

## 基于 VPN 通道下的加密流量分类算法

魏洁玲, 马秀丽, 金彦亮, 王 瑞

上海大学 通信与信息工程学院, 上海 200444

**摘 要:** 为了改善网络管理水平、加强网络安全监督, 针对虚拟专用网络 (virtual private network, VPN) 通道下流量加密性强、不透明度高的特点, 设计了加密流量数据的新构图方式, 提出了基于变体 ResNet18 网络的加密流量分类算法。为了验证算法有效性, 采集真实 VPN 通道下的热门 app 流量, 成功实现了多 VPN 通道下的多应用流量分类。所提算法最终在公有数据集与真实采集数据集上的分类准确率分别达到 98.1% 和 96.0%。实验结果表明, 该算法具有通用性且具有一定的实际价值。

**关键词:** 深度学习; 加密流量; 虚拟专用网络; 残差网络

中图分类号: TP309

文章编号: 0255-8297(2023)04-0646-11

## Encrypted Traffic Classification Algorithm Based on VPN Channel

WEI Jieling, MA Xiuli, JIN Yanliang, WANG Rui

*School of Communication and Information Engineering, Shanghai University,  
Shanghai 200444, China*

**Abstract:** This paper proposes a new encrypted traffic classification algorithm based on a variant ResNet18 network to improve network management and strengthen network security supervision. A three-channel image construction is designed to address the strong encryption and high opacity characteristics of traffic in virtual private network (VPN) channels. The proposed method successfully identifies different apps' traffic in different VPN channels, as validated using popular apps' traffic collected from real VPN channels. The algorithm achieves 98.1% and 96.0% classification accuracy on public and self-collected datasets, respectively. Experimental results demonstrate the algorithm's universality and practical value.

**Keywords:** deep learning, encrypted traffic, virtual private network (VPN), residual network

随着互联网应用的广泛普及与传播, 人们在访问互联网时产生了大量的网络流量。其中, 加密流量在网络中广泛使用, 所占总流量的份额正在逐年提升。因此, 加密流量分类算法逐渐

收稿日期: 2021-09-22

基金项目: 国家自然科学基金 (No. 61771299) 资助

通信作者: 马秀丽, 副教授, 研究方向为大数据和智能信息处理。E-mail: xlma@shu.edu.cn

受到学术界和工业界的广泛关注, 对加密流量进行正确识别分类已成为网络安全和网络管理中的重要课题。虚拟专用网络 (virtual private network, VPN) 是许多公司或个人常用的通信隧道, 可在公网上保证流量数据的私有性与完整性。在 VPN 通道下实现加密流量分类, 筛选出未知应用或非法应用具有极大的研究意义。

最初, 流量协议可以直接依照与之对应的端口号进行识别。比如 20 端口对应着应用文件传输协议 (file transfer protocol, FTP), 而 80 端口对应着超文本传输协议 (hypertext transfer protocol, HTTP)。然而随着动态端口技术和网络安全技术的发展, 基于端口号的流量分类算法开始不再适用。随后, 研究者们开始关注流量中的有效载荷, 这种方法又被称为深度包检测技术 (deep packet inspection, DPI)<sup>[1]</sup>。由于不同的加密协议或应用往往具有其特定的数据包格式, 比如含有固定字符串, DPI 可以从不同种类流量的有效载荷里, 探索归纳数据间的特征规律。

目前, 随着计算机技术的火热, 越来越多的科研工作者选择利用深度学习来解决加密流量分类问题。文献 [2] 提出了端到端的一维卷积神经网络 (convolutional neural network, CNN) 架构, 首次将图像分类运用到流量分类领域。文献 [3] 对加密流量进行先流后包的二次分割, 采用胶囊神经网络<sup>[4]</sup>进行分类训练, 更有效地学习了数据包的空间特征和数据包之间的时间序列特征。文献 [5] 提出 Seq2Img 在线流量分类算法, 该算法基于数据包的大小、间隔时间和前向后向信息, 利用再生核希尔伯特空间 (reproducing kernel Hilbert space, RKHS)<sup>[6]</sup> 将流量的原始部分序列转换为多通道图像, 并使用 CNN 将这些图像分类到不同的应用程序。文献 [7] 提出了一种循环神经网络 (recurrent neural network, RNN) 与 CNN 并行的网络结构 App-Net, 其中 RNN 用于从数据包长度序列中学习统计特征, 而 CNN 用于从数据包有效载荷中学习内容特征, 最终在 Dataset-80 上取得 91.05% 的准确率。

