

基于时频特征融合的伪造语音检测算法

袁程胜^{1,2}, 张雪原^{2,3}, 周志立⁴, 李欣亭⁵, 付章杰^{1,2}

(1. 南京信息工程大学计算机学院、网络空间安全学院, 南京 210044; 2. 南京信息工程大学数字取证教育部工程研究中心, 南京 210044; 3. 南京信息工程大学软件学院, 南京 210044; 4. 广州大学人工智能研究院, 广州 510006; 5. 国防科技大学外国语学院, 南京 210039)

摘要: 针对伪造语音检测精度不高和泛化性弱的难题, 提出一种基于时频特征融合的伪造语音检测算法。首先, 为了挖掘语音片段能量分布不均、基频波动异常, 以及提取语义连贯性的细微差别, 提出一种多分支特征融合网络, 分别从音高、音强以及能量分布来挖掘真假语音的差异痕迹, 以更好地表征真假语音的频率变化、振幅变化和峰值差异, 提高伪造语音检测的准确率。其次, 经典的坐标注意力机制未能对语音时频域的细粒度差异进行有效挖掘, 为此提出一种时频坐标注意力机制, 分别从时域和频域两个方向对能量分布和基频波动异常进行联合编码, 以更好地表征频谱图中的共性高频能量异常, 提升模型的泛化性。最后, 设计了一种自适应联合损失优化函数, 通过平衡不同分支网络的权重, 进一步提升模型对伪造语音中高频能量异常及语义不连贯性的学习能力。在 ASVspoof 2019 逻辑访问数据集上进行了性能评估, 实验结果表明, 与现有的工作相比, 所提方法在等错误率 (Equal error rate, EER) 和最小归一化串联检测代价函数 (Minimum normalized tandem detection cost function, min t-DCF) 两个指标上均取得较好性能, 分别降低了 0.34% 和 0.014。此外, 在应对极难检测的未知攻击 A17 时, 同样展现出较高的泛化性, 其中 EER 和 min t-DCF 分别下降了 3.952 2% 和 0.136 4。当应对未知类型的欺骗攻击时, 同样表现出较好的泛化性。

关键词: 伪造语音检测; 特征融合; 时频特征融合; 频谱图; 语音波形

中图分类号: TN912

文献标志码: A

Forged Speech Detection Algorithm Based on Time-Frequency Feature Fusion

YUAN Chengsheng^{1,2}, ZHANG Xueyuan^{2,3}, ZHOU Zhili⁴, LI Xinting⁵, FU Zhangjie^{1,2}

(1. School of Computer Science, School of Cyber Science and Engineering, Nanjing University of Information Science and Technology, Nanjing 210044, China; 2. Engineering Research Center of Digital Forensics, Ministry of Education, Nanjing University of Information Science and Technology, Nanjing 210044, China; 3. School of Software, Nanjing University of Information Science and Technology, Nanjing 210044, China; 4. Institute of Artificial Intelligence, Guangzhou University, Guangzhou 510006, China; 5. School of Foreign Languages, National University of Defense Technology, Nanjing 210039, China)

Abstract: To solve the problem of low accuracy and weak generalization of forged speech detection, a new algorithm based on time-frequency feature fusion is proposed. Firstly, in order to excavate the uneven energy distribution of speech fragments or the abnormal fundamental frequency fluctuation, and extract the subtle difference of semantic coherence, a multi-branch feature fusion network is proposed to excavate the

基金项目: 国家自然科学基金(U22B2062, U23B2023, 62102189); 南京市重大科技专项(202405002)。

收稿日期: 2024-10-13; **修订日期:** 2024-12-26

difference traces of true and false speech from the pitch, pitch intensity and energy distribution respectively, so as to better represent the frequency change, amplitude change and peak difference of true and false speeches, and improve the accuracy of forged speech detection. Secondly, the classical coordinate attention mechanism fails to effectively mine the fine-grained differences in the time-frequency domain of speech. Therefore, a time-frequency coordinate attention mechanism is proposed to jointly encode the energy distribution and fundamental frequency fluctuation anomalies from the time domain and the frequency domain respectively, so as to better characterize the common high frequency energy anomalies in the spectral graph and improve the generalization of the model. Finally, an adaptive joint loss optimization function is designed to balance the importance of different branch networks to further improve the model's ability to learn high frequency energy anomalies and semantic incoherence in forged speech. Performance is evaluated on the logical access (LA) dataset of ASVspoof 2019, and experimental results show that compared with the current methods, the proposed method achieves good performance in both EER (Equal error rate) and mint-DCF (Minimum normalized tandem detection cost function) indicators, which decrease by 0.34% and 0.014, respectively. In addition, when dealing with unknown attack A17, which is extremely difficult to detect, it also show good generalization, where EER and mint-DCF decrease by 3.952 2% and 0.136 4, respectively. When dealing with unknown types of spoofing attacks, it also shows better generalization.

