large feature maps, which overcame the difficulties of porting Transformer models to various dense prediction tasks and inherits the advantages of both CNN and Transformer. Unified transformer (Uniformer) [17] sought to integrate the merits of convolution and self-attention in a concise transformer format, which can tackle both local redundancy and global dependency. To achieve the balance between accuracy and efficiency, we adopt Uniformer as the default backbone network.

## 3 Methodology

The overall framework of our model is shown in Fig. 2. In this section, we review the baseline method CMC\_v1 [8] firstly and then introduce several simple and effective mechanisms to boost the detection performance.

### 3.1 Recap of CMC\_v1

CMC\_v1 employs the contrastive representation learning (CRL) as an auxiliary task to learn discriminative representations of COVID-19. CRL is comprised of the following components. 1) A stochastic data augmentation module  $A(\cdot)$ , which transforms an input CT  $x_i$  into a randomly augmented sample  $\tilde{x}_i$ . Two augmented volumes are generated from each input CT scan. 2) A base encoder  $E(\cdot)$ , mapping the augmented CT sample  $\tilde{x}_i$  to its feature representation  $r_i = E(\tilde{x}_i) \in \mathbb{R}^{d_e}$ . 3) A projection network  $P(\cdot)$ , used to map the representation

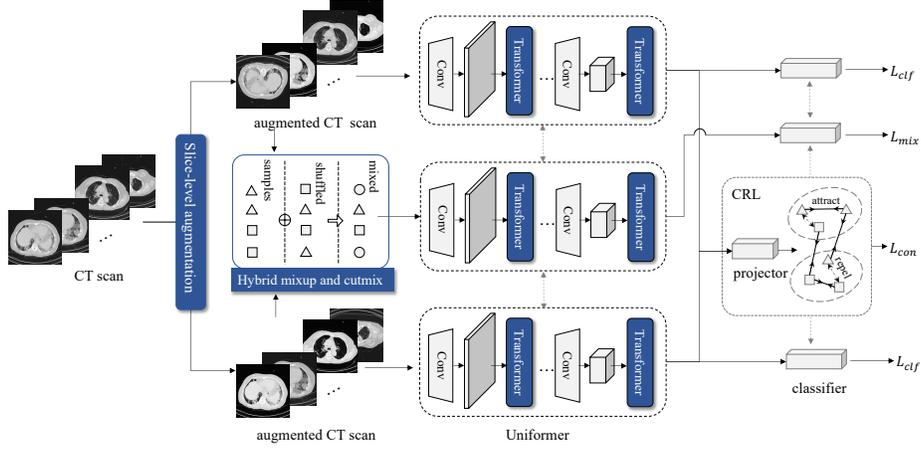


Fig. 2. Overview of our CMC\_v2 network for COVID-19 detection.

vector  $r_i$  to a relative low-dimension vector  $z_i = P(r_i) \in \mathbb{R}^{d_p}$ . 4) A classifier  $C(\cdot)$ , classifying the vector  $r_i \in \mathbb{R}^{d_e}$  to the final prediction.

**Contrastive representation learning.** Given a minibatch of  $N$  CT volumes and their labels  $\{(x_i, y_i)\}$ , we can generate a minibatch of  $2N$  samples  $\{(\tilde{x}_i, \tilde{y}_i)\}$  after data augmentations. Inspired by the supervised contrastive loss [11], we define the positives as any augmented CT samples from the same category, whereas the CT samples from different classes are considered as negative pairs. Let  $i \in \{1, \dots, 2N\}$  be the index of an arbitrary augmented sample, the contrastive loss function is defined as:

$$\mathcal{L}_{con}^i = \frac{-1}{2N_{\tilde{y}_i} - 1} \sum_{j=1}^{2N} \mathbb{1}_{i \neq j} \cdot \mathbb{1}_{\tilde{y}_i = \tilde{y}_j} \cdot \log \frac{\exp(z_i^T \cdot z_j / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{i \neq k} \cdot \exp(z_i^T \cdot z_k / \tau)}, \quad (1)$$

where  $\mathbb{1} \in \{0, 1\}$  is an indicator function, and  $\tau > 0$  denotes a scalar temperature hyper-parameter.  $N_{\tilde{y}_i}$  is the total number of samples in a minibatch that have the same label  $\tilde{y}_i$ .

