

## 多分支特征融合分类网络用于CXR图像识别

苏华强<sup>1</sup> 雷海军<sup>1</sup> 雷柏英<sup>\*2</sup>

(1. 深圳大学计算机与软件学院, 广东省普及型高性能计算机重点实验室, 广东深圳 518060; 2. 深圳大学医学部生物医学工程学院, 广东省生物医学信息检测和超声成像重点实验室, 广东深圳 518060)

**摘要:** COVID-19是由新型冠状病毒引起的一种传染性疾病,给全球公共卫生带来了巨大的挑战。在临床实践中,胸部X射线(Chest X-ray, CXR)检查是识别COVID-19感染和其他常见肺部疾病的重要手段,然而放射科医生对COVID-19患者进行检查需要耗费大量时间和精力,而且增加医生感染的风险。因此,能够从胸部X射线的图像中,自动识别COVID-19的算法显得尤为重要。本文提出了一种基于深度学习的CXR图像分类框架,该框架能够在有限的训练数据下生成更具判别力的特征。具体而言,首先通过残差神经网络(ResNet34和ResNet50)和Transformer组成多分支分类网络,其中ResNet分支通过深度残差结构,有效地提取丰富的语义信息和细腻的纹理信息;而Transformer分支则通过自注意力机制,捕捉图像的全局语义特征。随后,利用特征交互模块将ResNet分支提取丰富的语义信息和纹理信息,与Transformer提取的全局语义特征进行特征交互。最后,再通过特征融合模块来提取图像的多尺度语义特征。该方法能够在有限训练数据的条件下提取多尺度特征表示,以对COVID-19感染区域进行特征提取和定位。实验在公开DLAI3和COVIDx数据集上与15种方法进行了比较,相比于ResNet50的模型,准确率分别提高了1.37%和0.76%。本文提出的分类方法,结合ResNet和Transformer网络在特征提取上的优点,使得网络对CXR图像的识别结果更加准确。

**关键词:** 胸部X射线检查; 特征交互模块; 多分支分类网络; 残差神经网络; Transformer

**中图分类号:** TP301.6 **文献标识码:** A **DOI:** 10.12466/xhcl.2025.02.005

**引用格式:** 苏华强,雷海军,雷柏英.多分支特征融合分类网络用于CXR图像识别[J].信号处理,2025,41(2):253-266. DOI:10.12466/xhcl.2025.02.005.

**Reference format:** SU Huaqiang, LEI Haijun, LEI Baiying. Multi-branch feature fusion classification network for chest X-ray image recognition[J]. Journal of Signal Processing, 2025, 41(2): 253-266. DOI: 10.12466/xhcl.2025.02.005.

## Multi-Branch Feature Fusion Classification Network for Chest X-Ray Image Recognition

SU Huaqiang<sup>1</sup> LEI Haijun<sup>1</sup> LEI Baiying<sup>\*2</sup>

(1. College of Computer Science and Software Engineering, Shenzhen University, Guangdong Provincial Key Laboratory of Popular High Performance Computers, Shenzhen, Guangdong 518060, China;

2. Department of Biomedical Engineering, School of Medicine, Shenzhen University, Guangdong Key Laboratory for Biomedical Measurements and Ultrasound Imaging, Shenzhen, Guangdong 518060, China)

收稿日期: 2024-07-16; 修回日期: 2024-10-14

\*通信作者: 雷柏英 leiby@szu.edu.cn \*Corresponding Author: LEI Baiying, leiby@szu.edu.cn

基金项目: 国家自然科学基金(62276172);广东省自然科学基金(2023A1515011378);深圳市基础研究专项(JCYJ20230808105602005)

Foundation Items: The National Natural Science Foundation of China (62276172); Natural Science Foundation of Guangdong Province (2023A1515011378); Shenzhen Basic Research Project (JCYJ20230808105602005)

**Abstract:** COVID-19 is an infectious disease caused by the new coronavirus, which poses a significant challenge to global public health. In clinical practice, chest X-ray (CXR) examinations are an important means by which to identify COVID-19 infections and other common lung diseases. However, it is time-consuming and labor-intensive for radiologists to examine COVID-19 patients, and such procedures increase the risk of infection for doctors. Therefore, an algorithm that can automatically identify COVID-19 from chest X-ray images is particularly important. Therefore, this paper proposes a CXR image classification framework based on deep learning that can generate more discriminative features with limited training data. Specifically, a multi-branch classification network is first formed by residual neural networks (ResNet34 and ResNet50) and a Transformer. The ResNet branch effectively extracts rich semantic information and delicate texture information through a deep residual structure, whereas the Transformer branch captures the global semantic features of the image through a self-attention mechanism. Then, the feature interaction module is used to extract rich semantic and texture information from the ResNet branch, and the feature interaction is performed with the global semantic features extracted by the Transformer. Finally, the multiscale semantic features of the image are extracted through the feature fusion module. This method can extract multiscale feature representations under the condition of limited training data to extract features and locate COVID-19 infected areas. The experiment was compared with 15 methods on the public DLAI3 and COVIDx data sets, and the accuracy was improved by 1.37% and 0.76%, respectively, compared with the ResNet50 model. The classification method proposed in this paper combines the advantages of ResNet and Transformer networks in feature extraction to make the recognition results of the network more accurate for CXR images.

**Key words:** chest X-ray; feature interaction module; multi-branch classification network; residual neural network; Transformer

## 1 引言

2019冠状病毒病(COVID-19)是由新型冠状病毒引起的一种传染病,自从2019年冠状病毒疫情暴发,疫情的防控给全球医疗保健系统带来了巨大的负担。根据世界卫生组织报告的工作,截至2022年3月20日,全球累计COVID-19患者病例数高达4.699亿,累计死亡人数为610万<sup>[1]</sup>。COVID-19的主要常规临床诊断通常基于流行病学史、临床表现和各种实验室检测方法,包括核酸扩增试验、计算机断层扫描(Computed Tomography, CT)和血清学技术。然而,人工检测成本高昂、耗时,并且可能增加临床医生的感染风险。基于胸部CT图像的人工智能技术是一种常见且有效的诊断方法,它可以显示病理结果,例如磨玻璃样混浊、实变和结节,提供更多的病理信息和更高的诊断准确度。尽管多数基于CT的研究<sup>[2-3]</sup>重点关注COVID-19诊断,然而,对于儿童或孕妇,高剂量辐射限制了CT的使用<sup>[4]</sup>。

胸部X射线(Chest X-ray, CXR)作为一种医学成像技术,在诊断多种肺部疾病中表现出优异的性能,包括但不限于COVID-19的诊断<sup>[5-9]</sup>。CXR是一种低剂量辐射的便携式胸部检查,可最大限度地减少放射科医生交叉感染的风险,比CT更便宜且更广泛使用<sup>[10-11]</sup>。然而,放射科医生使用医学CXR图像手动诊断COVID-19存在以下几个问题<sup>[12]</sup>:1)手动诊断COVID-19要求很高、耗时且容易出现人为

错误。2)手动诊断需要大量放射科医生,来解释CXR图像上COVID-19放射学特征的细微表现。3)由于CXR图像中的影像学特征很大程度上是重叠的,因此很难将COVID-19与其他病毒性肺炎病例区分开来。为了解决上述问题,利用深度学习技术对医学影像进行自动COVID-19检测是有前景的<sup>[1]</sup>。为了应对这一流行病,研究人员提出了基于深度学习的方法从CXR图像诊断COVID-19,并取得了良好的效果<sup>[13-14]</sup>。基于深度学习分类的方法,可以直接从原始CXR图像中对COVID-19的性质进行分类<sup>[15-18]</sup>,这些网络主要是基于视觉网络VGG16<sup>[19]</sup>、ResNet<sup>[20]</sup>和DenseNet<sup>[21]</sup>等卷积神经网络(Convolutional Neural Network, CNN)框架,可以自动从CXR图像中学习特征表示,以实现图像的精准分类<sup>[18]</sup>。然而,由于CNN的几何结构固定,采样位置也是固定的,不能根据病变的复杂特征而改变<sup>[22-23]</sup>,因此很难从COVID-19的复杂病变中学习放射学特征。