Key words: forged speech detection; feature fusion; time frequency feature fusion; spectrogram; speech waveform

引言

移动互联网和智能穿戴设备的普及给人们的日常生活带来了诸多便捷。但近些年,因敏感内容保护不当而导致的数据泄露、隐私窃取问题愈发凸显,对个人信息的安全维护构成了严峻挑战。因此,验证用户身份的真实性和合法性成为保障用户隐私安全的首要防线。相较于指纹、人脸和行为等各类生物特征识别技术,语音识别技术凭借其非接触式交互、低认知负荷及场景普适性的特性,深受大众青睐^[1-2]。尤其是新冠疫情期间,其“零接触”特性更使其成为公共健康安全与数字身份核验平衡的优选技术。然而,最新的研究成果揭示,语音识别技术正遭受严峻的伪造语音欺诈攻击挑战。特别是,借助深度伪造技术,如今已能够生成在频谱特性、情感韵律乃至呼吸节奏方面都与真人无异的语音样本。更严峻的是,攻击者正将目标从通用语音库转向特定个体声纹克隆,仅需少量目标语音样本,就能构建出高精度的声纹模型,这使得“听觉欺骗”从技术可能演变为现实威胁。因此,鉴别待认证的语音是否源自自然人,对于提升语音识别系统的安全性和可靠性至关重要。鉴于伪造语音欺骗攻击的严峻态势,伪造语音检测技术应运而生。然而,如何切实有效地提升伪造语音检测任务的准确率与泛化能力,依旧是当前研究人员亟待解决的关键难题。

目前,伪造语音技术主要分为传统式和深度学习式两大类。前者主要依赖于波形拼接和参数统计两种方法。但是,该类方法伪造的语音时常面临韵律信息丢失(如重音、音调、节奏和停顿等)、噪声干扰以及信号失真等问题,使其在自然度和流畅度上表现不佳,且易被人耳辨识。而后者则是通过设计定制化的伪造模型,自动生成高度逼真的合成语音。为了抵御伪造语音的欺骗攻击,提升语音识别系统的安全性和可靠性,一系列伪造语音检测技术相继被提出。通过对伪造语音检测技术梳理归纳后发现,现有方法大多仅对单一频谱或波形特征进行分析,忽略了伪造语音片段间的语音连贯性,以及能量分布不均、基频波动异常情况,致使现有方法在检测精度方面仍然不高。此外,伪造语音在生成过程中

为了提高与原声的相似度,会将多个声道的语音叠加,从而导致声谱分布连续性较差,并且在高频区域容易出现异常的波峰。

本文受多模态特征融合思想的启发,提出一种基于时频特征融合的伪造语音检测方法,分别从音高、音强以及能量分布信息3个角度挖掘伪造语音中的差异信息,以更好地表征真假语音中的频率变化、振幅变化和峰值差异等表征真伪语音本质差异的信息。本文的主要贡献如下:

(1) 为了挖掘语音片段能量分布不均、基频波动异常,以及提取语义连贯性的细微差别,提出一种多分支特征融合网络(Multi-branch feature fusion network, MFFNet),分别从音高、音强以及能量分布信息角度挖掘语音中的高频能量异常以及语义不连贯信息,以更好地表征真假语音的频率变化、振幅变化和峰值差异。

(2) 经典的坐标注意力机制对语音的时频域处理能力较弱,未能提取频谱图中的细粒度伪造痕迹,为此提出一种时频坐标注意力机制(Time-frequency coordinate attention mechanism, TCAM),通过对时域和频域上的能量分布和基频波动异常进行联合编码,以更好表征真假语音的细微差异。

(3) 为了评估不同分支特征学习网络的重要性,剔除多分支中的冗余特征,还设计了一种自适应联合损失优化函数(Adaptive loss optimization function, ALOF),通过融合交叉熵损失和Multimargin-Loss,以平衡不同分支网络的权重,通过联合训练来挖掘不同分支下的高频能量共性特征。

(4) 在ASVspoof 2019 逻辑访问(Logical access, LA)数据集上进行了性能测试,实验结果显示,所提方法相较于现有工作,在等错误率(Equal error rate, EER)和最小归一化串联检测代价函数(Minimum normalized tandem detection cost function, min t-DCF)两项指标上分别为0.35%和0.011,表现出较好的检测精度。此外,在应对未知攻击A17时同样表现出极高的泛化性。

1 相关工作

通过对现有的工作梳理后发现,伪造语音检测技术主要分为4类:传统检测方法、基于深度特征学习的检测方法、基于端到端的检测方法和基于特征融合的检测方法。本节将对上述方法进行详细介绍。

1.1 传统检测方法

传统检测方法首先从原始语音信号中提取频谱特征,然后深入分析这些频谱特征,挖掘出频率的细微变化、峰值差异等关键信息,最后利用分类器对上述特征进行建模与测试。其中最具有代表性的频谱特征主要包含3类,即Mel频率倒谱系数(Mel frequency cepstral coefficients, MFCC)、常数Q倒谱系数(Constant Q cepstral coefficients, CQCC)和线性频率倒谱系数(Linear frequency cepstral coefficients, LFCC)。基于此,Patel等^[3]将其应用于语音识别任务中,提出了一种基于帧级特征的伪造语音检测方法,通过提取Mel频率倒谱系数、耳蜗滤波器倒谱系数(Cochlear filter cepstral coefficients, CFCC)以及耳蜗滤波器倒谱系数与瞬时频率变化的结合特征(Cochlear filter cepstral coefficients and change in instantaneous frequency, CFCC-IF)来捕捉跨帧的特征变化,降低噪声的干扰,实验结果表明,该方法在检测准确率上取得了不错的性能。在2016年举办的BTAS反欺骗语音挑战比赛中,Korshunov等^[4]提出了一种基于长期平均频谱和标准偏差信息的伪造语音检测方法,利用线性判别分析(Linear discriminant analysis, LDA)模型充分学习通道退化和语音自然性等信息,最终在鲁棒性测试中表现出不错的性能。Sahidullah等^[5]使用最大似然高斯混合模型(Gaussian mixture model with maximum likelihood, GMM-ML)和局部二值模式支持向量机(Local binary pattern support vector machine, LBP-SVM)分类器,分别对LFCC特征进行训练和测试,使模型更关注语音的高频区域。Todisco等^[6]将CQCC特征引入到语音伪造检测中,使用高斯混合模型捕获语音中的高低频伪造痕迹,取得当时的最佳性能(State of