**Mixup classification.** CMC\_v1 adopts the mixup [34] strategy during training to further boost the generalization ability of the model. For each augmented CT sample  $\tilde{x}_i$ , the mixup sample and its label are generated as:

$$\tilde{x}_i^{mix} = \lambda \tilde{x}_i + (1 - \lambda) \tilde{x}_p, \quad \tilde{y}_i^{mix} = \lambda \tilde{y}_i + (1 - \lambda) \tilde{y}_p, \quad (2)$$

where  $p$  is randomly selected indice;  $\lambda$  is the balancing coefficient. The mixup loss is defined as the cross-entropy loss of mixup samples:

$$\mathcal{L}_{mix}^i = \text{CrossEntropy}(\tilde{x}_i^{mix}, \tilde{y}_i^{mix}). \quad (3)$$

Different from the original design [34] where they replaced the classification loss with the mixup loss, we merge the mixup loss with the standard cross-entropy classification loss  $\mathcal{L}_{clf}^i = \text{CrossEntropy}(\tilde{x}_i, \tilde{y}_i)$  to enhance the classification ability on both mixup samples and raw samples.

The total loss is defined as the combination of the contrastive loss, mixup loss, and classification loss:

$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^{2N} (\mathcal{L}_{con}^i + \mathcal{L}_{mix}^i + \mathcal{L}_{clf}^i). \quad (4)$$

### 3.2 Improving COVID-19 detection with CMC\_v2

To boost the COVID-19 detection performance, we incorporate natural video priors into CMC\_v1 by adapting an efficient pre-trained video backbone to our task, and develop a hybrid data augmentation strategy to increase data efficiency.

**Transfer learning with a stronger backbone and pre-training.** In CMC\_v1, a 3D ResNeSt-50 model [33] is employed as the backbone network for feature extraction. Although 3D convnets capture local volume semantics efficiently, they are incapable of modeling long-range dependencies between spatial/temporal features explicitly. For simplicity, we refer ‘temporal’ to the relationship among different CT slices in this paper. Recent works on Vision Transformer [4] managed to encode long-range information using self-attention. However, global self-attention is computationally inefficient and transformer models only demonstrate superior results when huge data is available. Compared with natural image datasets, COVID-19 image datasets have a smaller scale and the model is prone to overfitting. To alleviate this issue, we adopt a video transformer named Uniformer [17], a novel hybrid CNN-transformer model which integrates the advantages of convolution and self-attention in spatial-temporal feature learning while achieving the balance between accuracy and efficiency. In particular, Uniformer replaces the naive transformer block with a Uniformer block, which is comprised of a Dynamic Position Embedding (DPE) layer, a Multi-Head Relation Aggregator (MHRA) layer, and a Feed-Forward Network (FFN).

Furthermore, we experimentally find that training the model from scratch leads to poor results. In transfer learning, it is a common practice to initialize the model on downstream tasks with weights pre-trained on a large-scale ImageNet dataset. To initialize the 3D model, CMC\_v1 inflated the ImageNet pre-trained 2D weights to the 3D model. This is achieved by either copying the 2D weights to the center of the 3D weights or repeating the 2D weights along the third dimension. However, these inflated 3D weights may not excel at modeling the temporal relationship between different slices. To address this issue, we directly initialize the model with 3D weights pre-trained on video action recognition datasets, i.e. k400 [1]. We empirically prove that k400 pre-training yields better results than inflated weight initialization in this task.