鉴于COVID-19放射学特征的多样性,ZHU等人<sup>[24]</sup>通过将自适应可变形卷积集成到ResNet中,开发了一种自适应可变形ResNet,以自适应提取COVID-19患者的放射学特征。该分类网络通过对空间采样位置应用额外的偏移来增强CNN建模能力,再通过灵活的采样位置,可以根据CXR图像上的形状和尺度自适应调整感受野,从而使特征提取过程适应感染区域的结构变化。WANG等人<sup>[25]</sup>提

出了一种新颖的协作学习框架,用于从分布不均的各种数据集中识别 COVID-19 特征。网络被分成两个分支:一个分支具有轻量级架构和四个不同的卷积层,另一个则具有更密集连接的学习块。然而,这些分类网络不能自适应的从形状复杂的感染区域进行全局的特征提取。PENG 等人<sup>[26]</sup>提出的 Conformer 分类网络,使用 CNN 和视觉变压器作为混合网络结构的两个分支,利用特征耦合单元提取图像的局部和全局特征。它们只是以一种相加的融合方式结合来自两个分支的特征,并将它们用于图像分类。但是这样并不能有效的提取到具有判别力的特征。在 COVID-19 中,典型的 CXR 成像结果包括磨玻璃样混浊和肺实变,通常具有各种不规则形状,如模糊、斑片、弥漫和网状结节<sup>[27]</sup>。此外,在感染的不同阶段和不同患者之间,病变区域的大小和位置差异很大,这增加开发 COVID-19 检测系统的难度。因此,COVID-19 识别的关键问题是如何提取 CXR 图像的全局和多尺度特征,以学习形状复杂的感染区域。

从上面的讨论来看,大多数现有的 COVID-19 识别方法主要是基于自适应可变形 ResNet 框架来学习特征表示<sup>[24, 28]</sup>,尽管该网络具有很强的自适应特征提取能力,但卷积核的局部性质使其无法捕获上下文信息,然而上下文信息通常对于更好地识别图像中的病灶区域特征至关重要<sup>[29, 30]</sup>, COVID-19 患者的病灶经常出现大规模且模糊的边界形态结构,而卷积网络仅在局部邻域中进行特征提取,忽略上下文信息<sup>[31]</sup>,这可能会限制其从复杂病变中学习具有判别力特征的能力。直观上,虽然并非所有 CXR 图像内容都有益于 COVID-19 检测,但分类网络能够提取到图像的上下文信息,就能更好的检测图像上的病变区域并忽略不相关的信息。此外,考虑到 CXR 图像上的 COVID-19 没有任何单一特征是特异性或诊断性的,放射科医生需要通过全局评估整个 CXR 图像并读取某些放射学特征来诊断 COVID-19 患者<sup>[32-33]</sup>。因此,COVID-19 的诊断应结合所有特征影像学表现和 CXR 背景信息进行整体评估。

针对上述问题,本文提出了一种多分支特征融合的分类方法,可以自动聚焦于肺实变和网状混浊等感兴趣区域,并提取鉴别性的放射学特征,以实现从 CXR 图像中对 COVID-19 患者进行诊断。总体来说,本文利用 ResNet 和金字塔视觉 Transformer (Pyramid Vision Transformer, PVT) 模型<sup>[34]</sup>构成多

分支特征融合分类网络,来对 COVID-19 进行识别。该网络通过将空洞空间金字塔池化(Atrous Spatial Pyramid Pooling, ASPP)<sup>[35]</sup>集成到网络中,以有效的提取 COVID-19 患者的放射学特征。ResNet 主要侧重于提取局部特征和纹理信息,而 Transformer 则更侧重于捕获图像的全局关系和上下文信息,将 ResNet 和 Vision Transformer 网络进行交互融合,使得分类网络可以根据 CXR 图像上的形状和尺度信息,提取到具有判别力的全局和多尺度特征,从而使得网络能够学习到丰富的上下文信息和局部特征,以此关注相关的感染区域和捕捉复杂病变的放射学特征。

## 2 方法

在本节中,将详细介绍所提出的基于多分支特征融合分类网络框架,整个框架的流程图如图 1 所示。该网络首先通过利用残差神经网络(ResNet34 和 ResNet50)和 Transformer 组成多分支分类网络,来提取在不同尺度下图像的深度特征表示。其次,使用特征交互融合模块,将 ResNet 分支提取丰富的语义信息和纹理信息,与 Transformer 提取的全局语义特征进行交互,随后,再通过特征融合模块,捕获图像的多尺度语义特征。最后,通过将 ResNet 和 Transformer 分支中捕获的高级特征进行拼接,并利用全连接层来对 CXR 图像进行分类。

### 2.1 多分支特征融合分类网络

CNN 模型的深层次架构对其强大的学习能力至关重要<sup>[20]</sup>,但是其不能有效的捕捉到全局特征。本文中采用了 HE 等人<sup>[20]</sup>提出的卷积神经网络(深度残差神经网络,ResNet34 和 ResNet50),与 WANG 等人<sup>[34]</sup>提出的 PVT 模型构成多分支特征提取网络,相较于典型的 CNN 架构相比,多分支网络引入残差连接和 Transformer 架构,能够在训练非常深的网络时解决梯度消失问题,以及能有效的捕捉图像的全局信息和局部信息。在本研究中,利用了 ResNet34, ResNet50, 和 PVT 构成多分支网络,通过利用特征交互模块来将 ResNet 捕捉到的局部信息与 PVT 网络分支提取的全局信息进行交互融合,同时利用特征融合模块来捕捉多尺度的信息。ResNet 的分支是由一系列残差块构成的,每个块都包含多个堆叠的卷积层(将修正线性单元(ReLu)层和归一化层视为卷积层的一部分)。这些残差块的结构可以被表述为:

$$F_{i+1} = \text{Relu}(F_i + \mathcal{F}(F_i, \mathbf{w}_i)) \quad (1)$$

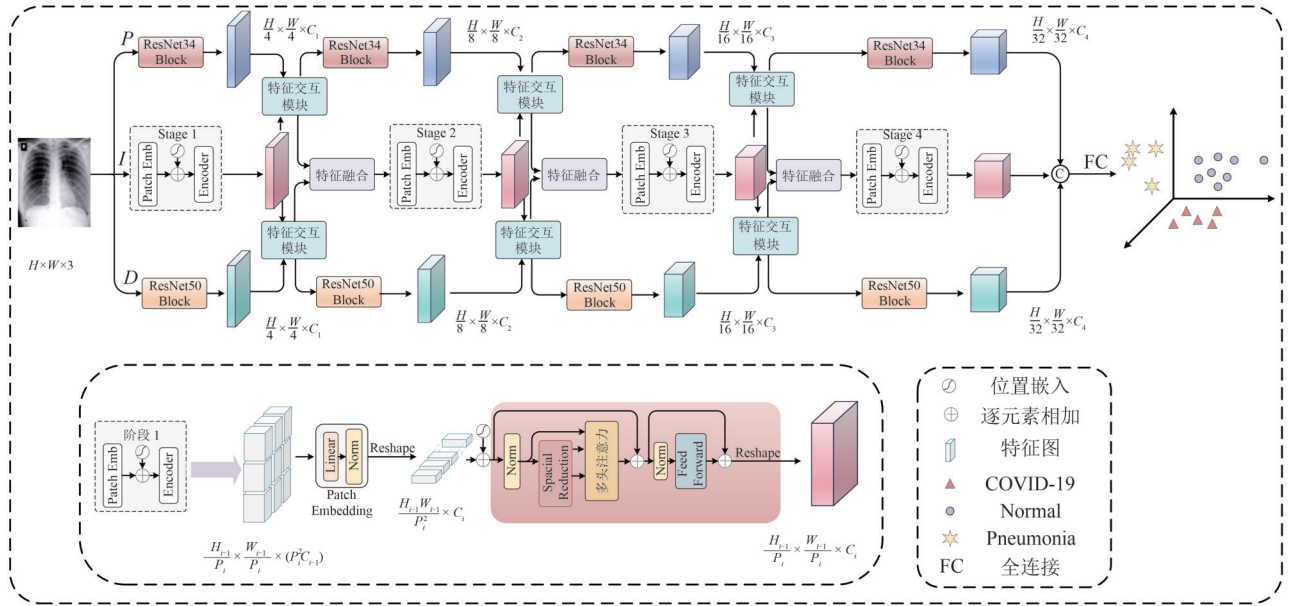


图1 多分支特征融合分类网络的COVID-19识别框架.该网络通过将CNN分支提取到的局部特征,与Transformer分支提取到全局特征进行交互融合,以增强分类模型的性能