上述方案普遍采用公有数据集进行验证, 针对的是单 VPN 与非 VPN 下的应用分类。综上, 本文提出了一种基于内容的加密流量分类算法, 对实验数据集进行了扩充, 实现了多 VPN 通道下的多应用分类。

## 1 基本原理

### 1.1 VPN 通信过程

如图 1 所示, VPN 网关为了保障用户的信息安全, 往往会提供访问控制、报文加密、报文认证、报文封装等功能。发送方发送的明文在经过 VPN 网关后会变成封装加密报文, 随后经过 IP 安全隧道进行传输, 而接收方会收到由接收方 VPN 网关解密后的明文报文。

报文在 VPN 通道下进行通信的关键技术一般包括三个方面: 隧道化协议、认证协议以及加密技术。隧道化协议是 VPN 实现内部网地址通信与多协议通信的重要技术, 它将报文进行分组封装。认证协议通过用户的用户名以及口令来验证该用户是否有权访问。加密技术能够对明文报文提供不同层面的安全保护。比如 IPsec 协议族中的封装安全载荷 (encapsulate security payload, ESP) 能够实现对 IP 数据项的可认证性、完整性以及机密性支持。

### 1.2 数据包结构

根据开放式系统互联 (open system interconnection, OSI) 参考模型, 流量数据会在网络传输的每一层中附加一个头部, 如图 2 所示。每一层的首部往往会含有一些重要的特征信息。数据链路层头部包含了源 MAC 地址, 目的 MAC 地址与类型信息; 网络层头部包含了版本、长度、传输层协议类型、源地址、目的地址等信息; 而传输层头部包含了流量的源端口、目的端口、头部长度、校验和等信息。五元组由源 IP 地址、目的 IP 地址、源端口、目的端口和传