**Hybrid mixup and cutmix strategy.** In CMC\_v1, the mixup strategy is introduced to generate diverse CT samples. These pseudo samples are beneficial for improving the model’s generalization ability. Similar to a mixup, cutmix [32] replaces a local region in the target image with the corresponding local region sampled in the source image. To combine the merits of both, we develop a hybrid mixup and cutmix strategy. In each iteration, we select one strategy with equal probability. This hybrid strategy works well on the traditional Data Parallel (DP) mechanism [21] in multi-GPU training. However, it’s challenging to scale to the modern Distributed Data Parallel (DDP) mechanism. The original batch size on each GPU is set to 1 in our case because the effective batch size is 4 after two-view augmentation and hybrid mixup and cutmix strategy, reaching the memory limit on each GPU (The shape of the mini-batch tensors on each GPU is  $4 \times T \times 3 \times H \times W$ ). As the DDP mechanism starts an individual process on each GPU, the hybrid strategy is directly employed on each GPU individually. Performing the hybrid mixup and cutmix strategy on the augmented views of the same image does not align with the original effect. To make it work, we gather all the samples from the GPUs, conduct the hybrid mixup and cutmix over all the samples, and dispatch the generated samples back to each GPU. It guarantees that the hybrid strategy is performed across different CT scans in the current mini-batch. The hybrid mixup and cutmix strategy boost the model’s generalization ability.

**Slice-level augmentation.** The data augmentation strategies used in CMC\_v1 are 3D rescaling, 3D rotation, and color jittering on all the slices. To further increase the data diversity, we follow the common practice in video data processing and perform different 2D augmentations on each slice, termed as SliceAug. SliceAug achieves slightly better performance than 3D augmentation while having a comparable pre-processing time.

**Small resolution training.** Prior works [4,23] have demonstrated the effectiveness of using small image resolution during training and large resolution during validation/testing. This mechanism bridges the gap between the image size mismatch caused by the random resized cropping during training and center cropping during testing [24]. Besides, the small resolution makes training more efficient. In the experiments, we use the resolution of  $192 \times 192$  and  $224 \times 224$  for training and testing, respectively.

## 4 Dataset

We evaluate our proposed approach on the COV19-CT-Database (COV19-CT-DB) [12]. The COV19-CT-DB contains chest CT scans marking the existence of COVID-19. It consists of about 1,650 COVID and 6,100 non-COVID chest CT scan series from over 1,150 patients and 2,600 subjects. In total, 724,273 slices correspond to the CT scans of the COVID category and 1,775,727 slices

correspond to the non-COVID category. Data collection was conducted in the period from September 1, 2020 to November 30, 2021. Annotation of each CT scan was obtained by 4 experienced medical experts and showed a high degree of agreement (around 98%). Each 3D CT scan includes a different number of slices, ranging from 50 to 700. This variation in the number of slices is due to the context of CT scanning. The database is split into training, validation, and testing sets. The training set contains 1,992 3D CT scans (1,110 non-COVID cases and 882 COVID cases). The validation set consists of 504 3D CT scans (289 non-COVID cases and 215 COVID cases). The testing set includes 5,281 scans and the labels are not available during the challenge.

## 5 Experiments

### 5.1 Implementation details

All CT volumes are resized from  $(T, 512, 512)$  to  $(128, 224, 224)$ , where  $T$  denotes the number of slices. For training, data augmentations include random resized cropping on the transverse plane, random cropping on the vertical section to 64, rotation, and color jittering. We employ the 3D ResNeSt-50 and Uniformer-S as the backbones in our experiments. The value of parameter  $d_e$  is  $2,048/512$  for ResNeSt-50/Uniformer-S, and  $d_p$  is set to 128. All networks are optimized using the Adam algorithm with a weight decay of  $1e-5$ . The initial learning rate is set to  $1e-4$  and then divided by 10 at 30% and 80% of the total number of training epochs. The networks are trained for 100 epochs. Our methods are implemented in PyTorch and run on eight NVIDIA Tesla A100 GPUs.