the arts, SOTA)。

虽然传统的检测方法在应对伪造语音检测任务取得了不错的性能,但该类方法极度依赖手工特征工程,以及核函数的选择和设计,普适性不高。

1.2 基于深度学习的检测方法

该类方法通过对手工特征进一步学习,以挖掘更深层次的伪造痕迹。鉴于深度学习在图像分类任务中的出色表现,部分研究人员也开始将其应用于伪造语音检测中^[7]。Lavrentyeva等^[8]提出了一种基于轻量化卷积神经网络(Lightweight convolutional neural network, LCNN)的伪造语音检测方法。该方法首先提取语音的频谱特征,然后通过LCNN对其中的伪造痕迹进行学习。尽管该方法在解决模型参数规模问题上表现出色,但随着模型层数的急剧增加,却极易陷入梯度消失的困境。为此,Lai等^[9]将残差神经网络模型引入伪造语音检测任务中,提出一种基于残差连接的伪造语音检测方法,实验结果表明,所提方法能够解决梯度消失问题。但是,该方法无法直接应用于序列性的语音检测任务。为此,Gs等^[10]提出一种基于双向长短期记忆网络(Bi-directional long short-term memory network, Bi-LSTM)的伪造语音检测方法,通过对时间序列和远程依赖进行建模,以捕获频谱特征中的序列性伪造痕迹。Gomez-Alanis等^[11]将轻量化的门控递归神经网络作为特征提取器,概率线性判别分析作为分类器,进一步提高伪造语音检测的性能。

当前,尽管基于深度特征学习的检测方法在提升模型检测性能方面展现出了一定的优势,但在频谱特征提取的过程中仍会不可避免地丢失部分关键信息。此外,分类器的选择过于依赖人工操作,难以实现特征提取与分类器之间的双向最优化,从而限制了检测性能的进一步提升。

1.3 基于端到端的检测方法

该类方法是直接构建一个能够鉴别语音真假的检测模型,无需执行一系列的复杂设计。为了提取更深层次的语义特征,Jung等^[12]设计了一个RawNet模型,提高了真假语音检测任务的准确率。Tak等^[13]提出一种基于图注意力网络的检测方法,能够更好地提取不同频率子带和时间段之间的关联性。此外,通过将图神经网络和其他深度学习模型融合,一定程度上弥补了单一网络模型检测精度不高的困扰。但是,该方法过于依赖语音中的相位信息。为此,Wang等^[14]利用RawNet2网络进行特征提取,并设计了一个名为SimAM(Simple attention module)的注意力模块来对特征图进行深入分析。同时,还引入了均方误差作为模型的元训练策略,有效提升了检测模型的泛化能力。在Wang等^[15]的另一项工作中,通过将正交卷积与时间卷积网络相结合,并输入RawNet网络,构建了TO-RawNet网络,能够进一步提高伪造语音的检测精度。

1.4 基于特征融合的检测方法

基于特征融合的伪造检测方法通过融合各种相同或不同的特征,以增强特征的表征能力,进一步提高任务的检测性能。Xue等^[16]提出了一种基于生理-物理特征融合的伪造语音检测方法,通过卷积神经网络提取音频中与人脸相关的生理特征(如性别、嘴型等)以及物理特征来识别语音欺骗攻击。Wang等^[17]分别使用隐藏单元BERT(Hidden-unit BERT, HuBERT)和Conformer模型提取语音的音素特征,并将其与Wav2Vec2模型提取的特征进行融合,通过深入分析原始语音信号的语调、音高等韵律信息来检测伪造语音。Liu等^[18]提出了一种创新方法,将原始语音转换为立体声,并从左右声道两个方面获取深层特征表示,最终将这两种深层特征进行融合,以获取更为全面且细致的伪造痕迹信息。

尽管基于特征融合的方法在提升伪造语音的检测性能上取得了一定成效,但仍然面临特征冗余和细粒度伪造痕迹难以挖掘等挑战。

2 基于时频特征融合的伪造语音检测模型

针对当前伪造语音检测任务所面临的诸多挑战,如难以精准捕获语音连贯性、难以检测高频区域出现的异常波峰,以及语音波形特征对相位和幅度信息过度依赖等不足,提出了一种新型伪造语音检测方法。本节首先将阐述多分支特征融合网络的设计过程;然后对现有的坐标注意力机制的局限性进行说明,并提出一种时频坐标注意力机制,能够增强差异特征的提取和表征能力;最后,为了评估不同分支特征的重要性并剔除冗余特征还设计了一个自适应联合损失优化函数,进一步提高伪造语音检测任务的准确性。

2.1 多分支特征融合网络

为解决当前特征学习方法在频谱信息丢失和挖掘不同伪造语音技术的共性特征方面难的问题,提出一种基于时频特征融合的伪造语音检测算法,具体实现流程如图1所示。与直接在物理层面进行特征融合不同,所提方法采用一种基于分类结果层面的均值融合策略,旨在避免因直接融合不同维度信息而导致的特征冗余,以及规避可能造成的特征维度灾难。通过将时域和频域特征在分类层面进行融合,能够使模型专注于各自特征空间的模式学习,充分利用不同特征维度的独特性。在分类结果融合阶段,整合了时域和频域特征所捕获的关键分类信息,有效提升了检测性能和泛化能力。