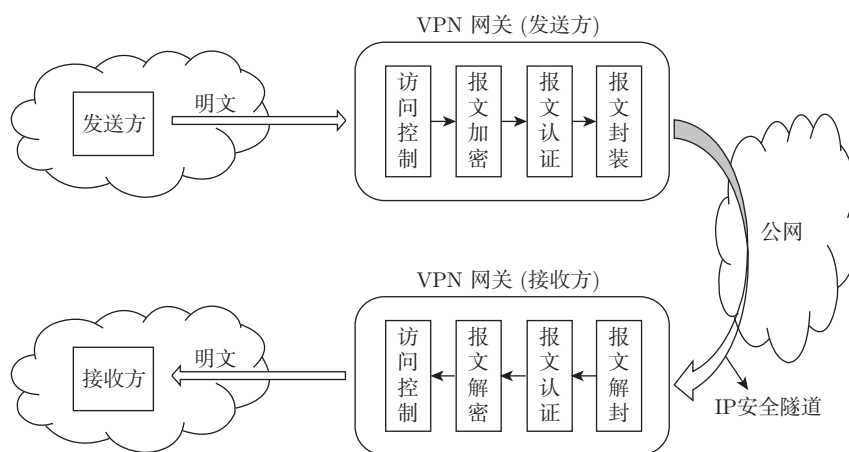


图 1 VPN 通信过程

Figure 1 VPN communication process

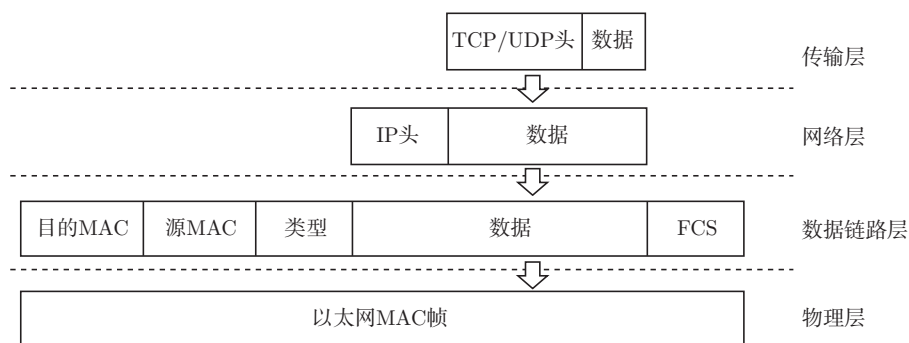


图 2 数据包结构

Figure 2 Packet structure

输层协议相同的一组流量集合构成,常被用于流量分类基础。

### 1.3 残差网络

残差网络<sup>[8]</sup>由微软提出,并获得了2015年ImageNet挑战赛图像分类任务冠军。在深度学习领域里,残差网络是继AlexNet<sup>[9]</sup>、VGG<sup>[10]</sup>、GoogleNet<sup>[11]</sup>后一大标志性创新分类神经网络,解决了近年来随着网络层数越来越深而导致的梯度消失难题,给分类神经网络的发展方向带来了新的思路。

与一般的分类神经网络相比,残差网络采用了跳层连接方式,如图3所示。

假设输入 $x$ 的目标函数为 $H(x)$ 。式(1)为普通神经网络的直接映射,式(2)为残差块中的残差映射。

$$H(x) = x \quad (1)$$

$$H(x) = G(x) + x \quad (2)$$

从表达式来看,拟合残差 $G(x) = H(x) - x$ 比直接拟合 $H(x)$ 更容易实现。事实上,只需调整 $G(x)$ 中的权重以及偏差使其为0,即可实现输入 $x$ 到目标函数 $H(x)$ 的恒等映射。

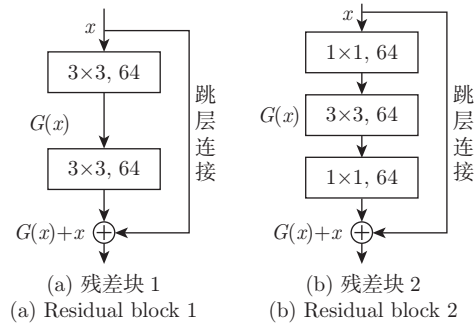


图 3 残差块

Figure 3 Residual block

针对不同的网络深度, 残差网络设计出两类跳层连接残差块, 主要由卷积、批归一化、和激活函数组成。随着网络深度的增加, 残差块添加了  $1 \times 1$  降维卷积层, 并在  $3 \times 3$  卷积完成后进行了  $1 \times 1$  升维卷积, 这样在保证精度的同时减少了计算的成本。

标准残差网络共有 5 种: ResNet18、ResNet34、ResNet50、ResNet101、ResNet152。

## 2 VPN 通道下的加密流量分类算法

### 2.1 分类网络模型

分类网络模型参考 ResNet18 网络, 包含深度为 64、128、256、512 的四大残差块。其中, 每一个残差块由多个残差单元组成, 而每一个残差单元又包含了 2 个普通  $3 \times 3$  的卷积。

ResNet18 网络默认输入为  $224 \times 224 \times 3$ 。根据加密流量数据集的大小, 本文的分类网络设计输入为  $32 \times 32 \times 3$ 。为了保持更多的细粒度特征, 采用  $3 \times 3$  卷积代替原本的  $7 \times 7$  卷积与最大池化层。改进后的变体 ResNet18 网络结构如图 4 所示。