### 5.2 Evaluation metrics

To evaluate the performance of the proposed method, we adopt the same official protocol of 2nd COVID-19 Competition as the evaluation metric. We report F1 Scores for non-COVID and COVID categories as well as the Macro F1 Score for overall comparison. The Macro F1 Score is defined as the unweighted average of the class-wise/label-wise F1 Scores. We also present ROC curves and Area Under Curve (AUC) for each category.

### 5.3 Ablation studies on COVID-19 detection challenge

We conduct ablation studies on the validation set of COVID-19 detection challenge to show the impact of each component of our proposed methods. We first analyze the effects of different backbones, and then we discuss the effectiveness of the CMC\_v1 framework and the choice of various pre-training methods. Finally, we investigate the impact of the new components in our CMC\_v2, i.e. slice-level augmentation (SliceAug), hybrid mixup and cutmix strategy (Hybrid), and small resolution training (SmallRes).

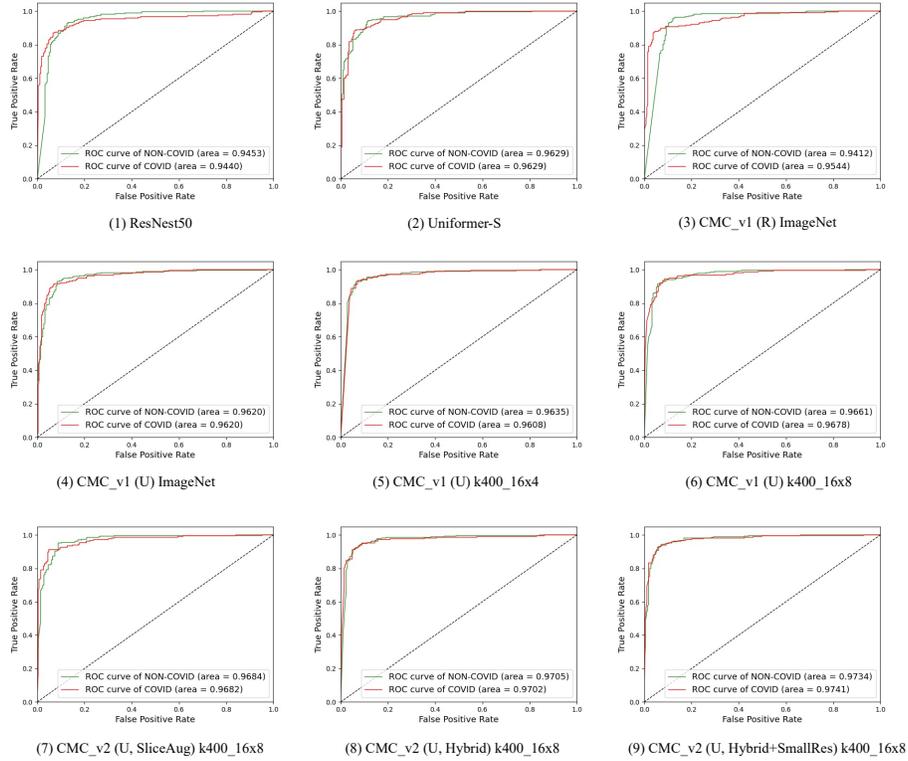
**Table 1.** The results on the validation set of COVID-19 detection challenge.

ID	Method	Param	FLOPs	Pre-train	Macro F1	F1	
						Non-COVID	COVID
1	ResNet50-GRU [12]	-	-	-	77.00	-	-
2	ResNeSt-50	52.8M	371.9G	ImageNet	89.89	91.27	88.52
3	Uniformer-S	21.2M	230.1G	ImageNet	90.98	92.08	89.88
4	CMC_v1 (R)	57.3M	371.9G	ImageNet	91.98	93.14	90.82
5	CMC_v1 (U)	21.5M	230.1G	ImageNet	92.26	93.11	91.42
6	CMC_v1 (U)	21.5M	230.1G	k400_16×4	92.48	93.28	91.67
7	CMC_v1 (U)	21.5M	230.1G	k400_16×8	92.70	93.41	91.99
8	CMC_v2 (U, SliceAug)	21.5M	230.1G	k400_16×8	93.07	93.94	92.20
9	CMC_v2 (U, Hybrid)	21.5M	230.1G	k400_16×8	93.29	94.12	92.45
10	CMC_v2 (U, Hybrid+SmallRes)	21.5M	169.1G	k400_16×8	93.30	94.07	92.52