Fig. 1 COVID-19 recognition framework based on a multi-branch feature fusion classification network. The network interacts and fuses the local features extracted by the CNN branch with the global features extracted by the Transformer branch to enhance the performance of the classification model

其中,  $F_i$  和  $F_{i+1}$  分别是第  $i$  个残差块的输入和输出;  $\text{Relu}(\cdot)$  是修正线性单元函数;  $\mathcal{F}(\cdot)$  表示残差映射函数,  $w_i$  是残差块的参数。

PVT 网络架构的概述如图 1 所示。与 ResNet 的主干网类似, 也分为四个阶段, 可生成不同尺度的特征图。所有阶段都共享相似的架构, 由补丁嵌入层和 Transformer 编码器层组成。在第一阶段, 给定大小为  $H \times W \times 3$  的输入图像, 首先将其分为  $HW/4^2$  个块。随后, 将展平的特征进行线性投影并获得大小为  $HW/4^2 \times C_1$  的嵌入补丁。最后, 将嵌入的补丁和位置嵌入通过 Transformer 编码器, 并将输出重塑为  $H/4 \times W/4 \times C_1$  大小的特征图  $F_1$ 。同样, 使用前一阶段的特征图作为输入, 可以得到以下特征图:  $F_2$ 、 $F_3$  和  $F_4$ 。

## 2.2 特征交互模块和特征融合模块

多分支特征融合分类网络利用预训练 ResNet 的分类模型, 对输入的 CXR 图像进行处理。所提出的特征交互模块目的是将 ResNet 和 PVT 分支特征进行挖掘和结合出最具鉴别性的特征通道, 并生成更丰富的特征表示, 如图 2 所示。具体的, 给定 ResNet 和 PVT 的第  $i$  个卷积块, 分别产生两个输入特征  $F_i^{\text{ResNet}}$  和  $F_i^{\text{PVT}}$ , 首先使用全局平均池化 (Global Average Pooling, GAP) 来获得  $F_i^{\text{ResNet}}$  和  $F_i^{\text{PVT}}$  特征图

中的全局特征向量。随后, 将这两个特征向量分别输入全连接 (Fully Connected, FC) 和激活函数  $\delta(\cdot)$ , 得到通道注意力向量  $\text{Att}_i^{\text{ResNet}}$  和  $\text{Att}_i^{\text{PVT}}$ , 分别表示  $F_i^{\text{ResNet}}$  特征和  $F_i^{\text{PVT}}$  特征的重要性。最后将通道注意力向量以通道乘法的方式应用于输入特征。这样, 特征交互模块将关注具有鉴别性的特征, 并抑制无关特征。此过程可定义为:

$$\text{Att}_i = \delta(w_i * \text{GAP}(F_i) + b_i) \quad (2)$$

其中,  $w_i$  和  $b_i$  表示第  $i$  个特征向量的 FC 层参数,  $\text{GAP}(\cdot)$  表示全局平均池化操作。生成通道增强特征  $\tilde{F}_i = \text{Att}_i \otimes F_i$ , 其中  $\otimes$  表示通道方向的乘法。

此外, 通道注意力向量  $\text{Att}_i^{\text{ResNet}}$  和  $\text{Att}_i^{\text{PVT}}$  通过  $\text{Max}(\cdot)$  函数对特征进行聚合, 以保留来自 ResNet 和 PVT 分支产生的具有判别力的特征通道, 将其进行归一化操作  $\alpha(\cdot)$ , 将输出特征归一化到 (0~1) 的范围。

从而得到了交互的融合通道注意力向量  $\text{Att}_i^{\text{CR}}$ 。此过程可定义为:

$$\text{Att}_i^{\text{CR}} = \alpha(\text{Max}(\text{Att}_i^{\text{ResNet}}, \text{Att}_i^{\text{PVT}})) \quad (3)$$

基于融合通道注意力向量  $\text{Att}_i^{\text{CR}}$ , 将  $\tilde{F}_i^{\text{ResNet}}$  和  $\tilde{F}_i^{\text{PVT}}$  深度特征与  $\text{Att}_i^{\text{CR}}$  增强后的特征进行相加的操作, 得到  $\check{F}_i^{\text{ResNet}}$  和  $\check{F}_i^{\text{PVT}}$ 。该过程可描述为:

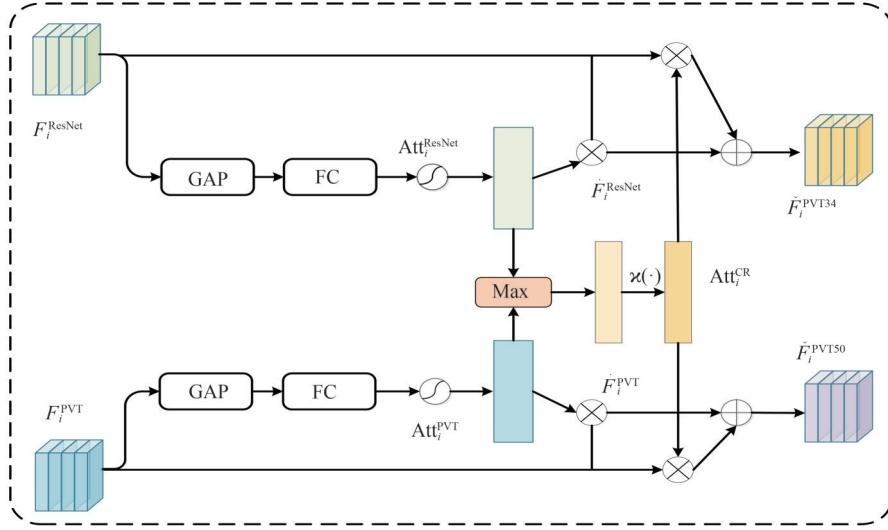


图 2 特征交互模块.FIM旨在通过挖掘和融合全局特征与局部特征之间最具判别力的通道,从而生成信息更加丰富的特征  
Fig. 2 Feature interaction module (FIM). The FIM aims to generate more informative features by mining and fusing the most discriminative channels between global and local features

$$\check{F}_i = \check{F}_i + \text{Att}_i^{\text{CR}} \otimes F_i \quad (4)$$

为了获得多尺度的特征信息,应首先对来自特征交互模块的  $\check{F}_i^{\text{PVT34}}$  和  $\check{F}_i^{\text{PVT50}}$  进行嵌入和压缩,以获得注意力向量。特征融合模块首先将特征沿着通道维度进行全局平均池化  $\text{GAP}(\cdot)$ ,然后再通过 ASPP 和多层感知器 (Multilayer Perceptron, MLP) 的操作来获得注意向量,上述操作可以表述为:

$$I = \text{ASPP}(\text{Cat}(\check{F}_i^{\text{PVT34}}, \check{F}_i^{\text{PVT50}})) \quad (5)$$

