

基于 SDN 的实际网络流中 Tor 网页复合特征提取方法

言洪萍, 周强, 王世豪, 姚旺, 何刘坤, 王良民

(江苏大学计算机科学与通信工程学院, 江苏 镇江 212013)

摘 要: 基于网站指纹(WF)攻击的 Tor 网页流量识别方法往往建立在分离好的 Tor 流量甚至是分离好的 Tor 网页流量的基础上, 但从实际网络的原始流中分离出 Tor 流量, 再从 Tor 流量中分离出 Tor 网页流量, 其计算量和困难程度远高于 Tor 网页流量的 WF 攻击本身。根据目前互联网的体系结构, 利用网络流量汇聚到区域中心节点的特点, 通过中心节点的 SDN 结构所提供的域内全局视角, 结合 Tor 网络公开的节点信息提出了一种区分 Tor 流量的双向统计特征(BSF), 可以有效分离 Tor 流量; 进而提出了一种基于 LSF 技术的网页流量隐藏特征提取方法, 从而获得了基于 BSF 和 LSF 的复合流量特征(CTTF); 在此基础上, 针对当前 Tor 流量训练数据较少的问题, 提出了一种基于平移的流量数据增强方法, 使增强后的流量数据与真实工作环境中捕获的 Tor 流量数据分布尽量一致。实验结果表明, 基于 CTTF 与仅使用原始数据特征相比, 识别率提高了 4%左右, 在训练数据较少时, 使用流量数据增强方法后分类效果提升更加明显, 并且可以有效降低误报率。

关键词: 流量发现; 流量识别; 统计特征; 数据增强

中图分类号: TP393

文献标志码: A

DOI: 10.11959/j.issn.1000-436x.2022056

Composite Tor traffic features extraction method of webpage in actual network flow based on SDN

YAN Hongping, ZHOU Qiang, WANG Shihao, YAO Wang, HE Liukun, WANG Liangmin

School of Computer Science and Communication Engineering, Jiangsu University, Zhenjiang 212013, China

Abstract: Website fingerprinting (WF) methods for Tor webpage traffic are often based on the separated Tor traffic or even the separated Tor webpage traffic. However, distinguishing Tor traffic from the original traffic of the actual network and Tor webpage traffic from the Tor traffic costs amount of computation, which is more difficult than the WF attack itself. According to the current architecture of the Internet and the characteristics of network traffic converging to regional central nodes, the bi-directional statistical feature (BSF) was proposed for distinguishing Tor traffic through the intra-domain global perspective provided by the SDN structure of the central node and the node information disclosed by the Tor network. Furthermore, a hidden feature extraction method for Web traffic based on lifted structure fingerprinting (LSF) was proposed, and a composited Tor-webpage-identification traffic feature (CTTF) was proposed based on BSF and LSF deep features. For solving the problem of traffic training data scarcity, a traffic data augmentation method based on translation was proposed, which made the augmented traffic data as consistent as the Tor traffic data captured in the real working environment. The experimental results show that the identification rate based on CTTF can be improved by about 4% compared with using only the original data features. When there is less training data, the classification accuracy is improved more obvious after using the traffic data augmentation method, and the false positive rate can be effectively reduced.

Keywords: traffic discovery, traffic classification, statistical feature, data augmentation

收稿日期: 2021-11-12; 修回日期: 2022-02-17

基金项目: 国家自然科学基金资助项目 (No.U1736216)

Foundation Item: The National Natural Science Foundation of China (No.U1736216)

0 引言

网站指纹(WF, website fingerprinting)攻击^[1-3]作为一类典型的基于流量特征识别的去匿名技术,通过挖掘加密网页流量中存在的流量特征来识别犯罪分子匿名访问的网站,是科学研究和打击网络犯罪实践的热点。但是,现有的WF攻击工作主要研究如何识别Tor流量对应的匿名网页,这在实验室环境下往往具有很好的效果。但是在实际应用中,首先还需从原始网络流量中精准地区分出Tor流量,否则现有的WF攻击效果将大大降低且不可用。因此,从实际网络流中识别Tor网页流量是WF攻击的研究成果走向实用的基础性和关键性课题。

在实际的流量环境中,面向Tor流量的WF攻击方法如果需要发挥实际用途,需要建立在区分Tor流量与其他流量的基础上,该工作的困难主要在于Tor流量相对于实际应用流量的比重非常小,因而训练数据的收集和特征的分离均非常困难。收集Tor流量的困难还有网络结构的问题,由于公网IPv4资源缺乏,目前Internet广泛使用网络端口地址转换技术解决内部网络地址访问外部网络资源的问题,内部网络的所有主机均共享一个合法外部IP地址,导致位于外部网络的审查者难以收集Tor流量。为此,本文充分发掘了现有网络体系结构中作为区域中心节点的数据中心的功能,这些作为网络中心节点的数据中心,不仅是Tor流量汇聚的关键位置,而且在实施中往往采用软件定义网络(SDN, software defined network)^[4]作为网络架构。充分利用数据中心在网络体系中的关键作用,并发挥数据中心的SDN架构将网络设备控制面和数据面分离的特点,能对网络流量进行灵活控制,本文以此作为提高Tor流量监控范围与收集能力的基础^[5],提出了一种基于SDN架构的数据中心下大范围Tor流量发现的方法,进而提出了有效提升Tor网页流量识别效果的复合特征表示与流量数据增强方法,其主要贡献包括如下三点。

1) 基于SDN架构获取能体现应用协议交互过程的网络流量信息,提出一种区分Tor流量和其他流量的双向统计特征(BSF, bi-direction statistic feature)。

2) 提出融合BSF和LSF(lifted structure fingerprinting)的复合流量特征(CTTF, composed

Tor-webpage-identification traffic feature),其中LSF是使用LS Loss^[6]训练深度神经网络提取的深度特征。

3) 提出一种对Tor流量数据进行增强的方法,通过最小化训练数据与真实工作环境中捕获的Tor流量数据分布差异获取增强数据,用于改进和提升Tor网页流量识别模型在训练数据较少情况下的识别效果。

本文对仿真环境下采集的数据进行了实验,并将CTTF结合当前相关方法进行对比实验,实验表明,其均能提升原方法的识别率;基于增强的Tor流量数据基础上提取的CTTF,在面向实际网络流识别时,即使训练数据较少,也能明显提升识别结果,并且可以有效降低误报率。

1 相关工作

本文研究面向实际流量的Tor网页特征识别,建立在数据中心是Tor流量汇聚的关键场景以及SDN架构提供了全局流量分析数据的基础上。为此,相关工作部分主要介绍SDN架构中的安全假设、SDN架构中的流量发现、网页流量识别和基于Tor网页识别的指纹攻击。

1.1 SDN架构中的安全假设

本文主要探讨SDN架构下对Tor的流量收集和流量识别,因此审查者需要收集Tor用户的流量用于进一步分析。本文做出以下两点假设。

1) 假设审查者为数据中心级别,即审查者位于数据中心内,且具有修改OpenFlow流表项的能力^[7]。这意味着即使审查者与Tor用户不在同一网段下,其仍然能通过操纵流规则将同属于一个数据中心下的用户流量重定向到未经授权的接收者或单纯地阻碍流的转发。