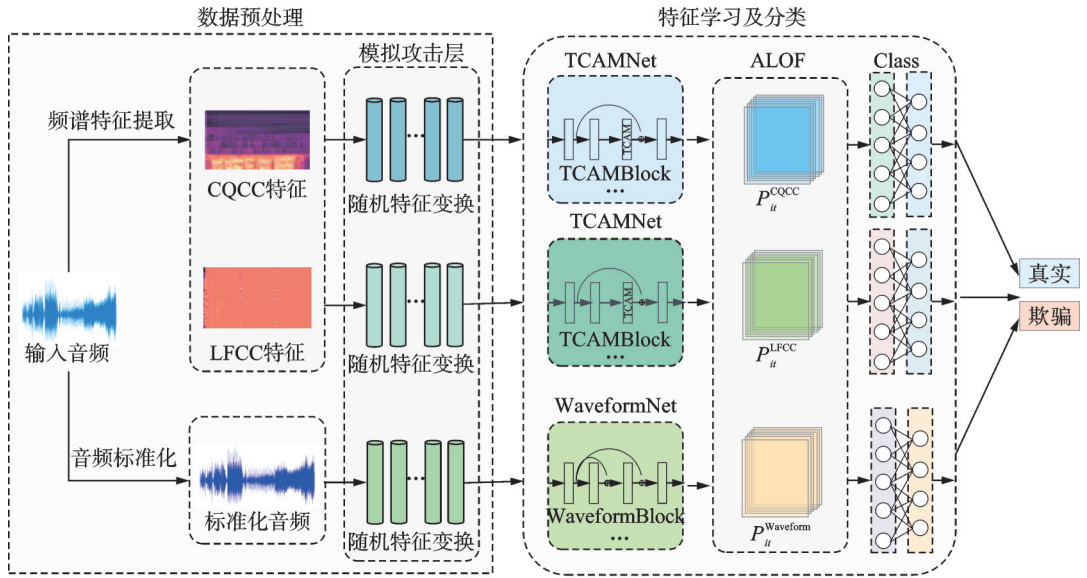


图1 多分支特征融合网络总体架构

Fig.1 Overall architecture of MFFNet

本文方法旨在从语音波形、LFCC和CQCC三个维度充分挖掘伪造语音的细粒度特征。首先,借助频谱变换技术,从原始的语音信号中提取出CQCC和LFCC两类频谱特征。随后,将这些特征输入到模拟攻击层(Simulation attack layer, SAL)之中,通过模拟随机噪声和信号增益等操作,在时频域内实施双重扰动策略,使得模型能够在信号干扰和传输等复杂环境下,能够精准地学习到伪造语音所独有的特征,频谱变换 $S(\bullet)$ 和时域波形变换 $W(\bullet)$ 实现如下

$$S_{\text{shift}}(t, f) = X(|(t - \Delta t) \bmod T|, |(f - \Delta f) \bmod F|) \quad (1)$$

首先,为了模拟语音信号在时间和频率上的特征偏移干扰,如式(1)所示,对二维频谱信号 $X(t, f)$

在时频两个维度上分别执行周期性偏移(偏移量分别为 Δt 和 Δf),并采用模运算 mod (T 和 F 分别表示时频维度的最大范围)来限制偏移范围。为了进一步模拟语音信号在时间或频率范围内的攻击场景,分别引入时域 $M_t(t_{\text{mask}})$ 和频域 $M_f(f_{\text{mask}})$ 掩码函数,即

$$S_{\text{mask}}(t, f) = X(t, f) \cdot M_t(t_{\text{mask}}) \cdot M_f(f_{\text{mask}}) \quad (2)$$

对频谱信号 $X(t, f)$ 部分区域(t_{mask} 和 f_{mask} 范围内)执行时频域的信号掩蔽操作。其中, $M(\cdot)$ 值为0表示该区域内信号被掩蔽,值为1表示该区域信号未被掩蔽,计算公式如下

$$M_t(t_{\text{mask}}) = \begin{cases} 0 & t_{\text{mask}} \in [t_{\text{start}}, t_{\text{start}} + L_t] \\ 1 & \text{其他} \end{cases} \quad (3)$$

$$M_f(f_{\text{mask}}) = \begin{cases} 0 & f_{\text{mask}} \in [f_{\text{start}}, f_{\text{start}} + L_f] \\ 1 & \text{其他} \end{cases} \quad (4)$$

式中: t_{start} 和 f_{start} 分别为指定掩蔽区域的起始位置, L_t 和 L_f 用来确定掩蔽区域的范围。通过随机掩蔽时频域中部分信息,增强模型对语音局部信息丢失的适应性,进而提升对低质量信号的感知能力。为了进一步增强模型在面临各种时频域攻击噪声时的鲁棒性,还在时频域引入随机噪声,即

$$S_{\text{noise}}(t, f) = X(t, f) + \sigma \cdot M \quad M \sim N(0, 1) \quad (5)$$

对频谱信号 $X(t, f)$ 添加服从标准正态分布 $N(0, 1)$ 的高斯噪声矩阵 M ,通过参数 σ 控制噪声强度。为了增强模型对时域信号变化时的鲁棒性,对频谱信号 $X(t, f)$ 执行线性变换来模拟信号增益,求解公式为

$$S_{\text{adjust}}(t, f) = X(t, f) \cdot \lambda \quad \lambda \sim U(1 - \epsilon, 1 + \epsilon) \quad (6)$$

式中: $\lambda \sim U(1 - \epsilon, 1 + \epsilon)$ 为服从连续型均匀分布的随机增益因子。为了进一步增强模型对语音信号时间轴偏移干扰的适应能力,还对语音波形信号 $X(t)$ 执行时间轴上的平移操作,即

$$W_{\text{shift}}(t) = X(|(t - \Delta t) \bmod T|) \quad (7)$$

此外,为增强模型应对不同波形信号强度下的检测能力,对语音波形信号 $X(t)$ 执行信号增益操作,随机增益因子 $\lambda \sim U(1 - \epsilon, 1 + \epsilon)$ 服从连续型均匀分布,计算公式如下

$$W_{\text{adjust}}(t) = X(t) \cdot \lambda \quad \lambda \sim U(1 - \epsilon, 1 + \epsilon) \quad (8)$$