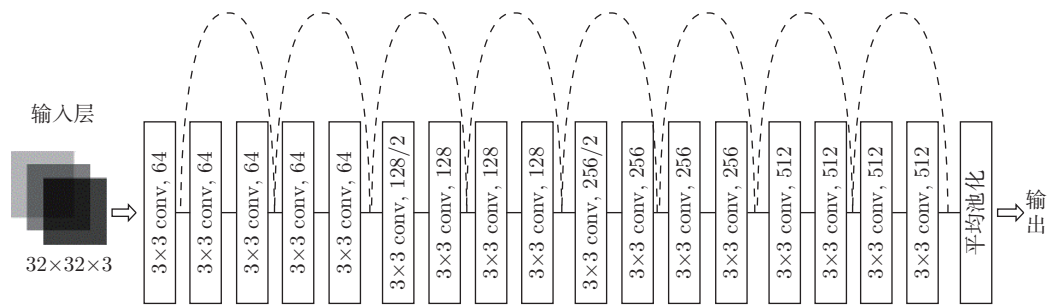


图 4 变体 ResNet18 网络结构图

Figure 4 Variant ResNet18 network structure

### 2.2 预处理算法

在流量分类领域里, 对流量数据的处理大致可分为 4 种粒度: 基于主机、基于流、基于数据包、基于比特。本文为实现多 VPN 通道下的多应用分类, 选择在数据包层面进行预处理。

对原始加密数据流进行以下操作: 数据清洗、基于数据包层面的特征提取、删除 Mac 和 IP 地址、数据包拼接、转换字节、归一化、堆叠生成三通道 RGB 图片、制作训练集与测试集。整个预处理过程如图 5 所示。

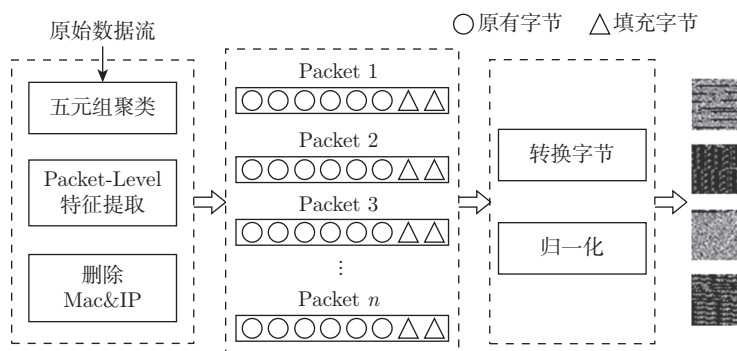


图5 预处理流程

Figure 5 Preprocessing procedure

1) 数据清洗。根据原始数据流的源 IP 地址、源端口、目的 IP 地址、目的端口和传输层协议进行五元组聚类, 过滤无关的流量噪声, 去除非目标用户流量的影响。

2) 删除 Mac 地址和 IP 地址。为了避免神经网络学习到的是固定地址特征, 导致模型过拟合, 对每条数据包中的 Mac 地址和 IP 地址进行了删除操作。

3) 数据包拼接。在每个数据包末尾填补一定长度的字节, 进行拼接。此外, 依据数据包的方向对填充字节设置了不同的数值: 上行数据包间隔填充 0, 下行数据包间隔填充 240。填充字节强调了数据包的统计信息, 如方向、频率等, 其与原始字节相结合的方式能更好地对流量的内容特征与统计特征进行表征。

4) 将流量字节转换成数组并归一化。以 3 072 字节为单位转换成相同大小的数组, 并进行归一化, 加快神经网络训练过程中的收敛速度。

5) 生成三通道 RGB 图像, 可视化如图 6 所示。采用 reshape 函数将归一化后的一维数组转换成三维数组, 数组中的每一位元素作为像素, 最终形成大小为  $32 \times 32$  的三通道图像。

6) 制作训练集与测试集。本文基于 TensorFlow 框架, 对生成的三通道 RGB 图片和标签转存为 TFRecords 标准格式, 训练集与测试集依照 9:1 划分。