2) 假设Tor用户同样存在于数据中心内,通过SDN转发设备与Tor网络建立三跳匿名连接访问不同的服务。

基于以上两点假设的具体网络结构如图1所示,在该结构内,SDN控制器通过北向应用程序接口(API, application program interface)实现与SDN应用交互^[8],通过南向控制-数据平面接口实时监控数据中心各交换机的流量情况以及网络拓扑情况;转发层的设备通过接收控制器的控制信息生成转发表,并对接收到的流量进行按表查找,若存在匹配的流表项,则根据流表项的指示转发数据包。默认情况下,系统内的Tor用户通过与Tor网络建立三跳连接从而匿名访问不同的站点。Tor用户首次

通过匿名连接访问网页产生的流量数据包将被 SDN 交换机捕获并询问控制器以选择合适的路径传输给目标地址。由于审查者具有修改 OpenFlow 流表项的能力，因此可以在流表项中添加一条将该用户流量复制转发到自己本地的动作，进一步分析该网络流量，从而判断用户的实际行为。

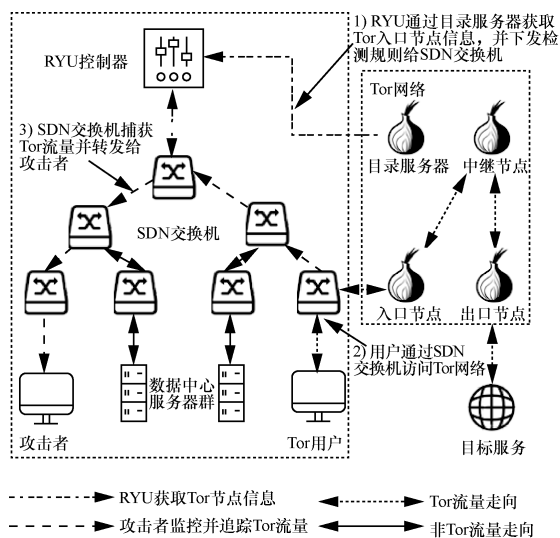


图 1 基于 SDN 结构的 Tor 服务示意

1.2 SDN 架构中的流量发现

SDN 架构通过开放的 API 和协议来动态管理和控制网络，由于控制平面与数据平面是解耦的，这样可使控制平面通过操纵数据平面中流量的路径和走向，从而给应用于各种匿名通信系统中的流量发现方案提供机会。

Oconnor 等^[9]通过在每个主机的内核中标记 IP 报头的服务类型字段，利用 SDN 的代理服务器跟踪 APT (advanced persistent threat) 流量的来源，以此在 SDN 交换机上检测出被标记的数据包。但是，这样就需要在每个主机上安装一个定制的内核，以便实现精确的基于标签的跟踪。由于在大型 SDN 中实现此条件较困难，该方案只局限于小型受控制的网络环境。

Ling 等^[10]提出了一种新颖实用的匿名流量发现技术来确定可疑服务器和用户之间的通信关系，利用目标服务器端的 SDN 交换机来拦截指向服务器的目标流量，并修改发布的 TCP 窗口大小，从而改变服务器端的流量速率。通过精心地改变流量速率，将一个秘密信号调制到流量中，承载该信号的流量通过匿名通信系统到达用户端的 SDN 交换机。然后从用户端的流量中检测出调制信号，以确定服

务器和用户之间的通信关系。文献[10]通过在 3 种流行的匿名通信系统 (SSH、Open VPN 和 Tor) 中进行了大量的实验，验证了技术的可行性和有效性。结果表明，SSH 和 Open VPN 的检测率接近 100%，Tor 的检测率接近 95%，而假阳性率则非常低，接近 0。

由于匿名通信系统 Tor 利用单跳或多跳代理服务服务器建立匿名加密隧道来中转用户流量，目的服务器只能观察到最后一跳代理服务器的 IP 地址，使跟踪工作变得复杂。因此，发现用户和入口节点之间的通信是非常重要的。

1.3 网页流量识别

网页流量识别方法可以概括为基于端口的方法、基于有效负载的方法、基于传统机器学习的方法和基于深度学习的方法。

基于端口的方法。IANA (Internet assigned number authority) 将已知的传输层端口分配给不同的协议，基于端口的分类器^[11]简单地从包头中提取端口号值，并将其与相应的协议相关联，但端口混淆、NAT、端口转发和协议嵌入会使该方法的准确率显著降低。

基于有效负载的方法。有效负载检测技术^[12]主要通过分析报文的应用层有效负载的内容进行流量识别，但在负载加密时会使它们的有效性降低。

基于传统机器学习的方法。在过去的几十年里，将机器学习技术应用于流量统计特征以进行网络流量识别的方法得到了大量关注^[13-14]。这种方法假设流量的最大包大小、最小包大小、数据包到达间隔时间、流量持续时间等统计特征对于每个应用来说几乎是唯一的。基于这一假设，大量的机器学习方法被应用于匿名流量分类。虽然将流量统计特征与各种机器学习算法结合起来在协议或业务级别的流量分类中取得了良好的效果，但在表征不同特定应用的流模式时，它们没有鉴别力。

基于深度学习的方法。近年来，深度学习技术由于其强大的特征表示能力，在计算机视觉领域取得了巨大的成功^[15-16]。因此，一些研究者开始探索深度学习技术在流量识别领域的应用^[17-22]，如 Shen 等^[20]仅使用加密数据包长度进行细粒度网站指纹攻击，Cadena 等^[21]针对现有的深度指纹攻击提出 Tor 流量分离机制，Hardegen 等^[22]基于深度学习和真实世界流量来预测网络流量特性。

与基于传统机器学习的方法相比，基于深度学习的方法使用深度神经网络（如堆叠式自动编码器或卷积神经网络）从原始流量数据中自动学习深度特征表示。该方法最大的优点是深度特征表示直接通过深度神经网络从原始数据中提取得到，而不涉及大量的工程技能和领域专家知识。此外，利用多个堆叠的特征提取层提取的深度特征往往比流统计特征更强大。

1.4 基于Tor网页识别的指纹攻击

在正确识别Tor网页流量的基础上，现有的WF攻击方法按照使用的技术可以分为两类：一类是基于传统机器学习的WF攻击方法，如基于流量相似度匹配的攻击方法^[13,23]、基于支持向量机（SVM, support vector machine）的攻击方法^[2,14,24]、基于K-近邻的攻击方法^[3]以及基于随机森林的攻击方法^[25]；另一类是基于深度学习的WF攻击方法，如深度指纹攻击方法^[17]。传统机器学习的WF攻击方法基于网页数据包大小、包时间间隔以及数据包方向等特征，受限于技术人员知识范围，人工提取的特征并不是最具分离性的特征；基于深度学习的WF攻击方法在提取特征时忽略了人工经验信息，因此，提取具有高分离性的特征对提升Tor网页流量指纹攻击方案的性能至关重要。

基于传统机器学习的WF攻击方法的优势是能够使用较少的样本数量达到一个不错的分类效果，而基于深度学习的WF攻击方法的优势是能够自动提取有效特征达到更好的分类效果，但其训练模型的过程需要消耗大量的训练样本。从本质上说，WF攻击的根本目的在于使用分类模型识别Tor网页流量所对应的具体网站，与Tor流量识别十分类似，WF攻击同样利用了流量中存在的特征，且特征同样分为人工特征和原始流量特征。在进行WF攻击时，基于原始流量特征和深度学习模型方案能够取得很好的效果，证明了原始流量特征包含丰富的信息。