接下来,模型需进一步提升对语音波形噪声的适应性。为此,向语音波形信号 $X(t)$ 添加随机噪声 I ,并通过参数 γ 调节噪声强度,即

$$W_{\text{noise}}(t) = X(t) + \gamma \cdot I \quad I \sim N(0, 1) \quad (9)$$

最终,为了模拟语音复杂干扰场景,通过双重变换实现组合攻击,具体如下

$$X_{\text{enhanced}} = \begin{cases} f_k(f_j(X)) & X \in \{\text{CQCC}, \text{LFCC}\} \\ g_k(g_j(X)) & X \in \text{Waveform} \end{cases} \quad (10)$$

$$f_k, f_j \in \{S_{\text{shift}}, S_{\text{mask}}, S_{\text{noise}}, S_{\text{adjust}}\}, g_k, g_j \in \{W_{\text{shift}}, W_{\text{adjust}}, W_{\text{noise}}\}$$

式中: $f(\cdot)$ 和 $g(\cdot)$ 分别表示单次特征变换($S(\cdot)$ 或 $W(\cdot)$ 变换),信号 X 经过组合特征变换后得到的 X_{enhanced} 可更真实地模拟复杂语音干扰场景。如果信号 X 为频谱特征,则对频谱信号应用两次频谱变换 f_k, f_j ;若信号 X 为语音波形特征,则应用两次波形变换 g_k, g_j 。

为了进一步挖掘不同伪造技术中的不变共性特征,将增强后的特征分别送入TCAMNet和WaveformNet两个专用网络进行处理。在CQCC特征提取分支中,利用如图2所示的TCAMNet骨干网络深入分析长期窗口变换的CQCC特征,以学习频谱图中各位置的能量分布模式。此外,该分支还能有效捕捉低频部分的高分辨率信息,从而更精确地刻画真伪语音在能量分布上的差异。WaveformNet网络

内部结构如图3所示,则从语音波形中提取深度特征。利用其内置的池化模块,能够对这些特征进行特征降维,以更精准地挖掘波形特征中的伪造痕迹。

在LFCC特征提取分支中,利用线性滤波器精准提取语音信号中的谐波结构、共振峰以及高频区域能量等关键信息。该分支与CQCC特征提取分支采用相同的网络结构和学习策略,以确保特征提取的一致性和有效性。选择频谱特征作为输入的优势在于有助于降低模型对相位信息的过度依赖,提高模型的检测精度和稳定性。最后,设计了ALOF模块,能够自适应地调整模型的联合损失,减少模型对语音中噪声以及特征中冗余信息的学习,进一步提升伪造语音检测的性能。

2.2 时频坐标注意力机制

坐标注意力机制在计算机视觉领域得到了广泛应用,核心思想是通过在通道注意力的基础上引入位置信息,以更好地捕捉特征的局部判别性信息^[19]。为了有效提取语音中的伪造痕迹和语义不连贯性特征,首次将坐标注意力机制引入到伪造语音检测任务中。文本转语音(Text-to-speech, TTS)作为当前主流的语音伪造技术,充分考虑了人耳对部分频带敏感度较低的特点,在生成语音时降低了部分频率的质量。这种处理方式致使合成语音在不同频带之间的能量分布不均匀,呈现出不正常的能量连贯性,如图4所示。伪造语音在部分频带之间普遍存在能量异常,具体体现在频谱图的不同频带部分出现能量损失或增益不平衡的现象。通过观察图4(b)中的伪造语音频谱图,可以明显发现部分频带的能量分布相较于真实语音存在明显的缺失,呈现出能量异常的断层。这些频带能量异常现象是鉴别伪造语音的重要特征,为检测模型提供了有力的判别依据。

坐标注意力机制此前主要应用于图像分类任务,但在处理音频的时频域方面,其能力相对受限,难以挖掘真假语音间细微的伪造特征,从而影响了语音伪造检测任务的性能。鉴于此,本文对其进行了优化改进,提出了一种专门用于伪造语音检测的时频坐标注意力机制。在具体实施中,首先通过一组卷积操作对输入的频谱图进行时频信息提取。然而,经典的坐标注意力机制采用的是平均池化,很难捕捉到某些能量异常区域,反而可能导致伪造痕迹信息被模糊处理,进而对伪造检测