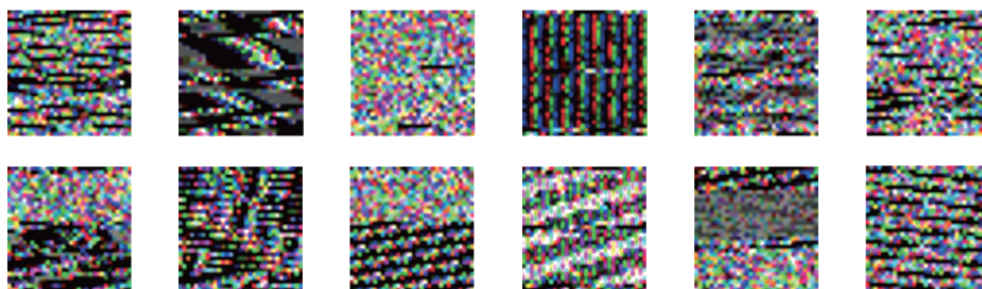


图6 RGB 图像可视化

Figure 6 Visualization of RGB images

### 3 实验

本文总体实验流程如图 7 所示。3.1 节介绍了采用的两种加密流量数据集以及采集过程。3.2 节将数据集输入变体残差网络进行了相关对照训练并分析结果, 证明了基于 VPN 通

道下的加密流量分类算法的有效性。

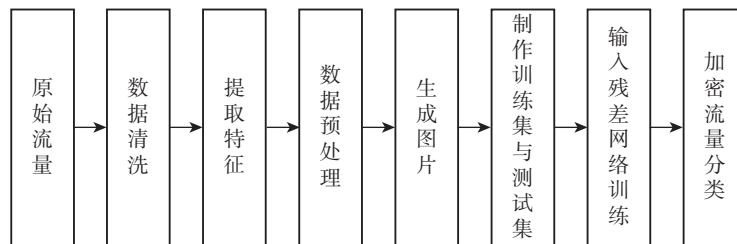


图 7 实验流程图

Figure 7 Experimental flowchart

### 3.1 数据集

为了验证算法的有效性和实用性。本文采用的数据集由两部分组成: ISCXVPN2016 公开数据集和自采数据集。

ISCXVPN2016<sup>[12]</sup> 由加拿大网络安全研究所提供。该数据集基于 OpenVPN 通道, 使用 Wireshark 和 TCPDUMP 捕获, 生成的数据总量约为 28 GB。该数据集包括了基于非 VPN 和 VPN 的 7 种加密类型流量: 聊天、电子邮件、流、文件传输、网页浏览器、语音传输、P2P, 并涉及到 ICQ、Facebook、Hangouts、YouTube、Vimeo、Skype 等众多应用软件。ISCXVPN 2016 公开数据集的具体流量分类类型如表 1 所示。其中, 实验所用的训练集与测试集具体样本数量如表 2 所示。

同时, 为了验证算法的实用价值, 本文使用 Wireshark 软件捕捉安卓手机执行脚本后所产生的流量数据。实验采集了基于 2 种 VPN 通道 (Ssr 和 Psiphon) 下 4 种 App 软件 (YouTube、Zalo、微博、微信) 共 8 种加密流量样本。在采集过程中, 为了保证样本的多样性和真实性, 实验使用了文本、视频、语音、静图、动图等多种消息格式, 同时每条消息采取随机时间间隔发送。真实采集数据集中每种加密流量类型采集的数据包大小如表 3 所示。

表 1 ISCXVPN 2016 流量分类类型

Table 1 Classes of traffic in ISCXVPN2016

流量类型	应用
(VPN-) 网页浏览器	Firefox, Chrome
(VPN-) 电子邮件	Email, Gmail
(VPN-) 聊天	ICQ, AIM, Skype, Facebook, Hangouts
(VPN-) 流	Netflix, Spotify, Vimeo, YouTube
(VPN-) 文件传输	Skype, FTPS, SFTP
(VPN-) 语音传输	Facebook, Skype, Hangouts voice calls
(VPN-) P2P	BitTorrent, uTorrent

### 3.2 实验结果与分析

分类实验采用 4 个评价指标: 准确率 Accuracy、精确率 Precision、召回率 Recall 与 F1。

表 2 ISCXVPN2016 样本数量

**Table 2** Number of samples in ISCXVPN2016

流量类型	训练样本	测试样本
聊天	6 573	731
电子邮件	2 586	288
文件传输	25 656	2 851
P2P	20 842	2 316
流	26 787	2 976
语音传输	26 700	2 967
(VPN) 聊天	7 182	798
(VPN) 电子邮件	2 047	227
(VPN) 文件传输	25 747	2 861
(VPN)P2P	24 324	2 703
(VPN) 流	26 678	2 964
(VPN) 语音传输	26 603	2 956