其中,  $\text{Cat}(\cdot)$  表示将特征沿着通道进行拼接的操作,  $\text{ASPP}(\cdot)$  表示空洞空间金字塔池化操作,用于提取输入特征的多尺度信息。深度特征的分支注意向量为:

$$W_{\text{PVT34}} = \delta(\text{MLP}(I)) \quad (6)$$

通过上述的操作,多分支特征融合网络可以充分地捕捉多尺度和细节特征,从而能够有效地抑制无关信息。然后,通过将特征  $\check{F}_i^{\text{PVT34}}$  和  $W_{\text{PVT34}}$  进行乘法的操作,可以获得一个噪声较小的特征表示  $F_i^{\text{PVT34}}$ :

$$F_i^{\text{PVT34}} = \check{F}_i^{\text{PVT34}} \otimes W_{\text{PVT34}} + \check{F}_i^{\text{PVT50}} \quad (7)$$

为了充分利用不同分支之间的互补性,需要进行特征的互补聚合。特征融合模块考虑了这两种分支的特征,并为  $\check{F}_i^{\text{PVT34}}$  和  $\check{F}_i^{\text{PVT50}}$  生成空间门,利用注意力机制控制每个特征图的信息流,如图 3 所示。为使空间门更加精确,首先将这两个特征  $F_i^{\text{PVT34}}$  和

$F_i^{\text{PVT50}}$  映射连接起来,以结合它们在空间中特定位置的特征。随后,定义了两个映射函数来将高维特征映射到两个不同的空间门:

$$\mathcal{F}_{\text{PVT34}}: F_{\text{concat2}} \rightarrow G_{\text{PVT34}} \quad (8)$$

$$\mathcal{F}_{\text{PVT50}}: F_{\text{concat2}} \rightarrow G_{\text{PVT50}} \quad (9)$$

其中,  $F_{\text{concat2}}$  为连接特征,  $G_{\text{PVT34}}$  为 PVT 和 ResNet34 经过特征交互后获得的特征图的空间门,  $G_{\text{PVT50}}$  为 PVT 和 ResNet50 经过特征交互后获得特征图,使用  $1 \times 1$  卷积来实现特征的映射,特征图上应用 Softmax:

$$A_{\text{PVT34}}^{(i,j)} = \frac{e^{G_{\text{PVT34}}^{(i,j)}}}{e^{G_{\text{PVT34}}^{(i,j)}} + e^{G_{\text{PVT50}}^{(i,j)}}} \quad (10)$$

$$A_{\text{PVT50}}^{(i,j)} = \frac{e^{G_{\text{PVT50}}^{(i,j)}}}{e^{G_{\text{PVT34}}^{(i,j)}} + e^{G_{\text{PVT50}}^{(i,j)}}} \quad (11)$$

其中,  $G_{\text{PVT34}}^{(i,j)}$  是分配给特征图  $\check{F}_i^{\text{PVT34}}$  中每个位置的权重,  $G_{\text{PVT50}}^{(i,j)}$  是分配给特征图  $\check{F}_i^{\text{PVT50}}$  中每个位置的权重。最终融合的特征  $M$  可以通过加权  $\check{F}_i^{\text{PVT34}}$  和  $\check{F}_i^{\text{PVT50}}$  映射得到。

### 3 实验与结果

#### 3.1 实验设置

1)数据集(DLAI3):DLAI3由以下来源和出版物收集:从门德利检索<sup>[36]</sup>的儿童肺炎 CXR 图像、COVID-19 CXR 数据集<sup>1</sup>, COVID-19 图像数据集<sup>[37]</sup>和 ChestXray8 数据库<sup>[38]</sup>。数据集 DLAI3 包含

<sup>1</sup> <https://github.com/clovaai/deep-text-recognition-benchmark>

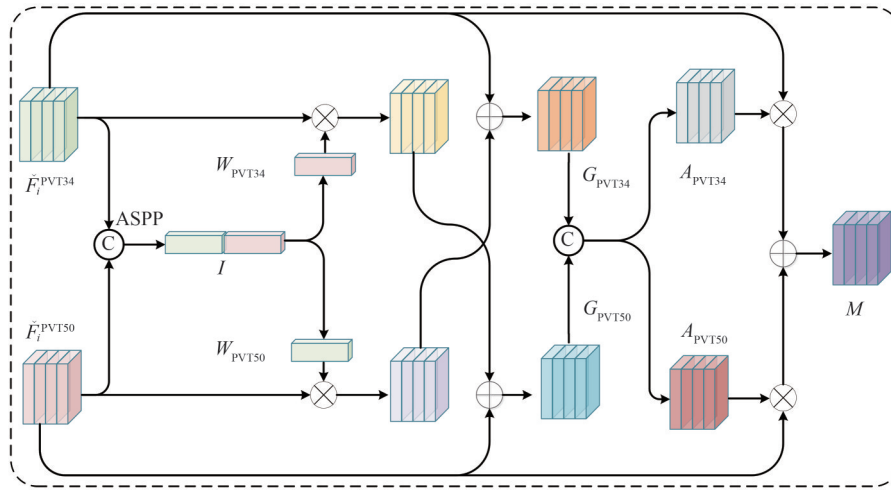


图3 特征融合模块.该模块通过ASPP和特征融合机制,实现不同分支特征的互补聚合

Fig. 3 Feature fusion module (FFM). This module realizes the complementary aggregation of different branch features through ASPP and the feature fusion mechanism

5507例CXR图像,依据相应的分类标签被划分为4407例训练数据集图像和1100例测试数据集图像。该数据集中有三种类别标签分别为:363例确诊的COVID-19阳性病例、1408例正常病例(Normal)和3736例普通肺炎病例(Pneumonia)。

2)数据集(COVIDx)<sup>2</sup>:依据相应的分类标签被划分为12122例训练数据集和3030例测试数据集。COVIDx数据集包含3615例确诊的COVID-19阳性病例、10192例正常病例和1345例普通肺炎病例。

3)数据集(COVID CT)<sup>3</sup>:依据相应的分类标签被划分为15749例训练数据集和3936例测试数据集。COVID CT数据集包括4001例肺炎(Positive CT, PCT)病例、9979例非新冠肺炎(Negative CT, NCT)病例,以及5705例非信息性CT图像(Non-informative CT, NiCT)病例。

4)评估:本研究的目的在于对CXR数据集中的COVID-19、正常、普通肺炎进行分类,为一个常规的三分类任务,通常使用准确率(Accuracy, Acc)、特异性(Specificity, Spe)、精确率(Precision, Pr)、召回率(Recall, Re)、F1-Score、Kappa等评价指标作为模型分类性能的评判标准。所有实验都在具有CPU Intel Xeon E5-2680 @ 2.70 GHz, GPU NVIDIA Quadro K4000和128G RAM的计算机上进行。

### 3.2 图像预处理和数据增强

考虑到数据集中图像特征的多样性和复杂性,

如图4所示,通过对图像进行预处理可以对深度特征网络的提取能力产生显著的影响。在图像分类任务中,图像预处理是一个关键步骤,其目的是通过一系列的操作来使图像更适合用于深度学习模型的训练。

1)图像尺寸变换:在分类网络中,通常会将图像调整为固定大小(如224×224像素)以作为分类网络的输入。通常情况下,为了进行训练和特征提取,会将所有图像调整为相同大小并裁剪为所需尺寸。在本研究中,选择了224×224像素的正方形图像作为分类网络的输入。

2)图像归一化和数据增强操作:在输入到分类网络之前,执行图像归一化的步骤是通过减去整个数据集的平均像素值来实现的,这样使得图像具有相似的统计特性(如零均值和单位方差)。这有助于模型更快地收敛,并提高模型的鲁棒性。对于图像的数据增强操作,则通过对图像进行旋转、平移、缩放、翻转变换,以生成更多样化的训练样本,这有助于减少过拟合,提高模型的泛化能力。

### 3.3 网络类型的实验

在分类对比实验中,除了对34和50层残差网络(ResNet34, ResNet50)进行COVID-19的分类实验外,还研究了其他几种不同深度的CNN模型和不同网络架构的Vision Transformer分类模型进行性能比较,包括8层的AlexNet<sup>[39]</sup>,16层的VGG16<sup>[19]</sup>,Mo-