**Backbone network.** To analyze the effects of architectures, we compare different backbone models, and the results are shown in the first three rows of Table 1. The reported result of the baseline approach ‘ResNet50-GRU’ [12] is 77.00% Macro F1 Score. This model is based on CNN-RNN architecture [14,16,15], where the CNN part performs local analysis on each 2D slice, and the RNN part combines the CNN features of the whole 3D CT scan. Compared to the baseline, our 3D ResNeSt-50 and Uniformer-S backbones achieve more than 12% improvements on the Macro F1 Scores. Specifically, the Uniformer-S achieves better performance on all the metrics, surpassing ResNeSt-50 by 1.09% Macro F1 Score, 0.81% and 1.36% F1 Scores for non-COVID and COVID classes. Besides, the Uniformer-S greatly reduces the network parameters and computational costs. The results demonstrate the long-range dependencies modeling ability of Uniformer-S, which is important to capture the relationships between different CT slices.

**Analysis of CMC\_v1.** We evaluate the effectiveness of the previous CMC\_v1 network. The 4th and 5th rows in Table 1 show the results of CMC\_v1 (R) and CMC\_v1 (U), where the R and U denote ResNeSt-50 and Uniformer-S backbones, respectively. CMC\_v1 on both backbones can achieve significant performance improvements. In particular, CMC\_v1 (U) obtains 92.26% on Macro F1 Score, 93.11% and 91.42% on F1 Scores for non-COVID and COVID categories. The results demonstrate the generality of the CMC\_v1, which can consistently improve the COVID-19 detection performance with different backbones.

**Pre-training schemes.** We compare three pre-training methods, namely ImageNet, k400\_16×4, and k400\_16×8. ImageNet pre-training inflates the 2D pre-trained weights to our 3D models. K400 pre-training denotes 3D weights pre-trained on the video action recognition dataset k400, where 16×4 and 16×8 indicate the sampling 16 frames with frame stride 4 and 8, respectively. It can



**Fig. 3.** The ROC curves and AUC scores of different networks.

be seen from the 5th to 7th rows in Table 1, CMC\_v1 (U) with k400\_16 $\times$ 8 pre-training weights outperforms the other two methods on all metrics. Based on the above results, we choose the Uniformer-S with k400\_16 $\times$ 8 pre-training weights as the default backbone for our proposed CMC\_v2.

**Analysis of CMC\_v2.** In this part, we investigate the impact of our newly proposed components in CMC\_v2, including slice-level augmentation (SliceAug), hybrid mixup and cutmix strategy (Hybrid), and small resolution training (Small-Res). The experimental results in the 8th row of Table 1 indicate that CMC\_v2 (U, SliceAug) can improve the performance on all metrics compared with the CMC\_v1 (U). The slice-level augmentation can further increase the data diversity and benefit COVID-19 detection performance. As for the hybrid mixup and cutmix strategy, the CMC\_v2 (U, Hybrid) achieves further improvement by 0.59% Macro F1 Score, 0.71% COVID F1 Score, and 0.46% non-COVID F1 Score, compared with the CMC\_v1 (U) that only employs the single mixup strategy. Our hybrid mixup and cutmix strategy generates diversified data for improving the model’s generalization ability in COVID-19 detection. When we adopt the small