表 3 自采数据集数据包数量

**Table 3** Number of packets in self-collected dataset

流量类型	数据包/条
Ssr + YouTube	100 341
Ssr + Zalo	107 500
Ssr + 微博	115 280
Ssr + 微信	85 926
Psiphon + YouTube	198 334
Psiphon + Zalo	205 062
Psiphon + 微博	240 600
Psiphon + 微信	171 277

计算过程为

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \quad (3)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (4)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (5)$$

$$F1 = \frac{2TP}{2TP+FP+FN} \quad (6)$$

式中: FN 表示预测结果为负, 预测错了的数量; TN 表示预测结果为负, 预测对了的数量; FP 表示预测结果为正, 预测错了的数量; TP 表示预测结果为正, 预测对了的数量。

首先, 实验研究了填充字节长度对于分类的影响。考虑到图像的边长为 32, 在预处理部分, 将填充字节设置了 3 个尺度: 0、24、48。当填充字节长度为 0 时, 所有数据包字节都是按顺序自然相连。由于图片尺寸是固定的, 在同一幅图片中, 填充的字节越多, 原始的数据包负载占比越低, 但不同数据包之间的间隔也会愈发明显。如表 4 所示, 当填充长度为 24 字节时, 分类结果表现最佳, 准确率达到了 98.1%。

表 4 不同填充字节长度的对比结果

Table 4 Comparison results under different fill byte scales

填充长度/字节	准确率	精确率	召回率	F1
0	0.97	0.96	0.95	0.95
24	0.98	0.96	0.96	0.96
48	0.97	0.95	0.95	0.95

其次, 为了验证数据头部信息是否影响了加密流量分类算法的有效性, 在其他预处理步骤相同的前提下, 实验进行了含有头部信息和不含头部信息构图的对比试验。如表 5 所示, 实验结果表明头部含有对加密流量分类有利的信息。为了解释其头部信息对分类的重要性, 采用了 Grad-CAM 算法可视化网络重点关注的区域, 如图 8 所示。实际上, 数据包的头部往往包含协议类型、数据包长度、版本号以及生存时间等有用的特征信息。本文认为在构图中增加头部信息在实际应用中可能会降低算法泛化性, 但却能够有效地提升算法的准确率与召回率。

表 5 不同构图方式的对比结果

Table 5 Comparison results of different composition methods

头部信息	准确率	精确率	召回率	F1
带有头部	0.98	0.96	0.96	0.96
不带头部	0.83	0.84	0.83	0.83

再次, 实验采用原有 ResNet18 网络和变体 ResNet18 网络对比训练。如表 6 所示, 实验结果表明, 在  $32 \times 32$  的小尺寸图片输入下, 变体 ResNet18 网络具有更好的分类效果。同时, 本文也与其他加密流量分类算法进行了实验对比, 如表 7 所示。在公有数据集的表现上, 本文具有更高的精确率以及召回率。

在真实采集数据集中, 实现了基于两类 VPN 下的四类应用分类, 识别准确率达到了 96.01%。混淆矩阵是常用于模型精度评估的一种标准格式。如图 9 和 10 所示, 以下表示了 ISCXVPN2016 公有数据集与真实采集数据集的具体分类情况。

实验结果表明, 在公有数据集 ISCXVPN2016 中大部分流量实现了 90% 以上的准确率, 而电子邮件类流量因原始样本过少准确率低于 90%。因此在采集真实应用流量的过程中, 本