总体来说，虽然基于深度学习的方法在进行网页流量识别和Tor网页识别的指纹攻击效果更好，但计算量也相对较大，对应实际需求的时效性往往不够，更困难的是，当面对实际网络中原始流量中Tor流量占比少、训练数据不足的情形时，当前方法基本失效。

2 基于SDN的Tor流量双向统计特征

本文工作直接面向数据中心收集的原始数据

流，首先基于SDN架构提供的丰富的网络流信息，提出可以区分Tor流量与其他流量的BSF；然后经过双向统计初步筛选的Tor嫌疑流量，利用深度网络模型作为特征提取器，并为模型设计更有效的损失函数，来增强对深度特征表示的判别能力，获取更有效的LSF；最后将BSF和LSF融合，形成用于Tor网页识别的CTTF。

2.1 基于SDN的Tor流量发现机制

发现Tor流量的基础在于找出Tor网络内提供服务的节点信息，而Tor的运行机制使本文能够获得Tor网络内绝大部分节点的信息，Tor网络内的节点主要由普通节点与网桥节点组成，普通节点的信息是完全公开的，而网桥节点的信息是半公开的。对于普通节点，Tor网络由多个权威目录服务器共同维护包含所有普通节点信息的共识文件。因此通过解析共识文件，本文能够获得所有的普通节点信息。而对于网桥节点，由于其半公开的性质，本文无法一次性获取所有节点的信息。但目前很多文献^[26-27]提出了发掘Tor网络中所有网桥节点的信息方法。因此本文可以掌握Tor网络中绝大多数节点的信息，这为本文发现Tor流量奠定了基础。

传统网络是分布式的网络，没有中心控制点，数据包的控制和转发均由路由设备负责。SDN将网络元素（如路由器和交换机）的路由和转发决策与数据平面分开，控制平面仅处理与逻辑网络拓扑相关的信息，数据平面则根据控制平面中已建立的配置来协调网络流量。

由于数据平面的转发设备不具备决策能力，当SDN控制器所辖SDN域内Tor用户生成新的网络流量时，数据平面的转发设备会将其转发至控制器并由控制器及上层的应用程序决定网络流量的转发方式。因此本文可以基于SDN架构的特性对Tor流量进行全局发现。应用程序层定义Tor流量发现规则，程序依据Tor流量发现规则维护Tor节点列表，当控制器接收到数据平面的设备传来的请求时，依据Tor节点列表中的IP进行匹配。若匹配失败，则正常转发该网络流量；若匹配成功，则将该网络流量复制转发到专用的服务器。

利用SDN结构，通过人工分析，掌握Tor网络中绝大多数节点的信息之后，就可以捕获一部分Tor流量，这为后续Tor流量统计特征的分析奠定了基础。

2.2 Tor 流量 BSF 提取方法

针对 SDN 结构中捕获人工标定的 Tor 流量和原始流量的对比分析, 本文发现原始流量特征中包含更加完整以及丰富的信息, 因为分类器能够从原始流量特征中获取不同的应用协议的交互信息。而当前基于去匿名化技术标定的人工特征相较原始流量特征虽然包含的信息量有一定程度的下降, 但是直接用原始特征进行识别, 由于匿名流量在原始流量中的占比往往具有较大的波动, 这个波动对深度学习模型造成巨大的影响, 因为深度学习模型对数据的轻微变动非常敏感, 往往会让 WF 攻击等网页分析方法效果显著下降。

本文提出了一种新的 Tor 流量序列表示方法 BSF。该方法同时利用人工特征和原始流量特征, 在提供丰富信息的同时, 保证了对基于伪数据包填充防御的稳健性, 也为后续的深度特征提取提供有力的支撑。

本文通过分析网页加载过程可知, 用户与服务器建立 TCP 连接需要进行少量交互, 在用户请求网页后, 服务器会向用户发送 HTML 文件, 从而引发大量发往用户的数据包, 而在浏览器解析 HTML 文件后会再次向外界发送少量数据包请求相应的图片、视频等其他资源数据, 继而再次引发服务器向用户发送大量数据包。音乐播放、视频播放、文件下载和邮件下载虽然都是服务器向用户发送大量数据包, 但是用户短时间内只需要向服务器发出一次请求。此外, 语音通话、聊天工具具有明显的交互特征, 因此 2 个用户之间的数据包传输量类似, 而邮件上传、文件上传时用户向外发送的流量更多。

因此本文提出了 BSF 来剔除捕获的非 Tor 流量, BSF 的基本理念在于展现 Tor 流量的前 N 个数据包或 TLS 记录的交互过程蕴含的应用协议握手信息, 这将同时进一步突出 Tor 网页流量与其他流量的差异性, 具体步骤如算法 1 所示。

算法 1 BSF 提取算法

输入 流量样本集 $T = \{item_1, item_2, \dots, item_n\}$, Tor 用户 IP 地址集 G , 相关参数 TLS (样本类型标志)、 N (流量子集长度)、weight (item 大小权重标志)

输出 双向统计特征集 M

- 1) if TLS == True then
- 2) 创建由流量样本集 T 前 N 条 TLS 记录组成的子集 T_1 ;

- 3) else
- 4) 创建由流量样本集 T 前 N 条数据包组成的子集 T_1 ;
- 5) end if
- 6) 创建列表 M 存储双向统计特征;
- 7) for item _{i} in T_1 do/*遍历流量样本子集*/
- 8) $d = 1$;
- 9) if the source IP of item _{i} in G then/*判断 item _{i} 方向*/
- 10) $d = -1$;
- 11) end if
- 12) get the size i of the item _{i} ;
- 13) if weight == False then/*判断是否使用 item _{i} 大小*/
- 14) $i = 1$;
- 15) end if
- 16) if no element in M then/*
- 17) add $d * i$ to list M ;/*计算列表的第一个元素*/
- 18) else
- 19) get the last element l of list M ;
- 20) add $d * i + l$ to list M ; /*乘积与前一个元素的和*/
- 21) end if
- 22) end for

3 Tor 网页流量 LSF 特征提取

从应用类型混杂的 Tor 流量中识别出网页流量, 包含 3 个方面的基本步骤。图 2 展示了网页流量识别模型的具体架构和工作流程^[28]。

如图 2 所示, 该流程由 3 个部分组成。1) 使用大量有标签 Tor 流量数据集 D_1 训练流量特征提取模型; 2) 使用训练好的流量特征提取模型来提取少量有标签 Tor 流量数据集 D_2 中流量样本的深度特征, 并将其用于训练流量分类模型; 3) 首先使用流量特征提取模型提取无标签 Tor 流量数据集 D_3 中流量样本的深度特征, 而后输入流量分类模型来确认该样本是否属于 Tor 网页流量。与其他方案^[29-30]不同, 本文将 Tor 网页流量特征提取任务和分类任务分开处理, 基于深度学习算法的 Tor 网页流量特征提取模型能够确保提取有效的深度特征, 基于机器学习算法的 Tor 网页流量分类模型能够保证分类器的灵活更新, 从而同时满足了有效性与灵活性。