**Table 2.** The leaderboard on the 2nd COVID-19 detection challenge.

Rank	Teams	Macro F1	F1	
			Non-COVID	COVID
1	FDVTS (Ours)	89.11	97.31	80.92
1	ACVLab	89.11	97.45	80.78
3	MDAP	87.87	96.95	78.80
4	Code 1055	86.18	96.37	76.00
5	CNR-IEMN	84.37	95.98	72.76
6	Dslab	83.78	96.22	71.33
7	Jovision-Deepcam	80.82	94.56	67.07
8	ICL	79.55	93.77	65.34
9	etro	78.72	93.48	63.95
10	ResNet50-GRU [12]	69.00	83.62	54.38

resolution training mechanism, the CMC\_v2 (U, Hybrid+SmallRes) achieves the best performance with minimal computational costs among all models. It obtains 93.30% on Macro F1 Score, 94.07% on non-COVID F1 Score, and 92.52% on COVID F1 Score. In particular, this model shows the superior recognition ability for the COVID-19 category among all other approaches.

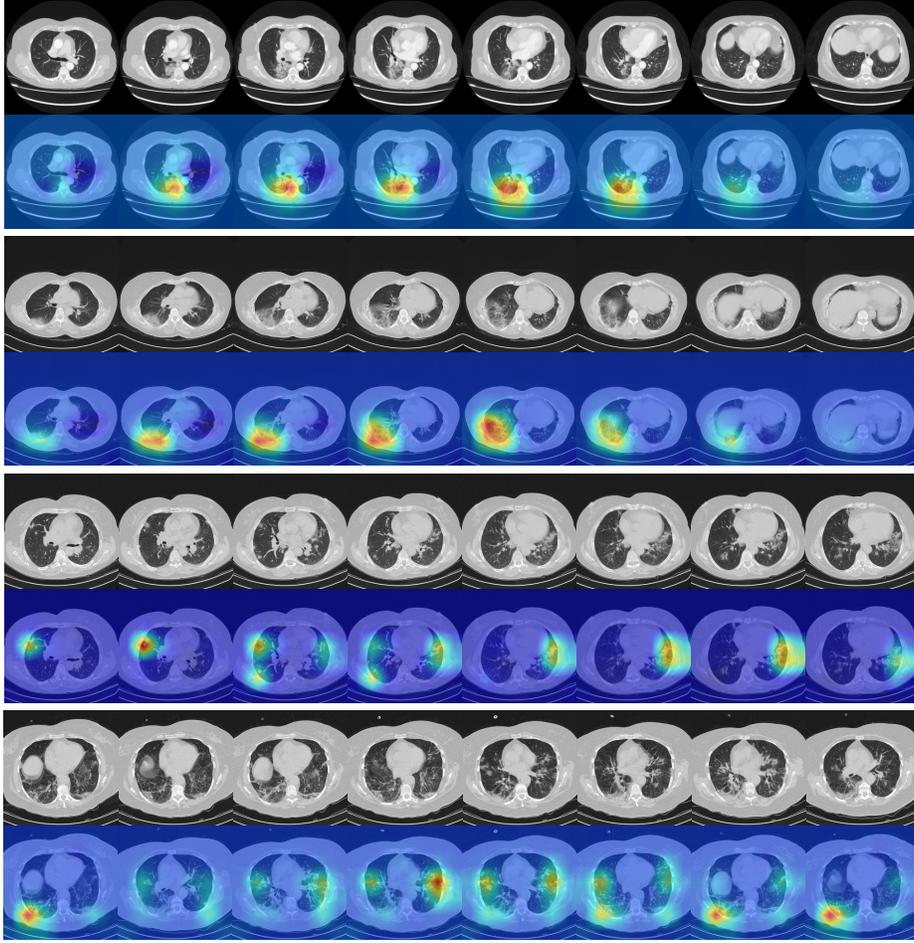
In addition, we present the ROC curves and AUC of our models in Fig. 3. The AUC results of all the models reach more than 0.94 for both non-COVID and COVID classes. Especially, the full version of CMC\_v2 (U, Hybrid+SmallRes) obtains the highest AUC Scores (0.9734 and 0.9741 for non-COVID and COVID, respectively) among all settings.