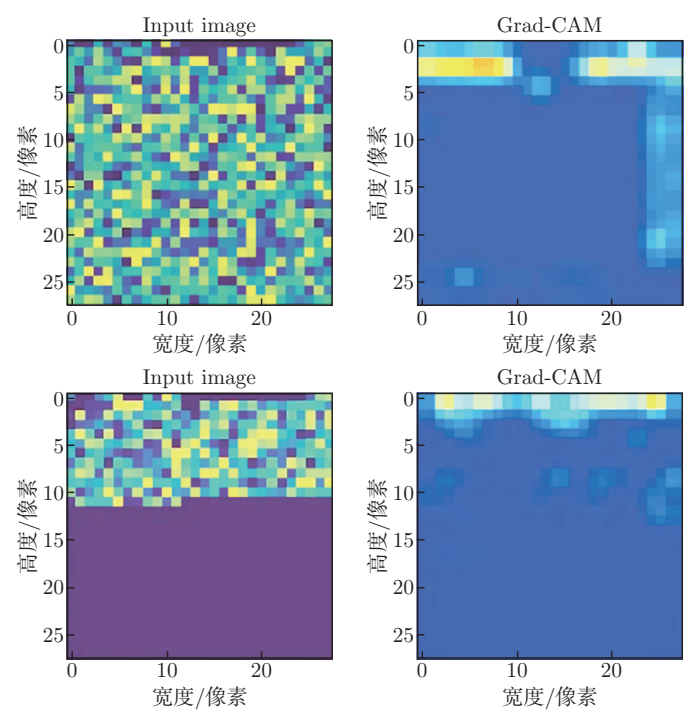


图 8 网络可视化解释

Figure 8 Visual explanations for network

表 6 不同网络结构的对比结果

Table 6 Comparison results of different network structures

网络	准确率	精确率	召回率	F1
ResNet18	0.90	0.85	0.84	0.84
变体 ResNet18	0.98	0.96	0.96	0.96

表 7 不同分类算法的对比结果

Table 7 Comparison results of different algorithms

算法	精确率	召回率	F1
1dCNN	0.85	0.86	0.86
SAE	0.92	0.92	0.92
CNN+LSTM	0.91	0.91	0.91
本文	0.98	0.96	0.96