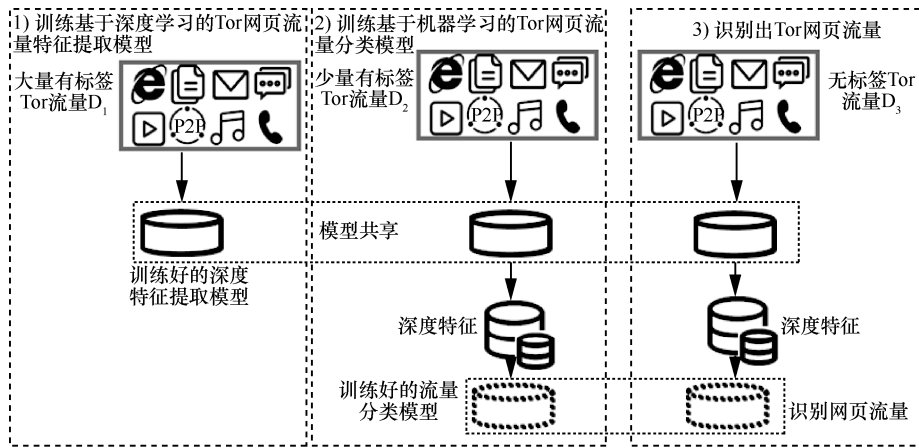


图 2 Tor 网页流量识别模型的具体架构和工作流程

本节 LSF 提取方法基于该模型设计，LSF 主要用于从经过 BSF 初步筛选的应用类型混杂的 Tor 流量中识别出网页流量，首先将深度神经网络作为特征提取器，然后使用光滑化方法对深度度量学习 (MDL, deep metric learning) 损失函数 LS Loss 进行优化，来增强网络对于深度特征表示的判别能力，最后使用该网络提取更有效的深度学习特征 LSF。

3.1 DML 模型方案

DML 依靠深度学习模型自动提取特征的能力，将原始数据映射到嵌入空间。在嵌入空间中，可以用常用的度量如欧氏距离或余弦距离评价样本之间的相似性。相较于原始空间中的样本，在嵌入空间中同类样本的相似度更高，异类样本的相似度更低，即样本在嵌入空间中更具可分离性。

模型的网络架构、样本选择策略和损失函数是 DML 领域最重要的 3 个因素。在 5.4 节实验中，本文通过实验选取了 DF^[17]作为基本网络架构。在损失函数方面，DML 领域已有众多非常成熟的损失函数，如 Contrastive Loss^[31]、Triplet Loss^[32]、Npair Loss^[33]等。

根据 Tor 网页流量的特点，本文选取 LS Loss 函数，利用其对类内样本数据结构的破坏性相对较低的特性，学习针对流量识别任务的有效特征，为后续进一步实行网站指纹攻击做准备。

本文基于小批次样本中所有的正负流量样本对来计算 LS Loss。具体来说，正流量样本指与每次所选取类型属于同一类型的流量样本，负流量样本指不同于此类型的流量样本。给定用于训练的小批次流量样本，LS Loss 定义为

$$L = \frac{1}{2|\hat{P}|} \sum_{(i,j) \in \hat{P}} \max(0, L_{i,j})^2 \quad (1)$$

其中， $L_{i,j}$ 为

$$L_{i,j} = \max(\max_{(i,m) \in \hat{N}} \alpha - D_{i,k}, \max_{(j,n) \in \hat{N}} \alpha - D_{j,l}) + D_{i,j} \quad (2)$$

其中， \hat{P} 是正流量样本对的集合， \hat{N} 是负流量样本对的集合， $D_{i,j}$ 是正流量样本对 $\{i, j\}$ 的距离。对于每一个正流量样本对 $\{i, j\}$ ，LS Loss 指导 DNN 模型分别挖掘相对 i 和 j 来说距离最近的负流量样本，假设分别为 m 和 n 。这些样本在原始空间中直接进行分类是相对困难的，因此选取它们作为困难训练样本来加速和改善模型的收敛，使通过模型将所有流量样本从原始空间映射到嵌入空间后，不同类型的样本可分离性得到提高。之后进一步通过比较 $D_{i,m}$ 和 $D_{j,n}$ ，选择两者中距离较小的负流量样本作为最难负样本，假设为 l 。最后，计算由流量样本 $\{i, j, l\}$ 确定的 Loss 来指导 DNN 模型执行反向传播。如图 3 所示，最难负流量样本针对的是正流量样本对的每一个样本。在这个包含 6 个流量样本的小批次样本中，正流量样本对中每一个样本都是独立地与所有其他负流量样本进行比较，并挖掘出最难负流量样本。

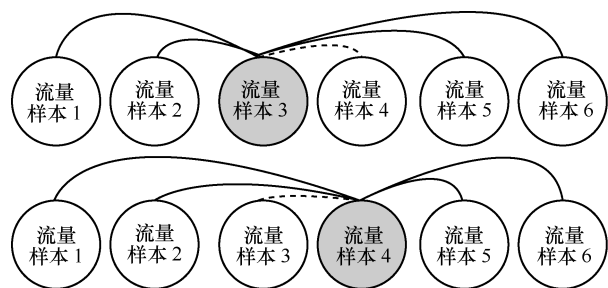


图 3 最难负流量样本挖掘方案

由于式(1)定义的 LS Loss 并不是光滑的函数,嵌套的 max 函数在实际训练过程中容易导致 DNN 模型执行反向传播时陷入局部最优的困境。因此,可以将其改进为一个光滑的函数,即

$$\tilde{L} = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{L}_{i,j})^2 \quad (3)$$

其中, $\tilde{L}_{i,j}$ 为

$$\tilde{L}_{i,j} = \log \left(\sum_{(i,m) \in N} \exp\{\alpha - D_{i,m}\} + \sum_{(j,n) \in N} \exp\{\alpha - D_{i,n}\} \right) + D_{i,j} \quad (4)$$

3.2 基于 DML 的 LSF 提取方法

神经网络不需要人工经验,可以端到端地学习到数据的高层次特征表示,从而更有效地提升流量数据的分类性能,本节详细介绍基于 DML 的 LSF 提取方法。

DML 方法将样本映射到特征空间中,通过使相同类别的特征相似性更高,不同类别的特征相似性更低,从而使不同类别的样本在特征空间中有更好的可分离性。在本文流量特征分类中,选取余弦距离来度量特征空间中不同特征的相似性,选取 LS Loss 函数来发掘类内样本数据结构的破坏性相对较低的特性,学习针对流量识别任务的有效特征,基于改进的光滑 LS Loss 函数来避免深度网络训练时陷入局部最优,具体步骤如算法 2 所示。

算法 2 LSF 深度特征提取算法

输入 流量样本集 $T = \{item_1, item_2, \dots, item_n\}$,

正流量样本 $P = \{T_{P_1}, T_{P_2}, \dots, T_{P_p}\}$, 负流量样本

$N = \{T_{N_1}, T_{N_2}, \dots, T_{N_n}\}$, 深度特征提取器 G

输出 LSF 深度特征

- 1) while not converge/*开始训练*/
- 2) for $0 \leq i < P_p$ /*遍历正流量样本子集*/
- 3) for $0 \leq j < N_n$ /*遍历负流量样本子集*/
- 4) $D_{ij} = \text{Cosine Distance}(G(T_{P_i}), G(T_{N_j}))$ /*计算特征空间正、负流量样本间距离*/
- 5) end for
- 6) end for
- 7) for $0 \leq k < N_n$ /*遍历负流量样本子集*/
- 8) for $0 \leq i < P_p$ /*遍历正流量样本子集*/
- 9) for $0 \leq j < P_p$ /*遍历剩余正流量样本子集*/