<sup>2</sup> <https://www.kaggle.com/tawsifurrahman/covid19-radiography-database>

<sup>3</sup> <https://www.kaggle.com/datasets/azaemon/preprocessed-ct-scans-for-covid19>

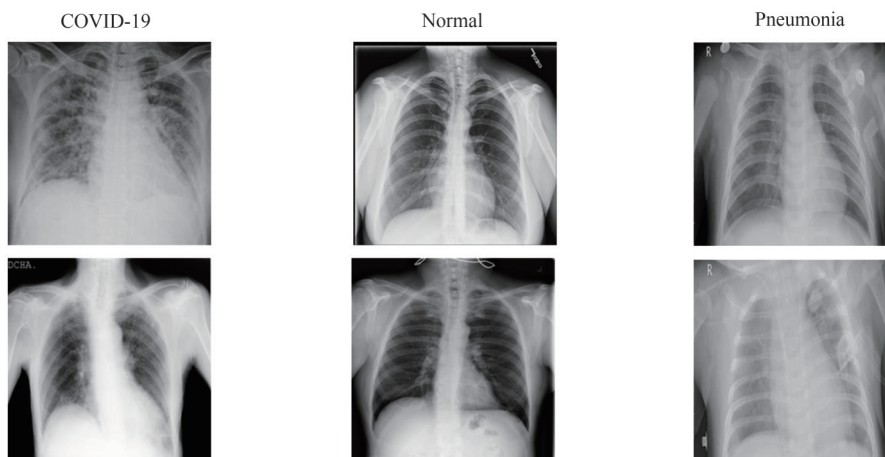


图4 COVID-19的胸部X射线检查图像.COVID-19,Normal,与Pneumonia之间的类内差异大,类间差异小

Fig. 4 Chest X-ray images of COVID-19. The within-class differences between COVID-19, Normal, and Pneumonia are large, whereas the inter-class differences are small

bileNet<sup>[40]</sup>, GoogleNet<sup>[41]</sup>, PVT, CaiT<sup>[42]</sup>, RdNet<sup>[43]</sup>, GC-ViT<sup>[44]</sup>, RepViT<sup>[45]</sup>, RMT<sup>[46]</sup>, SMT<sup>[47]</sup>, SwiftFormer<sup>[48]</sup>, SBCFormer<sup>[49]</sup>, CrossViT<sup>[50]</sup>, COVIDNet<sup>[25]</sup>, Conformer<sup>[26]</sup>。实验设置中采用相同的图像预处理方法,来对不同的分类模型进行公平的分类性能比较。从表1、表2的实验结果,可以观察到不同的分类网络模型对 COVID-19 的识别有很大影响,如基于 CNN 网络架构的 AlexNet 相比于基于 Vision Transformer 网络能够实现更好的分类性能。另外,对于在 COVIDx 数据集上的实验结果,基于 CNN 的分类方法也有不同的分类效果,如从较深的 ResNet50 中提取的特征优于从 AlexNet 网络中提取到的特征,对于评估指标 Acc 而言,两个网络大约有 2% 的差距。

为了评估网络架构对 COVID-19 识别精度的影响,在公开的数据集(DLAI3、COVIDx)上对多分支特征融合网络进行了消融实验,以及与不同的分类方法进行对比实验,来验证所提出的模块和分类网络的有效性,实验结果如表3、表4、表1、表2所示。从消融实验的结果中可以得到,特征交互模块能够有效将 ResNet 分支中捕捉到的局部特征与 PVT 捕捉到的全局特征进行交互融合,特征融合模块的主要作用是提取特征的多尺度信息,使得分类网络能够适应各种复杂病变特征,提高分类精度。同时,利用 Grad-CAM<sup>[51]</sup>可视化定位输入图像中对影响模型分类结果的重要区域,以此来提高模型的可解释性。从图5可得出,网络通过引入特征交互和特征融合模块,使网络对于病灶区域的响应更加敏感,

这说明所提出的网络能够定位到病变区域,从而网络能够准确地获得分类结果。

t-SNE 图(t-Distributed Stochastic Neighbor Embedding)用于可视化高维数据。通过 t-SNE 对比图,可以看到不同分类模型提取的特征在特征空间中的分布情况。同一类别的样本在 t-SNE 图中聚集在一起,而不同类别的样本分布较远。各个类别在 t-SNE 图中差距越明显,表明模型在学习特征表示方面的能力越强。从图6和图7的 t-SNE 可视化结果中可以观察到,所提出的分类网络在学习特征表示方面,相较于其他网络更强,这是由于基于多分支特征融合网络能够捕捉到多尺度特征以及能够利用特征交互的机制,有效地将来自 ResNet 中的局部信息和 PVT 分支提取到的全局信息进行交互融合。为了验证模型的泛化性,在不同的数据集上对模型进行验证,实验结果如表2所示。

此外,为了验证所提出分类模型在不同模态下的泛化能力,在 COVID CT 数据集上对基于双分支网络(COVIDNet、Conformer、CrossViT)进行了对比实验,实验结果如表5所示。所提出的多分支特征融合分类网络通过特征交互模块,有效结合了来自 CNN 的局部特征与 Transformer 的全局特征。虽然这种方法增加了模型的参数量,但相比参数量相近的多分支 Conformer 网络,通过特征融合模块聚合多尺度特征并抑制无关信息,从而有效提升了分类网络的性能。从 t-SNE 图的可视化效果可以看出,所提出的网络在有限的训练数据下,能够生成更多具有判别力的特征。再通过利用混淆矩阵来

表1 在数据集DLAI3上不同分类网络的分类结果(%)

Tab. 1 Classification results of different classification networks on the DLAI3 data set (%)

Network	Acc	Pr	Re	Spe	F1	Kappa
VGG16	94.91	94.43	93.27	96.98	93.69	91.86
AlexNet	95.00	93.77	93.19	96.30	93.48	90.20
GoogleNet	95.27	95.06	90.71	96.81	92.57	91.04
MobileNet	95.45	94.17	91.04	96.76	92.51	92.15
ResNet34	95.73	94.64	93.94	97.25	94.25	92.29
ResNet50	96.18	95.22	91.65	97.29	93.27	91.61
PVT	88.27	84.00	74.52	90.80	78.16	73.70
CaiT	89.82	81.08	81.88	93.21	81.41	78.38
RdNet	84.18	70.61	64.44	89.56	66.63	65.67
GCViT	87.55	77.33	85.63	92.66	79.68	75.89
RepViT	88.18	77.84	84.44	93.11	80.22	78.69
RMT	91.73	87.85	85.50	94.97	86.43	85.80
SMT	90.18	83.11	83.97	93.18	83.38	82.19
SwiftFormer	88.82	81.13	78.11	92.67	79.44	79.66
SBCFormer	92.73	92.65	85.97	94.56	88.80	86.08
<b>Ours</b>	<b>97.55</b>	<b>97.16</b>	<b>95.40</b>	<b>98.11</b>	<b>96.25</b>	<b>94.75</b>

表2 在数据集COVIDx上不同分类网络的分类结果(%)

Tab. 2 Classification results of different classification networks on the COVIDx data set (%)

Network	Acc	Pr	Re	Spe	F1	Kappa
VGG16	95.02	95.14	91.94	95.67	93.43	90.52
AlexNet	95.87	95.60	93.18	96.56	94.66	92.26
GoogleNet	96.30	94.72	95.30	97.15	94.90	92.28
MobileNet	95.94	93.54	96.72	97.86	95.04	92.78
ResNet34	96.61	94.94	95.29	97.47	95.08	93.27
ResNet50	97.89	97.73	96.36	98.23	97.03	95.60
PVT	81.91	75.35	84.68	90.82	78.76	66.27
CaiT	78.09	76.07	68.17	81.32	70.72	53.78
RdNet	73.07	69.87	59.01	77.27	62.60	40.31
GCViT	81.85	78.73	76.46	87.93	77.24	63.46
RepViT	85.64	82.07	82.65	90.51	82.31	71.31
RMT	91.68	88.75	90.37	94.19	89.48	83.39
SMT	92.77	91.54	90.39	95.20	90.89	84.27
SwiftFormer	87.89	85.52	83.44	90.90	84.42	74.72
SBCFormer	94.36	92.35	92.73	96.03	92.53	88.14
<b>Ours</b>	<b>98.65</b>	<b>97.97</b>	<b>98.15</b>	<b>99.03</b>	<b>98.06</b>	<b>97.29</b>