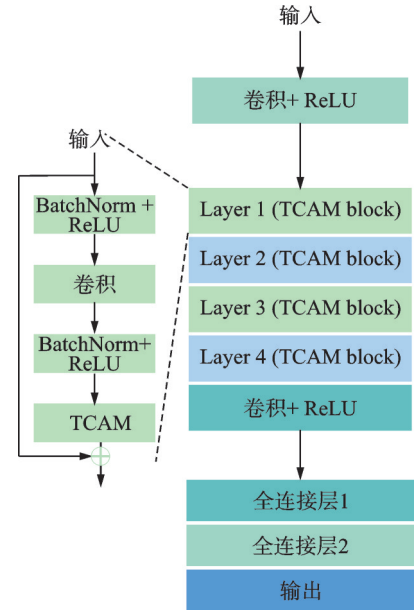


图2 TCAMNet学习架构

Fig.2 TCAMNet learning architecture

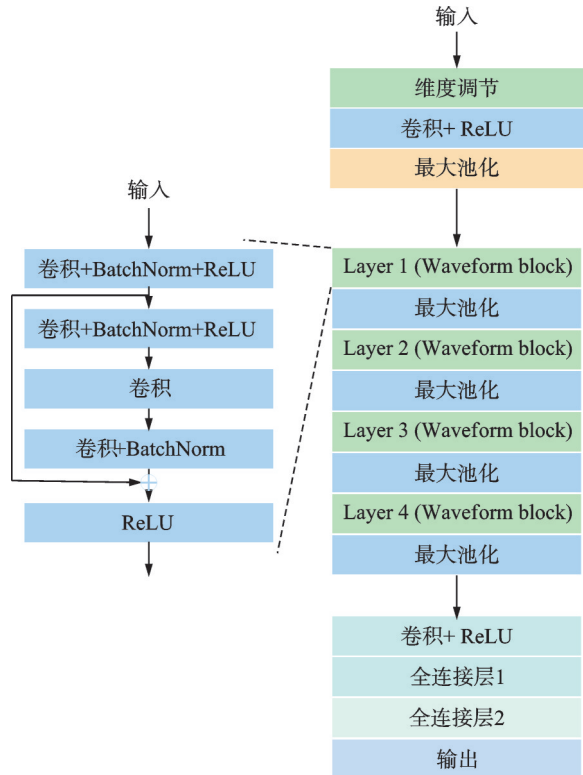


图3 WaveformNet学习架构

Fig.3 WaveformNet learning architecture

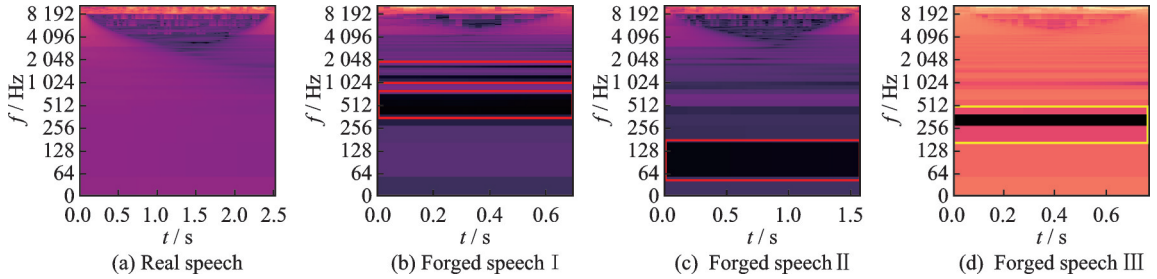


图4 不同语音伪造痕迹在CQCC频谱图中的位置

Fig.4 Positions of different speech forgery traces in CQCC spectrogram

的性能产生不利影响。为了挖掘全局特征中的局部变化,改用了针对于时频域的最大自适应池化操作,以捕获全局特征中的局部变化,如频谱中峰值变化的尖锐区域和能量不连续的跳变区域。通过该操作能够在保证全局信息不受损失的条件下,准确提取伪造语音中微小能量变化和频谱失真信息,避免微小的时频特征被忽略。随后,池化后的时频域特征会经过特征连接和分离等操作,进一步增强对不同频带及时间段特征能量信息的提取。具体而言,通过特征连接操作,将来自时频域尺度的特征信息融合在一起,构建一个多尺度的特征表示。该融合能够让模型同时捕捉到频谱中不同频率段的全局信息和局部细节,提升模型的时频特征表达能力。特征分离操作则通过将频域和时域的特征执行分离操作,有助于在时频域内分别强化伪造痕迹的表征,特别是语音信号的高频异常和时序上的语义不连贯性。在注意力中间层,将最大通道数固定为64,以确保在特征提取过程中保持足够的特征丰富度。此外,考虑到注意力网络中深层的通道数量会逐渐衰减,通过设计适当的通道衰减策略来避免模型的特征学习能力过早下降。通过在每一层中加入一定衰减比例的约束通道数,使得时频信息的压缩和提炼保持在合理的范围内,进而最大程度地提高检测精度。为了进一步提高语音检测模型的稳定性和泛化能力,在该注意力机制激活函数的选择上还进行了优化。采用ReLU4激活函数,它是一种调整后的ReLU函数,能够限制激活值的上界为4。相较于标准的ReLU,ReLU4在伪造语音检测模型收敛过程中表现出更强的稳定性,特别是在高频特征和能量异常区域的检测中,可大幅减少特征值过大导致的检测误差。

机制内部结构如图5所示,输入 X 经过TCAM后,首先通过自适应最大池化层分别对频谱图的时间方向和频率方向进行池化操作,得到两个池化后的频域和时域张量 x_f 和 x_t 。然后将这两个张量拼接为一个新的张量 y ,并对其执行双层卷积、批量归一化和激活等操作,得到处理后的 y 。接着,将 y 进行特征分离并再次生成张量 x_f 和 x_t ,对其再次执行时频域卷积操作。最后,利用激活函数分别计算两个注意力模块的权重,并将这两个注意力权重与输入特征 X 相乘,作为该模块的最终输出。