- 10) $L = \underset{k}{\operatorname{argmin}}(D_{ik}, D_{jk})$ /*距离正样本对 (i, j) 最近的负样本序号*/
- 11) end for
- 12) end for
- 13) end for
- 14) 基于改进的式(3)、式(4)计算 LS Loss
- 15) 使用随机梯度下降算法最小化 LS Loss 训练深度神经网络
- 16) end while

4 CTF 复合特征与流量数据增强

4.1 CTF 特征和网页特征识别流程

SDN 架构下,基于实际流量的 Tor 网页流量分类分为 2 个层次进行,如图 4 所示,第一个层次针对 SDN 数据中心采集到的原始流量,提取原始流量的 BSF,将原始流量分为 Tor 流量和非 Tor 流量,为后续进一步细粒度分类奠定基础。

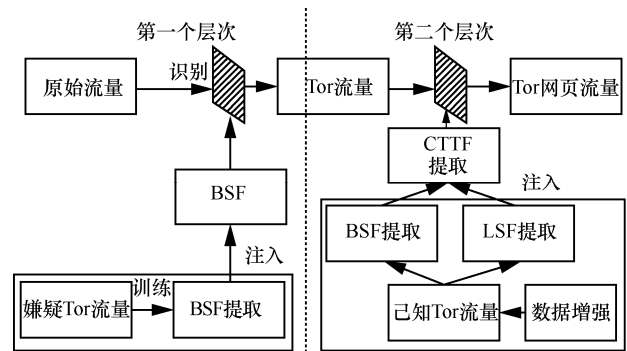


图 4 CTF 特征和网页特征识别流程

第二个层次是将上述 Tor 流量中的网页流量识别出来。针对已知的 Tor 流量,本文通过对 Tor 流量提取相应的 BSF,基于 LS Loss 训练提取 LSF,二者融合获得 Tor 流量的 CTF,再基于标记的训练数据训练 KNN 分类器,识别 Tor 流量数据中的网页流量。

针对训练数据稀缺的问题,本文提出基于平移的流量数据增强方法,对增强后的流量提取 BSF 和 LSF 得到 CTF。基于上述得到的特征,本文提出的流量数据增强方法也能有效解决训练数据稀缺的问题,提升 Tor 网页流量分类准确率。

4.2 基于平移的流量数据增强方法

通过对 Tor 流量识别任务的观察,本文提出了一种基于平移的流量增强方法,该方法利用移位的

数据增强手段来弥补可能遇到的训练样本不足的问题。数据增强^[34]是计算机视觉中常用的减少过拟合问题的方法，通过数据增强提高训练数据规模和质量，可以为机器学习算法提供容量更大、更接近真实分布的训练数据，令训练集和实际需要识别的数据之间的分布尽量匹配。

因此本文探索使用小数据集进行实验，最终选取离线增强方式进行数据扩充，具体的Tor数据增强策略算法的描述如算法3所示。

算法3 左右平移的数据增强算法

输入 流量样本集 $T = \{item_1, item_2, \dots, item_n\}$ ，每条流开始收集时间集 start，每条流结束收集时间集 end

输出 数据增强后的流量样本集 T

- 1) for $item_i$ in T do/*遍历流量样本子集*/
- 2) $s_time_i =$ timestamp of the first packet of $item_i$
- 3) $e_time_i =$ timestamp of the last packet of $item_i$
- 4) if $s_time_i < first_i$ then/*判断 $item_i$ 起始位置是否提前*/
- 5) for chi_item_j in $item_i$ do/*从 $item_i$ 前向后遍历*/
- 6) if timestamp of $chi_item_j \geq s_time_i$ /*判断 $item_i$ 中最接近流量开始收集时间的数据包时间戳*/
- 7) if $j \geq 100$ then/*限制删除流 $item_i$ 最开始的数据包个数不超过100个*/
- 8) remove first random(1,100) packets of $item_i$
- 9) else
- 10) remove first random(1, j) packets of $item_i$
- 11) end if
- 12) end if
- 13) break
- 14) end for
- 15) end if
- 16) if $e_time_i > end_i$ then/*判断 $item_i$ 结束位置是否推后*/
- 17) for chi_item_k in reverse of $item_i$ do/*从 $item_i$ 后向前遍历*/
- 18) if timestamp of $chi_item_k \leq e_time_i$ /*判断 $item_i$ 中最接近流量结束收集时间

的数据包时间戳*/

- 19) if $k \geq 100$ then
- 20) remove last random(1,100) packets of $item_i$ /*限制删除流 $item_i$ 最后面的数据包个数不超过100个*/
- 21) add random(1,100) random direction packets to the front of $item_i$
- 22) else
- 23) remove last random(1, k) packets of $item_i$
- 24) add random(1, k) random direction packets to the front of $item_i$
- 25) end if
- 26) end if
- 27) break
- 28) end for
- 29) end if
- 30) end for

5 实验评估

本文首先在SDN仿真环境中验证了第2节提出的Tor流量发现机制的有效性，基于收集的流量数据验证了基于CTTF的Tor网页流量识别方法的有效性。在此基础上，本文研究了流量增强机制对Tor网页流量识别效果的影响。

5.1 Tor流量发现机制有效性验证

为了验证了本文提出的Tor流量发现机制的有效性，本文首先基于物理机器搭建了具备4个节点的私有Tor网络，在SDN仿真环境Mininet^[35]中进行了实验。本文基于Mininet搭建了具有2个交换机和2个PC的数据平面，并使用RYU作为SDN控制器，在RYU中以应用程序的形式实现本文所需的Tor流量发现功能，实验中使用如图1所示的网络结构拓扑。

首先，应用程序周期性地向权威目录服务器请求获取Tor网络内的节点信息。当Tor用户访问Tor网络时生成新的网络流量数据，则SDN交换机会向RYU控制器请求该网络流量数据的处理方法。控制器经过对比发现该网络流量数据为Tor流量，则在正常转发该流量的前提下将该网络流量数据复制转发到服务器进行保存，以完成对Tor流量的发现。接下来，令Tor用户不通过Tor网络向网站上传文件，20s后再通过Tor网络向网站上传文件，

在此过程中审查者记录收集到的 Tor 流量。

如图 5 所示, 在前 20 s 虽然 Tor 用户并未使用 Tor 进行传输数据, 但是审查者仍能够捕获少量 Tor 流量, 这是因为 Tor 用户在维护与入口节点之间的链路, 而 20 s 后 Tor 用户使用 Tor 进行通信, 此时审查者能够捕获大量的 Tor 流量。

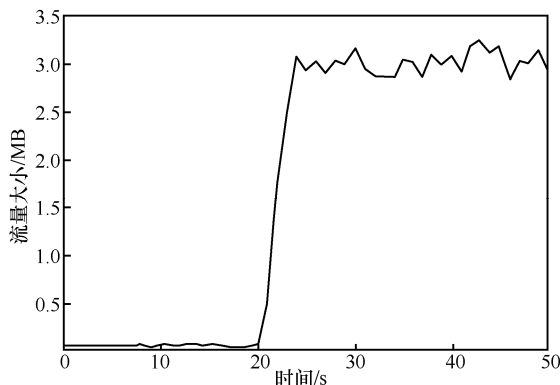


图 5 不同使用状态下的 Tor 流量大小

5.2 网页流量识别机制有效性验证

为了验证本文提出的网页流量识别机制的有效性, 本文收集了 8 种类型的流量用于实验, 包括网页流量、在线音乐、网络视频、语音通话、即时聊天、文件传输、电子邮件和 P2P。