直观地显示多分支特征融合网络对三个数据集的分类结果,如图8所示。实验结果表明,所提出的分类模型在DLAI3, COVIDx和COVID CT数据集中具有较高的准确率,这是由于该模型能够结合来自ResNet分支捕捉到的语义信息和纹理信息,以及PVT分支提取的全局性特征的优势,并且通过特征

融合模块中的ASPP的操作来提取多尺度特征,使得分类网络能够捕捉COVID-19患者的各种复杂病变特征。

#### 4 讨论

本文提出了一种新颖有效的分类方法,用于从

表 3 在数据集 DLAI3 上多分支特征融合网络的消融实验结果 (%)

**Tab. 3 Ablation experimental results of the multi-branch feature fusion network on the DLAI3 data set (%)**

ResNet34	ResNet50	PVT	FIM	FFM	Acc	Pr	Re	Spe	F1	Kappa
√					95.73	94.64	93.94	97.25	94.25	92.29
	√				96.18	95.22	91.65	97.29	93.27	91.61
√	√	√			96.64	95.46	94.19	97.87	94.77	93.87
√	√	√	√		97.00	96.67	95.18	97.71	95.91	94.76
√	√	√	√	√	<b>97.55</b>	<b>97.16</b>	<b>95.40</b>	<b>98.11</b>	<b>96.25</b>	<b>94.75</b>

表 4 在数据集 COVIDx 上多分支特征融合网络的消融实验结果 (%)

**Tab. 4 Ablation experimental results of the multi-branch feature fusion network on the COVIDx data set (%)**

ResNet34	ResNet50	PVT	FIM	FFM	Acc	Pr	Re	Spe	F1	Kappa
√					96.61	94.94	95.29	97.47	95.08	93.27
	√				97.89	97.73	96.36	98.23	97.03	95.60
√	√	√			98.15	97.62	96.95	98.49	97.27	96.12
√	√	√	√		98.32	97.89	97.82	98.82	97.86	96.75
√	√	√	√	√	<b>98.65</b>	<b>97.97</b>	<b>98.15</b>	<b>99.03</b>	<b>98.06</b>	<b>97.29</b>

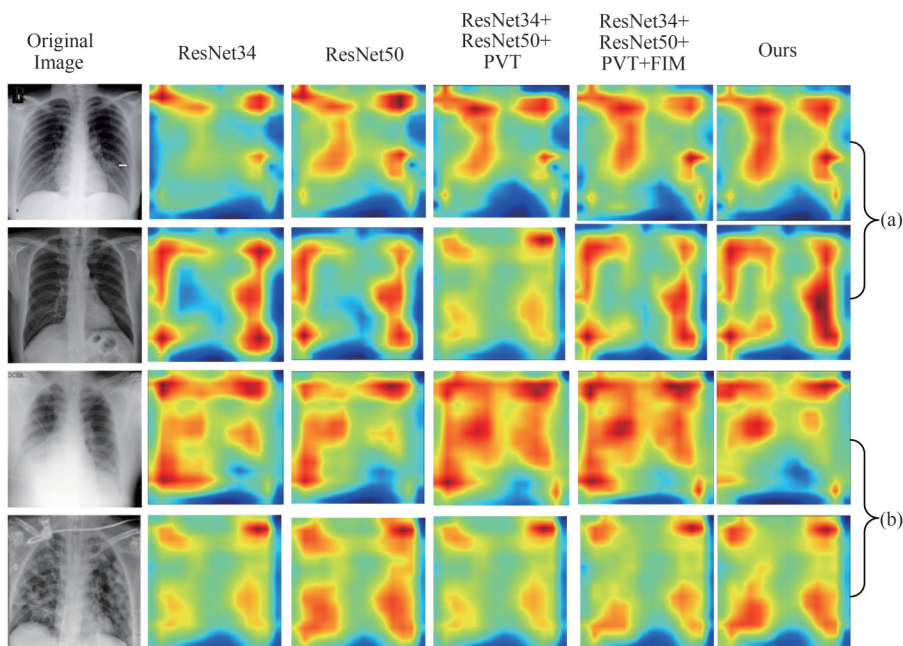


图 5 多分支特征融合网络在两个数据集上的 Grad-CAM 图：(a) DLAI3 数据集，(b) COVIDx 数据集

Fig. 5 Grad-CAM graph of a multi-branch feature fusion network on two data sets: (a) DLAI3 data set, (b) COVIDx data set

CXR 检查图像中自动准确地识别 COVID-19, 在公开的数据集 (DLAI3、COVIDx、COVID CT) 进行了大量的实验来验证其对 COVID-19 分类的有效性。多分支特征融合网络的分类方法启发于基于 CNN 框架能够有效捕捉到图像的分类特征。例如, 图 5、图 8 和表 1、表 2、表 3、表 4 所示, 在预训练的

ResNet34 和 ResNet50 中, 将从分类模型提取到的输出特征画成 t-SNE 图, 来验证模型能够从图像中, 提取到具有判别力的特征。对于卷积层, 通常用成百上千的卷积核生成大量的激活图, 以捕获图像的各种特征, 但是由于 CNN 可能无法通过在不同空间位置的卷积层中使用共享权重来有效地



图6 在DLAI3数据集上不同分类方法的t-SNE可视化

Fig. 6 t-SNE visualization of different classification methods on the DLAI3 data set

捕获全局特征,以及针对在CXR图像感染的区域大小不一致的问题,CNN可能无法捕捉到多尺度的特征,因此所提出的多分支特征融合分类网络通过将ResNet和Transformer进行交互融合,来有效地对全局特征和局部特征的提取,以及利用ASPP来捕捉多尺度特征,以此从CXR图像中提取具有判别力的特征表示。与基于CNN和Transformer现有分类方法相比,多分支特征融合分类网

络使用预训练的ResNet和PVT组成的多分支特征融合网络作为特征提取器,从而能更好提取不同尺度的特征。为了验证所提出的特征交互模块和特征融合模块的有效性,在有限的训练样本下,对各模块进行消融实验,并且获得了良好的分类效果。

多分支特征融合分类网络虽然在数据集(DLAI3、COVIDx、COVID CT)上分类性能有一定

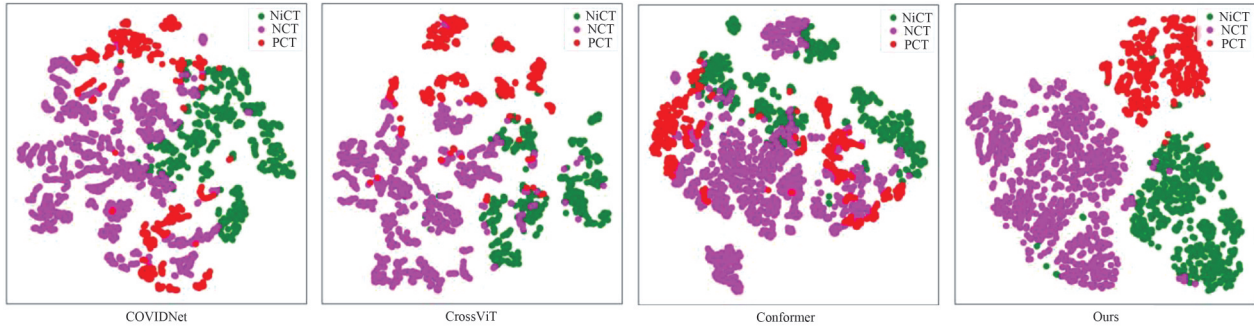


图 7 在 COVID CT 数据集上不同分类方法的 t-SNE 可视化

Fig. 7 t-SNE visualization of different classification methods on the COVID CT data set