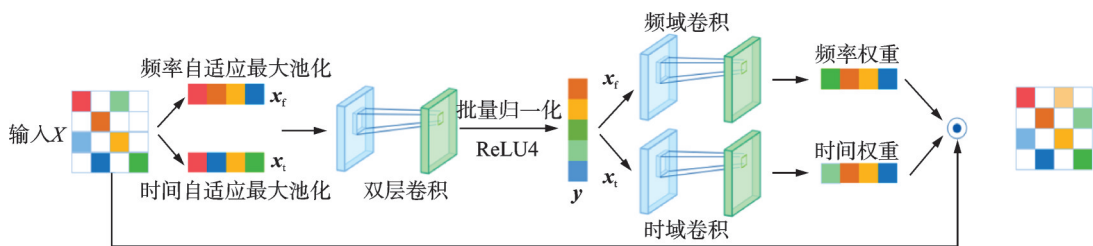


图5 时频坐标注意力机制

Fig.5 Time-frequency coordinate attention mechanism

2.3 自适应联合损失优化函数

为了评估不同分支特征的重要性,以便充分学习细粒度伪造痕迹,设计了一种自适应联合损失优化函数,其包括 MultimarginLoss 函数^[20]和交叉熵损失函数^[21]。在训练初期,采用交叉熵损失函数,能够有效处理类别的分布,引导模型学习真假语音的差异性纹理,实现模型训练的快速收敛。一旦模型趋向收敛状态,则切换到 MultimarginLoss 函数,并通过 margin 对模型的分类边界进行微调。MultimarginLoss 函数能够通过设置类别间的距离,帮助模型细化决策边界,提高对复杂伪造类型样本的分类。目前主流的伪造语音数据的数据样本往往分布不够均衡,因此采用 MultimarginLoss 进行微调也能够使模型对不同类型伪造语音特征进行更为全面的学习,缓解正负样本失衡带来的分类偏差问题。其中,交叉熵损失函数的计算公式为

$$\text{Loss} = -\frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 y_{it} \ln P_{it} \quad (11)$$

式中: S 表示样本数量, y_{it} 表示真实标签的类别, P_{it} 表示模型的预测类别, t 表示样本标签数。通过运用交叉熵损失函数来比较预测类别与真实标签之间的差异,可以使模型初步达到收敛的状态。其中 MultimarginLoss 函数的定义为

$$\text{Loss} = \frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 \max(0, \text{margin} - P_{it} + y_{it})^p \quad (12)$$

式中: P_{it} 代表模型的预测值, y_{it} 对应真实的标签类别, p 为计算损失时的指数值,用来约束对模型错误分类的惩罚力度,取值只能为1或2, i 表示数据编号。 p 值越大,模型对错误分类所施加的惩罚力度也就越强。margin 是指在分类任务中,正确分类得分与其他类别得分之间的最小间隔。当正确类别的得分超过一个阈值,并且超过所有类别的得分之和时,损失函数的取值为0;反之,则为正值。当 margin 较大时,会增加标签值和其他类别间的间隔,能在一定程度上提高模型的预测能力。而当 margin 值过大时,也会使得模型产生过拟合,而若 margin 值较小,会缩小真实标签和预测错误值间的间隔,使得模型的决策边界更趋于平滑,泛化性能更强。但若 margin 过小,则会导致样本学习不充分。为了评估不同分支特征的重要性,提出了一种自适应联合损失优化函数 ALOF,通过调节不同分支网络在任务中的权重,提升模型对不同特征伪造痕迹的学习能力,其公式如下

$$\text{Loss} = \begin{cases} \text{Loss}_{\text{front}} & \text{Loss}_{\text{dev}} \geq \text{Threshold} \\ \text{Loss}_{\text{rear}} & \text{其他} \end{cases} \quad (13)$$

$$\text{Loss}_{\text{front}} = \max \left(-\frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 y_{it} \ln P_{it}^{\text{CQCC}}, -\frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 y_{it} \ln P_{it}^{\text{LFCC}}, -\frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 y_{it} \ln P_{it}^{\text{waveform}} \right) \quad (14)$$

$$\begin{aligned} \text{Loss}_{\text{rear}} = \max & \left(\frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 \max(0, \text{margin} - P_{it}^{\text{CQCC}} + y_{it})^p, \frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 \max(0, \text{margin} - P_{it}^{\text{LFCC}} + y_{it})^p, \right. \\ & \left. \frac{1}{S} \sum_{i=0}^{S-1} \sum_{t=0}^1 \max(0, \text{margin} - P_{it}^{\text{Waveform}} + y_{it})^p \right) \end{aligned} \quad (15)$$

有关 ALOF 的推理过程如式(13~15)所示。其中在式(13)中, Loss 表示真伪语音分损失, Threshold 为控制损失函数的阈值,将其设为0.05。当 Loss_{dev} 低于阈值时进入后期训练阶段, t 代表分类的类别数。在式(14~15)中, y_{it} 表示真实标签的类别, S 表示样本数量, P_{it}^{CQCC} 、 P_{it}^{LFCC} 和 P_{it}^{Waveform} 分别表示不同分支下的预测类别, p 表示对不同类型分支进行调节的惩罚项。通过自适应地调节联合损失,能够解决模型的过拟合问题,提高模型的检测精度。

3 实验和分析

3.1 数据集

为了构建一个共享的性能测试平台,推动学术界和产业界的研究人员在伪造语音检测领域的交流与合作,英国爱丁堡大学联合多个研究机构成功举办 ASVspoof 反欺骗挑战赛,至今已连续举办多届,吸引了包括中山大学在内的众多研究人员的参与。在每届挑战赛启动前,主办方都会发布一批最新的数据集,其涵盖语音合成、语音转换等多种攻击类型的数据。同时,为了便于不同算法之间的性能对比和优秀算法的筛选,主办方还会给出相应的评价指标。本文采用的数据集为主流的 ASVspoof 2019 LA 数据集,其详细信息如表 1 所示。ASVspoof 2019 LA 数据集在无通道或噪声干扰的环境下,从语音克隆工具箱(Voice cloning toolkit, VCTK)语料库中收集了 46 名男性和 61 名女性的真实语音样本,并通过多种语音合成和语音转换技术生成伪造语音。该数据集分为 Train、Dev 和 Eval 三部分,其中 Train 和 Dev 主要包含 6 种已知的伪造语音攻击,而 Eval 则涵盖了 2 种已知攻击和 11 种未知攻击的伪造语音,以全面检验算法的性能^[22]。ASVspoof 2021 数据集包括两种类型:LA 数据集和深度伪造(Deep fake, DF)数据集。LA 数据集中的伪造语音由 13 种不同的语音合成和语音转换攻击算法生成,在生成过程中引入了编码和传输干扰,进一步增加了伪造的复杂性。相比之下,DF 数据集的伪造语音由 110 种不同的欺骗性伪造算法生成,并通过有损编码器来处理真实语音和伪造语音,以模拟语音失真场景。本文方法在 ASVspoof 的训练集中进行训练,在 2019 LA、2021 LA 以及 2021 DF 的测试集中进行测试^[23]。