由于本文的主要任务是识别网页流量, 因此本文将网页流量与其他流量的比例调整为 1:1。对于网页类型的 Tor 流量, 本文将 Alexa 网页排名榜单中排行前 2 000 的网页作为研究对象。同时使用火狐和谷歌 2 种浏览器访问网页, 每种浏览器获取每个网页的 5 个 Tor 流量样本, 最终获得了 20 000 条网页类型的 Tor 流量样本。对于其他 7 种类型的流量, 每个类型抓取 3 000 个 Tor 流量样本, 共 21 000 个样本。

本文将所有类型的流量样本均分为 2 个部分, 即 Dataset1 数据集和 Dataset2 数据集, 每个数据集包含 20 500 条流量。Dataset1 被 LS Loss 用来指导本文深度特征提取器的训练, Dataset2 流量用于各方案分别提取特征来训练分类器和测试其识别效果。

本文将 Wang 等^[29]提出的方案命名为 J48, Lotfollahi 等^[30]中提出的方案命名为 DP。J48 方案使用 Tranalyzer2 提取 Tor 流量中的 79 种特征, 并结合 J48 分类器得到了较好的流量识别效果。本文将 J48 提出的特征定义为 F1, 通过 F1 对 Dataset2 处理后得到特征集 Dataset2_{F1}, 分别在 Dataset2_{F1} 上

训练 J48、DP 分类器, 并对它们的识别效果进行评估。本文将 3.2 节提出的 BSF 定义为 F2, 通过 F2 对 Dataset2 进行处理后得到特征集 Dataset2_{F2}。利用 LS Loss 指导 DNN 模型在 Dataset1 上训练得到 LSF 提取器 E1, 利用 E1 对 Dataset2 进行提取特征进一步得到深度特征集 Dataset2_{E1}, 将 Dataset2_{F2} 和 Dataset2_{E1} 融合得到 Dataset2 的 CTTF, 用于训练 KNN 分类器, 并对其识别效果进行评估。通过控制参与训练的流量数据占比, 实验分别得出了 J48、DP 以及 CTTF 方案的准确率、真正正确率及假正确率随训练数据占比变化的曲线。

网页流量识别效果对比如图 6 所示。首先随着训练数据占比的增加, 3 种模型的识别效果都呈从快速提升到缓慢提升的变化趋势, 其中, 本文提出的 CTTF 方案相较其他 2 种方案在数据量更少时识别效果提升更加明显, 在 40% 的数据参与训练的情况下, CTTF 方案已经达到了收敛。然而, 基于手工特征集和机器学习分类算法的 J48 方案并未达到收敛。基于深度学习模型的 DP 方案由于训练数据量不足, 识别效果最差。图 6 实验结果显示了本文方案能够通过提取更加稳健、有效的流量特征使分类器以更少的训练数据达到更好的识别效果。随着参与训练的数据量进一步增大, 当 80% 的数据参与训练时, 巨大的训练数据量使基于深度学习模型的 DP 方案识别效果进一步提升, DP 方案与本文提出的 CTTF 方案的识别率达到同一水平。然而由于本文设计的 CTTF 方案为负样本对设计了更高的损失值权重, 误报更少, FPR 相较 DP 方案更低, 在实际应用中其他类型流量被误识别为网页流量的概率更小。

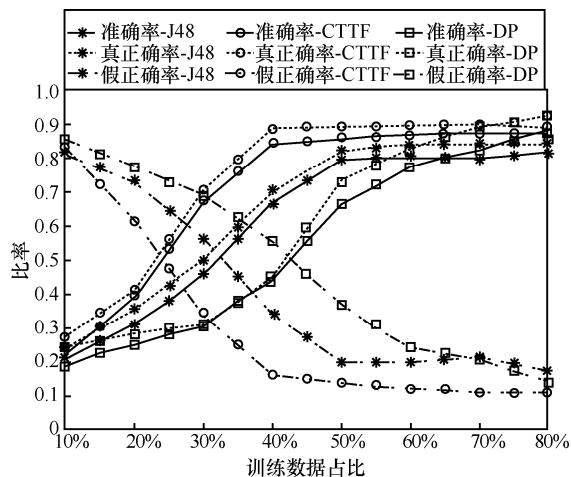


图 6 网页流量识别效果对比

5.3 数据增强机制有效性验证

为了验证数据增强机制的有效性，本节对本文提出的方案在数据增强机制辅助下的识别效果进行了评估。由于审查者很难掌握 Tor 流量的具体位置，但是仍能大致判断 Tor 流量的起始位置。因此，本文对流量数据的平移操作被限制到了 100 个数据包的范围，通过左右平移流量形成新的流量数据。本文采用随机平移的方式对涉及训练特征提取器和分类器的数据进行在线扩充至原来的 5 倍。

如图 7 所示，本文将未采用数据增强机制的方法称为 CTF，只在测试数据上执行数据增强的方法称为 Non，在训练数据和测试数据同时进行数据增强的方法称为 Arg。当参与训练的数据占比超过 40% 时，经过数据增强后分类的准确率和 TPR 都有小幅提升，FPR 有小幅下降，但是均无明显改善。当参与训练的数据占比小于 40% 时，经过数据增强后分类的准确率和 TPR 都有较大幅度的提升，FPR 也有较大幅度的下降。当训练数据较少时，通过平移操作对 Tor 流量进行数据增强，能够使训练数据与实际需要识别的数据更加相似，从而提高分类器的识别效果。当训练数据充足时，数据增强虽然会少量增加分类器的识别效果，但与此同时也会增加训练分类器的时间消耗，因此，这种情况下需要谨慎采用数据增强技术。

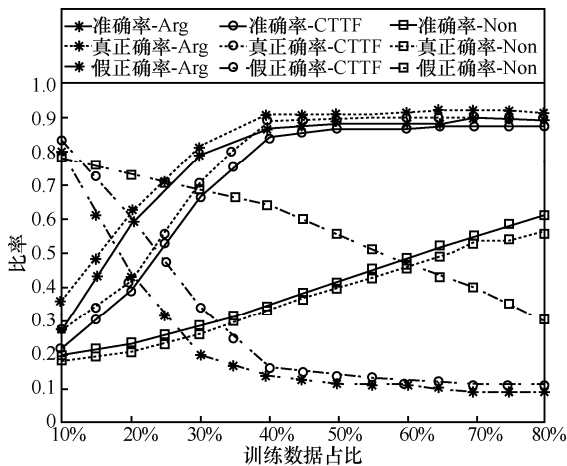


图 7 数据增强机制的效果

5.4 模型超参数的选取

通常情况下，模型训练的 epoch 越多，预测效果越好。使用不同 DNN 子网络模型的性能曲线如图 8 所示。DF 作为子网络时，在 10 个 epoch 时就达到 82% 的精确率和 85% 的召回率。随着

epoch 不断增加，模型性能逐渐变好，最终在 40 个 epoch 后趋于平稳。本文发现，DF 作为子网络时模型在训练 epoch 最少的情况下就能达到收敛，性能仍表现最好，在流量识别的任务中再次展现了强大的特征提取能力。因此，本文最终选择 DF 作为 LSF 提取模型中的基础网络，且模型的训练 epoch 为 40。