#### 5.4 Results on COVID-19 detection challenge leaderboard

Table 2 shows the results of our method and other participants on the testing set of 2nd COVID-19 detection challenge. Our method ensembles all the CMC\_v2, including CMC\_v2 (U, SliceAug), CMC\_v2 (U, Hybrid), and CMC\_v2 (U, Hybrid+SmallRes) following the strategy in [8]. The final prediction of each CT scan is obtained by averaging the predictions from individual models. We also adopt a test time augmentation (TTA) operation to boost the generalization ability of our models on the testing set. It can be seen from Table 2 that our proposed method ranks first in the challenge with 89.11% Macro F1 Score. Compared to other methods, our model achieves significant improvement on the F1 Score for the COVID category (80.92%), indicating the ability to distinguish COVID cases from non-pneumonia and other types of pneumonia correctly.

#### 5.5 Visualization results

To verify the interpretability of our model, we visualize the results using Class Activation Mapping (CAM) [35]. As illustrated in Fig. 4, we select four COVID-19 CT scans from the validation set of COV19-CT-DB dataset. In each group,



**Fig. 4.** The visualization results on the COVID-19 CT scans.

the upper row shows the series of CT slices, and the lower row presents the corresponding CAM results. In the first group, it can be seen that the attention maps focus on the local infection regions accurately. In the second group, the wide range of infection regions can also be covered. In the third and fourth groups, the infections in bilateral lungs can also be located precisely. These attention maps provide convincing interpretability for the COVID-19 detection results, which is helpful for real-world clinical diagnosis.

## 6 Conclusions

In this paper, we propose a novel and practical solution winning COVID-19 detection at the 2nd COVID-19 Competition. Based on the CMC\_v1 network, we

further develop the CMC\_v2 network with substantial improvements, including the CNN-transformer video backbone, hybrid mixup and cutmix strategy, slice-level augmentation, and small resolution training mechanism. The experimental results demonstrate that the new components boost the COVID-19 detection performance and the generalization ability of the model. On the testing set, our method ranked 1st in the 2nd COVID-19 Competition with 89.11% Macro F1 Score among 14 participating teams.

## Acknowledgement

This work was supported by the Scientific & Technological Innovation 2030 - “New Generation AI” Key Project (No. 2021ZD0114001; No. 2021ZD0114000), and the Science and Technology Commission of Shanghai Municipality (No. 21511104502; No. 21511100500; No. 20DZ1100205). Yuejie Zhang, Xiaobo Zhang, and Rui Feng are corresponding authors.

## References

1. Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6299–6308 (2017)
2. Chen, J., Wu, L., Zhang, J., Zhang, L., Gong, D., Zhao, Y., et al.: Deep learning-based model for detecting 2019 novel coronavirus pneumonia on high-resolution computed tomography. *Scientific Reports* **10**(1), 1–11 (2020)
3. Chung, M., Bernheim, A., Mei, X., Zhang, N., Huang, M., Zeng, X., Cui, J., Xu, W., Yang, Y., Fayad, Z.A., et al.: Ct imaging features of 2019 novel coronavirus (2019-ncov). *Radiology* (2020)
4. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
5. Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C.: Multiscale vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6824–6835 (2021)
6. Gao, X., Qian, Y., Gao, A.: Covid-vit: Classification of covid-19 from ct chest images based on vision transformer models. arXiv preprint arXiv:2107.01682 (2021)
7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)
8. Hou, J., Xu, J., Feng, R., Zhang, Y., Shan, F., Shi, W.: Cmc-cov19d: Contrastive mixup classification for covid-19 diagnosis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 454–461 (2021)
9. Hou, J., Xu, J., Jiang, L., Du, S., Feng, R., Zhang, Y., Shan, F., Xue, X.: Periphery-aware covid-19 diagnosis with contrastive representation enhancement. *Pattern Recognition* **118**, 108005 (2021)
10. Jin, S., Wang, B., Xu, H., Luo, C., Wei, L., Zhao, W., et al.: Ai-assisted ct imaging analysis for covid-19 screening: building and deploying a medical ai system in four weeks. *MedRxiv* (2020)

11. Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D.: Supervised contrastive learning. In: Annual Conference on Neural Information Processing Systems 2020 (2020)
12. Kollias, D., Arsenos, A., Kollias, S.: Ai-mia: Covid-19 detection & severity analysis through medical imaging. arXiv preprint arXiv:2206.04732 (2022)
13. Kollias, D., Arsenos, A., Soukissian, L., Kollias, S.: Mia-cov19d: Covid-19 detection through 3-d chest ct image analysis. arXiv preprint arXiv:2106.07524 (2021)
14. Kollias, D., Bouas, N., Vlaxos, Y., Brillakis, V., Seferis, M., Kollia, I., Sukissian, L., Wingate, J., Kollias, S.: Deep transparent prediction through latent representation analysis. arXiv preprint arXiv:2009.07044 (2020)
15. Kollias, D., Tagaris, A., Stafylopatis, A., Kollias, S., Tagaris, G.: Deep neural architectures for prediction in healthcare. *Complex & Intelligent Systems* **4**(2), 119–131 (2018)
16. Kollias, D., Vlaxos, Y., Seferis, M., Kollia, I., Sukissian, L., Wingate, J., Kollias, S.D.: Transparent adaptation in deep medical image diagnosis. In: TAILOR. pp. 251–267 (2020)
17. Li, K., Wang, Y., Zhang, J., Gao, P., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unifying convolution and self-attention for visual recognition. arXiv preprint arXiv:2201.09450 (2022)
18. Li, L., Qin, L., Xu, Z., Yin, Y., Wang, X., Kong, B., et al.: Artificial intelligence distinguishes covid-19 from community acquired pneumonia on chest ct. *Radiology* **296**, 200905 (2020)
19. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
20. Park, S., Kim, G., Oh, Y., Seo, J.B., Lee, S.M., Kim, J.H., Moon, S., Lim, J.K., Ye, J.C.: Vision transformer using low-level chest x-ray feature corpus for covid-19 diagnosis and severity quantification. arXiv preprint arXiv:2104.07235 (2021)
21. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019)
22. Song, Y., Zheng, S., Li, L., Zhang, X., Zhang, X., Huang, Z., et al.: Deep learning enables accurate diagnosis of novel coronavirus (covid-19) with ct images. *MedRxiv* (2020)
23. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. pp. 10347–10357. PMLR (2021)
24. Touvron, H., Vedaldi, A., Douze, M., Jégou, H.: Fixing the train-test resolution discrepancy. *Advances in neural information processing systems* **32** (2019)
25. Wang, S., Kang, B., Ma, J., Zeng, X., Xiao, M., Guo, J., et al.: A deep learning algorithm using ct images to screen for corona virus disease (covid-19). *European radiology* pp. 1–9 (2021)
26. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 568–578 (2021)
27. Wang, X., Deng, X., Fu, Q., Zhou, Q., Feng, J., Ma, H., et al.: A weakly-supervised framework for covid-19 classification and lesion localization from chest ct. *IEEE Transactions on Medical Imaging* **39**(8), 2615–2625 (2020)

28. Wang, Z., Xiao, Y., Li, Y., Zhang, J., Lu, F., Hou, M., et al.: Automatically discriminating and localizing covid-19 from community-acquired pneumonia on chest x-rays. *Pattern Recognition* **110**, 107613 (2020)
29. WHO: Coronavirus disease (covid-19) pandemic. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019> (2022)
30. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1492–1500 (2017)
31. Xu, X., Jiang, X., Ma, C., Du, P., Li, X., Lv, S., et al.: A deep learning system to screen novel coronavirus disease 2019 pneumonia. *Engineering* **6**(10), 1122–1129 (2020)
32. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6023–6032 (2019)
33. Zhang, H., Wu, C., Zhang, Z., Zhu, Y., Lin, H., Zhang, Z., Sun, Y., He, T., Mueller, J., Manmatha, R., et al.: Resnest: Split-attention networks. arXiv preprint arXiv:2004.08955 (2020)
34. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412 (2017)
35. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2921–2929 (2016)