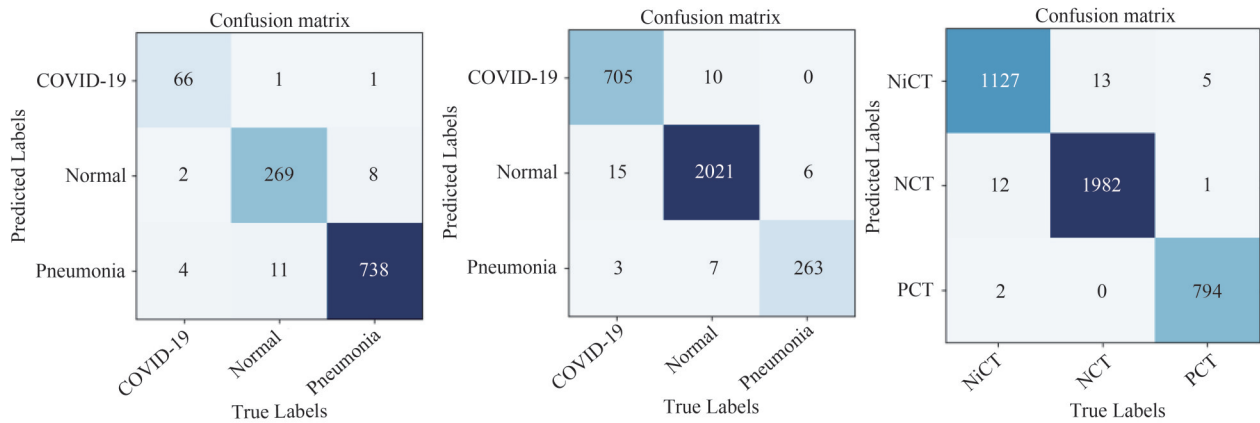


图 8 多分支特征融合网络在三个数据集上的混淆矩阵:DLAI3 数据集(左),COVIDx 数据集(中),COVID CT 数据集(右)  
Fig. 8 Confusion matrix of the multi-branch feature fusion network on three data sets: DLAI3 data set (Left), COVIDx data set (Medium), COVID CT data set (Right)

表 5 在数据集 COVID CT 上不同分类网络的分类结果

Tab. 5 Classification results of different classification networks on the COVID CT data set

Network	Acc	Pr	Re	Spe	F1	Kappa	Params/M
COVIDNet	96.93%	96.65%	96.69%	98.34%	96.67%	94.56%	48.68
CrossViT	97.33%	96.76%	97.70%	98.67%	97.20%	96.51%	27.48
Conformer	97.54%	97.88%	96.84%	98.53%	97.33%	96.48%	101.74
<b>Ours</b>	<b>99.16%</b>	<b>99.18%</b>	<b>99.12%</b>	<b>99.54%</b>	<b>99.15%</b>	<b>98.59%</b>	<b>102.47</b>

的提升,但仍然存在一些限制。首先,针对标注数据较少的问题,本文仅在 CXR 和 CT 等标注的图像数据集上进行了监督学习,而没有应用半监督学习方法来结合未标注数据进行 COVID-19 的识别研究。另外,多分支特征融合分类网络利用了 ImageNet 上预先训练过的网络,但是对于其他数据集上的预训练模型没有做实验。

### 5 结论

本文提出一个多分支特征融合网络,用于从

CXR 识别 COVID-19 病例。该网络通过 ResNet 分支捕捉局部特征,PVT 分支捕捉全局特征,进行交互融合,从而实现对感染区域及复杂结构的精确定位。特征融合模块通过 ASPP 来扩大感受野,适应病变区域的不同形状,以提取 COVID-19 的放射学特征。在数据集(DLAI3、COVIDx、COVID CT)上的实验结果表明,所提出的分类网络取得了最好的分类结果。未来,将多分支特征融合网络尝试扩展到,用 CXR 图像进行 COVID-19 联合诊断和严重程度预测。

## 参考文献

- [1] LIU Yanbei, LI Henan, LUO Tao, et al. Structural attention graph neural network for diagnosis and prediction of COVID-19 severity [J]. *IEEE Transactions on Medical Imaging*, 2023, 42(2): 557-567.
- [2] ZHOU Longxi, LI Zhongxiao, ZHOU Juexiao, et al. A rapid, accurate and machine-agnostic segmentation and quantification method for CT-based COVID-19 diagnosis [J]. *IEEE Transactions on Medical Imaging*, 2020, 39(8): 2638-2652.
- [3] YAO Qingsong, XIAO Li, LIU Peihang, et al. Label-free segmentation of COVID-19 lesions in lung CT [J]. *IEEE Transactions on Medical Imaging*, 2021, 40(10): 2808-2819.
- [4] WANG Zheng, XIAO Ying, LI Yong, et al. Automatically discriminating and localizing COVID-19 from community-acquired pneumonia on chest X-rays [J]. *Pattern Recognition*, 2021, 110: 107613.
- [5] PALURU N, DAYAL A, JENSSEN H B, et al. Anamnet: Anamorphic depth embedding-based lightweight CNN for segmentation of anomalies in COVID-19 chest CT images [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2021, 32(3): 932-946.
- [6] ZHANG Kang, LIU Xiaohong, SHEN Jun, et al. Clinically applicable AI system for accurate diagnosis, quantitative measurements, and prognosis of COVID-19 pneumonia using computed tomography [J]. *Cell*, 2020, 181(6): 1423-1433.e11.
- [7] HARMON S A, SANFORD T H, XU Sheng, et al. Artificial intelligence for the detection of COVID-19 pneumonia on chest CT using multinational datasets [J]. *Nature Communications*, 2020, 11(1): 4080.
- [8] MEI Xueyan, LEE H C, DIAO Kaiyue, et al. Artificial intelligence-enabled rapid diagnosis of patients with COVID-19 [J]. *Nature Medicine*, 2020, 26(8): 1224-1228.
- [9] OUYANG Xi, KARANAM S, WU Ziyang, et al. Learning hierarchical attention for weakly-supervised chest X-ray abnormality localization and diagnosis [J]. *IEEE Transactions on Medical Imaging*, 2021, 40(10): 2698-2710.
- [10] VARELA-SANTOS S, MELIN P. A new approach for classifying coronavirus COVID-19 based on its manifestation on chest X-rays using texture features and neural networks [J]. *Information Sciences*, 2021, 545: 403-414.
- [11] JACOBI A, CHUNG M, BERNHEIM A, et al. Portable chest X-ray in coronavirus disease-19 (COVID-19): A pictorial review [J]. *Clinical Imaging*, 2020, 64: 35-42.
- [12] SHOEIBI A, KHODATARS M, JAFARI M, et al. Automated detection and forecasting of COVID-19 using deep learning techniques: A review [J]. *Neurocomputing*, 2024, 577: 127317.
- [13] SHI Feng, WANG Jun, SHI Jun, et al. Review of artificial intelligence techniques in imaging data acquisition, segmentation, and diagnosis for COVID-19 [J]. *IEEE Reviews in Biomedical Engineering*, 2021, 14: 4-15.
- [14] YU Xiang, LU Siyuan, GUO Lili, et al. ResGNet-C: A graph convolutional neural network for detection of COVID-19 [J]. *Neurocomputing*, 2021, 452: 592-605.
- [15] MARQUES G, AGARWAL D, DE LA TORRE DÍEZ I. Automated medical diagnosis of COVID-19 through EfficientNet convolutional neural network [J]. *Applied Soft Computing*, 2020, 96: 106691.
- [16] MINAEE S, KAFIEH R, SONKA M, et al. DeepCOVID: Predicting COVID-19 from chest X-ray images using deep transfer learning [J]. *Medical Image Analysis*, 2020, 65: 101794.
- [17] NOUR M, CÖMERT Z, POLAT K. A Novel Medical Diagnosis model for COVID-19 infection detection based on Deep Features and Bayesian Optimization [J]. *Applied Soft Computing*, 2020, 97: 106580.
- [18] ISMAEL A M, ŞENGÜR A. Deep learning approaches for COVID-19 detection based on chest X-ray images [J]. *Expert Systems with Applications*, 2021, 164: 114054.
- [19] SIMONYAN K, ZISSERMAN A. Very deep convolutional networks for large-scale image recognition [EB/OL]. 2014: 1409.1556. <https://arxiv.org/abs/1409.1556v6>.
- [20] HE Kaiming, ZHANG Xiangyu, REN Shaoqing, et al. Deep residual learning for image recognition [C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Las Vegas, NV, USA. IEEE, 2016: 770-778.
- [21] HUANG Gao, LIU Zhuang, VAN DER MAATEN L, et al. Densely connected convolutional networks [C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Honolulu, HI, USA. IEEE, 2017: 2261-2269.
- [22] DAI Jifeng, QI Haozhi, XIONG Yuwen, et al. Deformable convolutional networks [C]//2017 IEEE International Conference on Computer Vision (ICCV). Venice, Italy. IEEE, 2017: 764-773.
- [23] JIN Qiangguo, MENG Zhaopeng, PHAM T D, et al. DUNet: A deformable network for retinal vessel segmentation [J]. *Knowledge-Based Systems*, 2019, 178: 149-162.
- [24] ZHU Xizhou, HU Han, LIN S, et al. Deformable ConvNets V2: More deformable, better results [C]//2019

- IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach, CA, USA. IEEE, 2019: 9300-9308.
- [25] WANG Zhao, LIU Quande, DOU Qi. Contrastive cross-site learning with redesigned net for COVID-19 CT classification [J]. *IEEE Journal of Biomedical and Health Informatics*, 2020, 24(10): 2806-2813.
- [26] PENG Zhiliang, HUANG Wei, GU Shanzhi, et al. Conformer: Local features coupling global representations for visual recognition [C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada. IEEE, 2021: 357-366.
- [27] STOGIANNOS N, FOTOPOULOS D, WOZNITZA N, et al. COVID-19 in the radiology department: What radiographers need to know [J]. *Radiography*, 2020, 26(3): 254-263.
- [28] ROBERTS M, DRIGGS D, THORPE M, et al. Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans [J]. *Nature Machine Intelligence*, 2021, 3: 199-217.
- [29] ZHANG Ning. Learning adversarial transformer for symbolic music generation [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2023, 34(4): 1754-1763.
- [30] BELLO I, ZOPH B, LE Q, et al. Attention augmented convolutional networks [C]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). Seoul, Korea (South). IEEE, 2019: 3285-3294.
- [31] SHI Baoguang, YANG Mingkun, WANG Xinggang, et al. ASTER: An attentional scene text recognizer with flexible rectification [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019, 41(9): 2035-2048.
- [32] WONG H Y F, LAM H Y S, FONG A H T, et al. Frequency and distribution of chest radiographic findings in patients positive for COVID-19 [J]. *Radiology*, 2020, 296(2): E72-E78.
- [33] CHUNG M, BERNHEIM A, MEI Xueyan, et al. CT imaging features of 2019 novel coronavirus (2019-nCoV) [J]. *Radiology*, 2020, 295(1): 202-207.
- [34] WANG Wenhai, XIE Enze, LI Xiang, et al. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions [C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada. IEEE, 2021: 548-558.
- [35] CHEN L C, PAPANDREOU G, SCHROFF F, et al. Rethinking atrous convolution for semantic image segmentation [EB/OL]. 2017: 1706.05587. <https://arxiv.org/abs/1706.05587v3>.
- [36] KERMANY D S, GOLDBAUM M, CAI Wenjia, et al. Identifying medical diagnoses and treatable diseases by image-based deep learning [J]. *Cell*, 2018, 172(5): 1122-1131.
- [37] COHEN J P, MORRISON P, DAO Lan, et al. COVID-19 image data collection: Prospective predictions are the future [EB/OL]. 2020: 2006.11988. <https://arxiv.org/abs/2006.11988v3>.
- [38] WANG Xiaosong, PENG Yifan, LU Le, et al. ChestX-Ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases [C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Honolulu, HI, USA. IEEE, 2017: 3462-3471.
- [39] KRIZHEVSKY A, SUTSKEVER I, HINTON G E. ImageNet classification with deep convolutional neural networks [J]. *Communications of the ACM*, 2017, 60(6): 84-90.
- [40] HOWARD A G, ZHU Menglong, CHEN Bo, et al. MobileNets: Efficient convolutional neural networks for mobile vision applications [EB/OL]. 2017: 1704.04861. <https://arxiv.org/abs/1704.04861v1>.
- [41] SZEGEDY C, LIU Wei, JIA Yangqing, et al. Going deeper with convolutions [C]//2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Boston, MA, USA. IEEE, 2015: 1-9.
- [42] TOUVRON H, CORD M, SABLAYROLLES A, et al. Going deeper with image transformers [C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada. IEEE, 2021: 32-42.
- [43] KIM D, HEO B, HAN D. DenseNets reloaded: Paradigm shift beyond ResNets and ViTs [EB/OL]. 2024: 2403.19588. <https://arxiv.org/abs/2403.19588v2>.
- [44] HATAMIZADEH A, YIN Hongxu, KAUTZ J, et al. Global context vision transformers [C]//2023 International Conference on Machine Learning (ICML). Honolulu, HI, United states. PMLR, 2023: 12633-12646.
- [45] WANG Ao, CHEN Hui, LIN Zijia, et al. Rep ViT: Revisiting mobile CNN from ViT perspective [C]//2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA. IEEE, 2024: 15909-15920.
- [46] FAN Qihang, HUANG Huaibo, CHEN Mingrui, et al. RMT: Retentive networks meet vision transformers [C]//2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Seattle, WA, USA. IEEE, 2024: 5641-5651.
- [47] LIN Weifeng, WU Ziheng, CHEN Jiayu, et al. Scale-aware modulation meet transformer [C]//2023 IEEE/

- CVF International Conference on Computer Vision (ICCV). Paris, France. IEEE, 2023: 5992-6003.
- [48] SHAKER A, MAAZ M, RASHEED H, et al. SwiftFormer: Efficient additive attention for transformer-based real-time mobile vision applications [C]//2023 IEEE/CVF International Conference on Computer Vision (ICCV). Paris, France. IEEE, 2023: 17379-17390.
- [49] LU Xiangyong, SUGANUMA M, OKATANI T. SBCFormer: Lightweight network capable of full-size ImageNet classification at 1 FPS on single board computers [C]//2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV). Waikoloa, HI, USA. IEEE, 2024: 1112-1122.
- [50] CHEN C F R, FAN Quanfu, PANDA R. CrossViT: Cross-attention multi-scale vision transformer for image classification [C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Montreal, QC, Canada. IEEE, 2021: 347-356.
- [51] SELVARAJU R R, COGSWELL M, DAS A, et al. Grad-CAM: Visual explanations from deep networks via gradient-based localization [J]. International Journal of Computer Vision, 2020, 128(2): 336-359.

## 作者简介



**苏华强** 男, 1994年生, 海南东方人。深圳大学在读博士研究生, 主要研究方向为医学图像处理。  
E-mail: 18233271771@163.com



**雷海军** 男, 1968年生, 湖南郴州人。深圳大学副教授, 主要研究方向为医学图像处理。  
E-mail: lhj@szu.edu.cn



**雷柏英** 女, 1982年生, 湖南郴州人。深圳大学副教授、博士生导师, 主要研究方向为医学图像处理。  
E-mail: leiby@szu.edu.cn

(责任编辑: 刘建新)