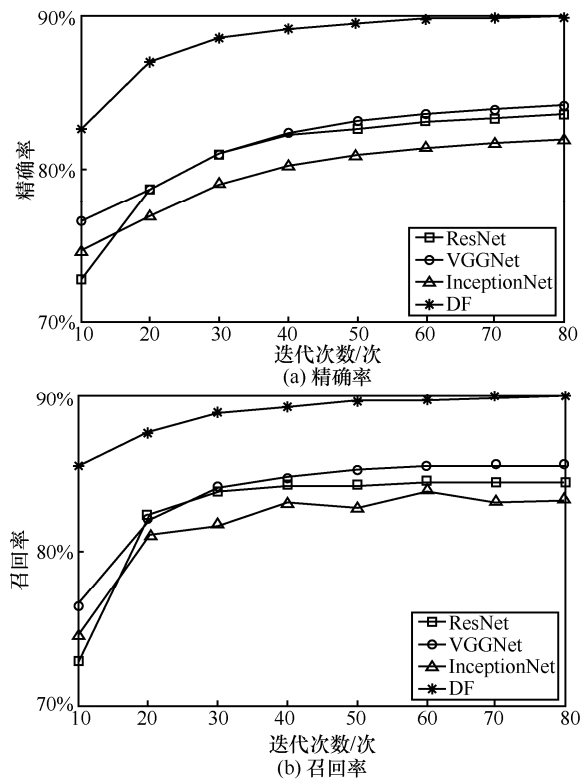


图 8 使用不同 DNN 子网络模型的性能曲线

通过对比测试，本文选择余弦距离作为 LSF 提取模型中的距离度量，依靠其独有的特性来衡量流量样本间的相似度。此外，Adam 优化器^[28]结合了 SGDM^[36]的一阶动量和 RMSProp^[37]的二阶动量，在梯度下降的过程中加入了惯性，并实现了自适应的学习率调整。因此，本文选择 Adam 作为 LSF 提取模型训练过程中的优化器，以获得更佳的性能。

通过基于 DML 训练得到的特征提取器，将流量样本在高维空间中的原始稀疏特征向量映射到低维空间中稠密特征向量后，同类流量样本间距离减少，异类流量样本间距离增加。由于由 DML 指导训练得到的深度学习模型通常配合 KNN 一起完成分类任务，因此本文同样采用 KNN 作为目标分类器，以 BSF 和 LSF 二者融合

获得的 CTTF 特征训练 KNN。与普通的分类算法不同, KNN 通过在欧氏空间中从参与训练分类器的流量样本中找出与需要预测的样本距离最接近的 K 个流量样本, 然后在它们当中找出某一标签对应样本数量最多的标签作为预测样本的类别, 该思想类似于“投票法”。 K 值的选取在一定程度上会影响 KNN 模型的预测结果。如果选择较小的 K 值, 只有与输入流量样本较近或相似的训练样本才会对预测结果起作用, 但意味着模型容易发生拟合; 如果选择较大的 K 值, 就相当于利用较多训练时的流量样本进行预测, 但此时与输入流量样本不相似的训练样本也会对预测结果起作用, 导致预测错误。通过交叉验证, 本文对 KNN 主要采取以下超参数指导分类器进行训练和预测。

- 1) 采用余弦距离来衡量样本间的相似度。
- 2) 每个参与训练样本的权重设置为与距离成反比。
- 3) 每一次查找 10 个与输入流量样本最近的样本, 也就是 K 为 10。

针对 K 的取值问题, 本文采用测试多个 K 值, 并从中选取使分类性能达到最佳的值作为最终的 K 值, 如图 9 所示, 当 K 取值为 9、10、11 时, 模型取得最高的准确率, 本文选取中间值 10 作为 K 的最终取值。

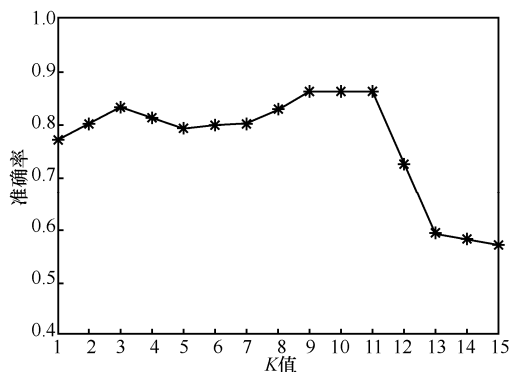


图 9 使用不同 K 值对识别准确率的影响

6 结束语

本文基于 SDN 架构特性, 在数据中心环境下发现和收集 Tor 流量, 据此分析获得了 Tor 流量的 BSF; 进一步基于卷积网络提取深度特征的流量特征表示, 使用 LS Loss 训练深度特征提取

器, 获得了基于深度度量学习的 Tor 网页 LSF; 针对训练数据不足的情形, 提出了流量数据增强的方法。实验表明, 基于 BSF+LSF 的复合特征 CTTF 能针对实验模拟的数据中心的原始流量进行 Tor 网页流量识别, 相比现有方法, 识别率提升了 4%, 达到 85.9%, TPR 达到 88.7%, FPR 降至 10.9%; 其中基于 Tor 流量数据进行增强的方法可将分类器的分类效果在训练数据较少的情况下得到有效提升。

参考文献:

- [1] HERRMANN D, WENDOLSKY R, FEDERRATH H. Website fingerprinting: attacking popular privacy enhancing technologies with the multinomial naïve-Bayes classifier[C]//Proceedings of the 2009 ACM Workshop on Cloud Computing Security. New York: ACM Press, 2009: 31-42.
- [2] PANCHENKO A, NIESSEN L, ZINNEN A, et al. Website fingerprinting in onion routing based anonymization networks[C]//Proceedings of the 10th Annual ACM Workshop on Privacy in the Electronic Society. New York: ACM Press, 2011: 103-114.
- [3] WANG T, CAI X, NITHYANAND R, et al. Effective attacks and provable defenses for website fingerprinting[C]//23rd USENIX Security Symposium. Berkeley: USENIX Association, 2014: 143-157.
- [4] SCOTT-HAYWARD S, O'CALLAGHAN G, SEZER S. SDN security: a survey[C]//Proceedings of 2013 IEEE SDN for Future Networks and Services. Piscataway: IEEE Press, 2013: 1-7.
- [5] 魏松杰, 孙鑫, 赵茹东, 等. SDN 中 IP 欺骗数据分组网络溯源方法研究[J]. 通信学报, 2018, 39(11): 181-189.
- [6] WEI S J, SUN X, ZHAO R D, et al. Tracing IP-spoofed packets in software defined network[J]. Journal on Communications, 2018, 39(11): 181-189.
- [7] SONG H O, XIANG Y, JEGELKA S, et al. Deep metric learning via lifted structured feature embedding[C]//Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition. Piscataway: IEEE Press, 2016: 4004-4012.
- [8] 祝现威, 常朝稳, 朱智强, 等. 基于身份属性的 SDN 控制转发方法[J]. 通信学报, 2019, 40(11): 1-18.
- [9] ZHU X W, CHANG C W, ZHU Z Q, et al. SDN control and forwarding method based on identity attribute[J]. Journal on Communications, 2019, 40(11): 1-18.
- [10] BENZEKKI K, FERGOUGUI A E, ELALAOUI A E. Software-defined networking (SDN): a survey[J]. Security and Communication Networks, 2016, 9(18): 5803-5833.
- [11] OCONNOR T, ENCK W, PETULLO W M, et al. PivotWall: SDN-based information flow control[C]//Proceedings of the Symposium on SDN Research. [S.l.:s.n.], 2018: 1-14.
- [12] LING Z, LUO J Z, XU D N, et al. Novel and practical SDN-based traceback technique for malicious traffic over anonymous networks[C]//Proceedings of IEEE INFOCOM 2019 - IEEE Conference on Computer Communications. Piscataway: IEEE Press, 2019: 1180-1188.
- [13] MOORE A W, PAPAGIANNAKI K. Toward the accurate identification of network applications[C]//International Workshop on Passive and Active Network Measurement, Berlin: Springer, 2005: 41-54.
- [14] FINSTERBUSCH M, RICHTER C, ROCHA E, et al. A survey of