文均匀地收集了不同类型的流量，成功地解决了样本不均衡的问题，实现了所有流量类型的高准确率分类。

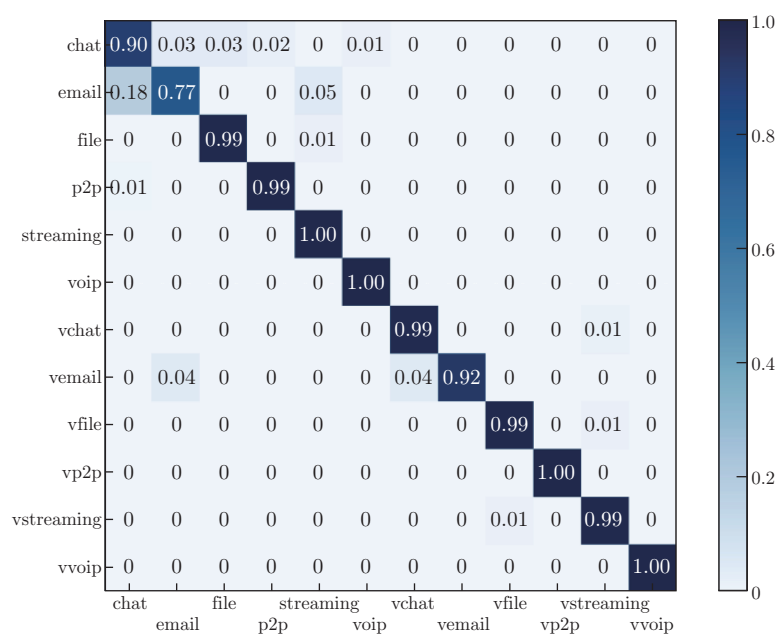


图 9 ISCXVPN2016 上的混淆矩阵

Figure 9 Confusion matrix on ISCXVPN2016

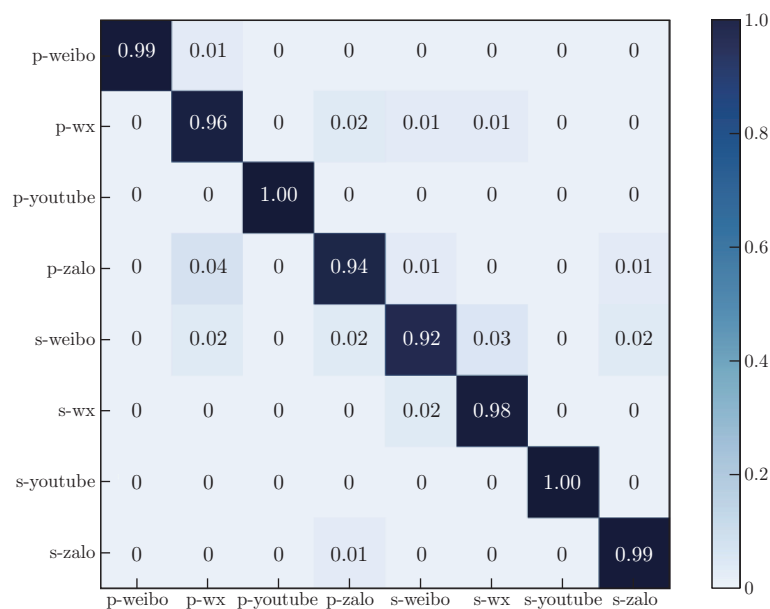


图 10 真实采集数据集上的混淆矩阵

Figure 10 Confusion matrix on self-collected dataset

## 4 结 语

本文提出了一种有效的基于 VPN 通道的加密流量分类算法, 实现了在不同 VPN 通道下

不同 app 加密流量的分类,并在公有数据集与真实采集数据集上均实现 95% 以上的分类准确率。与传统图像分类任务不同的是,加密流量分类任务中特征的选择与处理极为关键。在未来的工作中,将基于 VPN 通道和 App 的加密通信原理上去思考更多的预处理构图方式。另外,在 VPN 通道和 App 频繁更新版本时的互联网背景下,实现可扩展、可迁移数据集的有效分类也是亟待解决的问题。

### 参考文献:

- [1] DHARMAPURIKAR S, KRISHNAMURTHY P, SPROULL T, et al. Deep packet inspection using parallel Bloom filters [C]//11th Symposium on High Performance Interconnects, 2003: 44-51.
- [2] WANG W, ZHU M, WANG J L, et al. End-to-end encrypted traffic classification with one-dimensional convolution neural networks [C]//2017 IEEE International Conference on Intelligence and Security Informatics (ISI), 2017: 43-48.
- [3] CUI S S, JIANG B, CAI Z Z, et al. A session-packets-based encrypted traffic classification using capsule neural networks [C]//2019 IEEE 21st International Conference on High Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS), 2019: 429-436.
- [4] RAJASEGARAN J, JAYASUNDARA V, JAYASEKARA S, et al. DeepCaps: going deeper with capsule networks [C]//2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019: 10717-10725.
- [5] CHEN Z T, HE K, LI J, et al. Seq2Img: a sequence-to-image based approach towards IP traffic classification using convolutional neural networks [C]//2017 IEEE International Conference on Big Data (Big Data), 2017: 1271-1276.
- [6] ROSIPAL R, TREJO L J. Kernel partial least squares regression in reproducing kernel Hilbert space [J]. Journal of Machine Learning Research, 2001, 2: 97-123.
- [7] WANG X, CHEN S H, SU J S. App-net: a hybrid neural network for encrypted mobile traffic classification [C]//IEEE INFOCOM 2020 - IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS), 2020: 424-429.
- [8] HE K M, ZHANG X Y, REN S Q, et al. Deep residual learning for image recognition [C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016: 770-778.
- [9] KRIZHEVSKY A, SUTSKEVER I, HINTON G E. ImageNet classification with deep convolutional neural networks [J]. Communications of the ACM, 2017, 60(6): 84-90.
- [10] SIMONYAN K, ZISSERMAN A. Very deep convolutional networks for large-scale image recognition [DB/OL]. 2014 [2021-09-22]. <https://arxiv.org/abs/1409.1556>.
- [11] SZEGEDY C, LIU W, JIA Y Q, et al. Going deeper with convolutions [C]//2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015: 1-9.
- [12] DRAPER-GIL G, LASHKARI A H, MAMUN M S I, et al. Characterization of encrypted and VPN traffic using time-related features [C]//International Conference on Information Systems Security and Privacy (ICISSP), 2016: 407-414.

(编辑: 王 雪)