- payload-based traffic classification approaches[J]. IEEE Communications Surveys & Tutorials, 2014, 16(2): 1135-1156.
- [13] WANG T, GOLDBERG I. Improved website fingerprinting on tor[C]//Proceedings of the 12th ACM Workshop on Privacy in the Electronic Society. New York: ACM Press, 2013: 201-212.
- [14] CAI X, ZHANG X C, JOSHI B, et al. Touching from a distance: website fingerprinting attacks and defenses[C]//Proceedings of the 2012 ACM Conference on Computer and Communications Security. New York: ACM Press, 2012: 605-616.
- [15] HE K M, ZHANG X Y, REN S Q, et al. Deep residual learning for image recognition[C]//Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition. Piscataway: IEEE Press, 2016: 770-778.
- [16] ZHU J Y, PARK T, ISOLA P, et al. Unpaired image-to-image translation using cycle-consistent adversarial networks[C]//Proceedings of 2017 IEEE International Conference on Computer Vision. Piscataway: IEEE Press, 2017: 2242-2251.
- [17] SIRINAM P, IMANI M, JUAREZ M, et al. Deep fingerprinting: undermining website fingerprinting defenses with deep learning[C]//Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2018: 1928-1943.
- [18] RIMMER V, PREUVENEERS D, JUAREZ M, et al. Automated website fingerprinting through deep learning[C]//Proceedings of 2018 Network and Distributed System Security Symposium. Reston: Internet Society, 2018: 1-16.
- [19] BHAT S, LU D, KWON A, et al. Var-CNN: a data-efficient website fingerprinting attack based on deep learning[J]. Proceedings on Privacy Enhancing Technologies, 2019, 2019(4): 292-310.
- [20] SHEN M, LIU Y T, ZHU L H, et al. Fine-grained webpage fingerprinting using only packet length information of encrypted traffic[J]. IEEE Transactions on Information Forensics and Security, 2021, 16: 2046-2059.
- [21] CADENA W D L, MITSEVA A, HILLER J, et al. TrafficSliver: fighting website fingerprinting attacks with traffic splitting[C]//Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security. New York: ACM Press, 2020: 1971-1985.
- [22] HARDEGEN C, PFÜLB B, RIEGER S, et al. Predicting network flow characteristics using deep learning and real-world network traffic[J]. IEEE Transactions on Network and Service Management, 2020, 17(4): 2662-2676.
- [23] SHI Y, MATSUURA K. Fingerprinting attack on the tor anonymity system[C]//International Conference on Information and Communications Security. Berlin: Springer, 2009: 425-438.
- [24] PANCHENKO A, LANZE F B, ZINNEN A, et al. Website fingerprinting at Internet scale[C]//Proceedings of 2016 Network and Distributed System Security Symposium. Reston: Internet Society, 2016: 1-15.
- [25] HAYES J, DANEZIS G. K-fingerprinting: a robust scalable website fingerprinting technique[J]. arXiv Preprint, arXiv: 1509.00789, 2015.
- [26] MATIC S, TRONCOSO C, CABALLERO J. Dissecting tor bridges: a security evaluation of their private and public infrastructures[C]//Proceedings of 2017 Network and Distributed System Security Symposium. Reston: Internet Society, 2017: 1-15.
- [27] LING Z, LUO J Z, YU W, et al. Extensive analysis and large-scale empirical evaluation of tor bridge discovery[C]//2012 Proceedings of IEEE INFOCOM. Piscataway: IEEE Press, 2012: 2381-2389.
- [28] KINGMA D P, BA J. Adam: a method for stochastic optimization[J]. arXiv Preprint, arXiv: 1412.6980, 2014.
- [29] WANG L M, MEI H T, SHENG V S. Multilevel identification and classification analysis of tor on mobile and PC platforms[J]. IEEE Transactions on Industrial Informatics, 2021, 17(2): 1079-1088.
- [30] LOTFOLLAHI M, SIAVOSHANI M J, ZADE R S H, et al. Deep packet: a novel approach for encrypted traffic classification using deep learning[J]. Soft Computing, 2020, 24(3): 1999-2012.
- [31] CHOPRA S, HADSELL R, LECUN Y. Learning a similarity metric discriminatively, with application to face verification[C]//Proceedings of 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. Piscataway: IEEE Press, 2005: 539-546.
- [32] SCHROFF F, KALENICHENKO D, PHILBIN J. FaceNet: a unified embedding for face recognition and clustering[C]//Proceedings of 2015 IEEE Conference on Computer Vision and Pattern Recognition. Piscataway: IEEE Press, 2015: 815-823.
- [33] SOHN K. Improved deep metric learning with multiclass n-pair loss objective[C]//Advances in Neural Information Processing Systems. [S.l.:s.n.], 2016: 1857-1865.
- [34] PEREZ L, WANG J. The effectiveness of data augmentation in image classification using deep learning[J]. arXiv Preprint, arXiv: 1712.04621, 2017.
- [35] 兰巨龙, 张学帅, 胡宇翔, 等. 基于深度强化学习的软件定义网络QoS优化[J]. 通信学报, 2019, 40(12): 60-67.
- LAN J L, ZHANG X S, HU Y X, et al. Software-defined networking QoS optimization based on deep reinforcement learning[J]. Journal on Communications, 2019, 40(12): 60-67.
- [36] SUTSKEVER I, MARTENS J, DAHL G, et al. On the importance of initialization and momentum in deep learning[C]//International Conference on Machine Learning. [S.l.:s.n.], 2013: 1139-1147.
- [37] ZEILER M D. ADADELTA: an adaptive learning rate method[J]. arXiv Preprint, arXiv: 1212.5701, 2012.

[作者简介]



言洪萍(1985-), 男, 江苏常州人, 江苏大学博士生, 主要研究方向为机器学习、匿名流量分析。

周强(1992-), 男, 安徽安庆人, 博士, 江苏大学讲师, 主要研究方向为机器学习、匿名流量分析。

王世豪(1996-), 男, 山东济南人, 江苏大学硕士生, 主要研究方向为机器学习、匿名流量分析。

姚旺(1997-), 男, 江苏滨海人, 江苏大学硕士生, 主要研究方向为机器学习、匿名流量分析。

何刘坤(1996-), 男, 安徽潜山人, 江苏大学硕士生, 主要研究方向为机器学习、匿名流量分析。

王良民(1977-), 男, 安徽潜山人, 博士, 江苏大学教授、博士生导师, 主要研究方向为密码学与安全协议、物联网安全、大数据安全。