表 1 ASVspoof 数据集
Table 1 ASVspoof dataset

数据集		真实语音	欺骗语音	
		数量	数量	欺骗攻击
ASVspoof 2019 LA	训练集	2 580	22 800	A01-A06 (6 种)
	开发集	2 548	22 296	A01-A06 (6 种)
	测试集	7 355	63 882	A07-A19 (13 种)
ASVspoof 2021 LA	测试集	14 816	133 360	A07-A19 (13 种)
ASVspoof 2021 DF	测试集	14 869	519 059	110 种

3.2 评价指标

为了确保对不同算法的性能评估的公正性和一致性,采用 7 个公开的评价标准用于性能的评估,包括等错误率、最小归一化串联检测代价函数、召回率(Recall)、 F_1 得分(F_1 -score)、检测误差权衡曲线(Detection error tradeoff, DET)和接受者操作特征曲线(Receiver operating characteristic, ROC)。此外,为了便于观察所提算法的分类效果,本文还采用 t 分布-随机邻近嵌入(t-distributed stochastic neighbor embedding, t-SNE)指标,对模型的输出结果进行可视化展示。EER 是指在错误接受率与错误拒绝率相等时的特定数值,作为一个综合性指标,能够全面反映模型的检测性能。其计算公式为

$$EER = P_{fa,cm} = P_{fr,cm} \tag{16}$$

式中: $P_{fa,cm}$ 表示错误接受率, $P_{fr,cm}$ 为错误拒绝率。最小归一化串联检测代价函数 min t-DCF 作为一种伪造检测的评价指标,能够更好地评估模型的性能。值越小,则表明方法性能愈佳。其计算公式为

$$\min t\text{-DCF} = C_{fr} \cdot P_{fr,cm} + C_{fa} \cdot P_{fa,cm} \tag{17}$$

式中: C_{fr} 为错误拒绝率的风险系数,本文取值为 2.405 95, C_{fa} 表示错误接受率的风险系数,取值为 1。

min t-DCF 的值越小,说明算法的性能越好。精确率是指所有模型预测为正样本的样本数量中,模型预测为正样本且分类正确的样本数量的概率。值越大,表示算法的性能越好。其计算公式为

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (18)$$

式中:TP代表实际为正样本且被模型正确预测为正样本的数目,FP则指实际为负样本却被模型错误地预测为正样本的数目,而FN则是实际为正样本却被模型错误地预测为负样本的数目。作为评估指标之一,召回率表示在所有真实正样本中,被模型正确预测为正样本的比例,其计算公式为

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (19)$$

虽然准确率和召回率是用于评估不同算法性能的重要指标,但其仍然不能全面反映算法的表现。例如,某算法可能在准确率方面表现出色,但在召回率上却表现不佳。为了解决这一问题,引入了 F_1 得分作为一个综合评价指标,以平衡不同指标的影响,从而更全面地评估算法性能。 F_1 得分越高,代表着算法的性能越好。其计算公式为

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (20)$$

此外,t-SNE作为一种经典的降维可视化技术,能够充分考虑数据点之间的相对距离关系,从而将高维空间的模型分类结果映射至低维空间。这一过程不仅保留了数据的局部结构特征,还使得研究者能够直观地观察和理解模型的分类效果,为后续的模型优化提供有力的支持^[24]。DET和ROC曲线作为评估二元分类模型性能的指标,能够直观反映模型分类性能。其中,DET曲线的横轴表示假正例率(False positive rate, FPR),即负样本被错误分类为正样本的概率^[25];纵轴表示假负例率(False negative rate, FNR),即正样本被错误分类为负样本的概率。DET曲线越接近图像的左下角,其性能越好。曲线下的面积与模型性能成反比,面积越大,则方法性能越低。ROC曲线的横轴表示假正例率^[26],纵轴表示真正例率(True positive rate, TPR),即真实样本被正确分类为真实样本的概率。若曲线越靠近左上角,则表明方法性能越好。

3.3 实验环境与参数设置

本文所采用的网络模型是基于Pytorch 1.7.1深度学习框架,运行在Ubuntu 20.04.5 LTS操作系统上。硬件配置包括12th Gen Intel(R) Core(TM) i5-12400处理器以及2块NVIDIA Tesla P40 24 GB显卡。在骨干网络的训练过程中,采用Adam优化器进行参数的优化,学习率设置为0.000 1,并将反向传播参数lambda设为0.05。对于骨干网络中的TCAMNet,频谱图的特征长度设置为750,编码维度设置为256,WaveformNet的输入音频设置为6 s。

3.4 实验结果与分析

本文首先在ASVspoof数据集上进行了攻击成功率测试和泛化性测试,详细实验结果如表2所示。

通过与若干现有的方法对比可知,所提算法在EER和min t-DCF两个指标上均取得最佳性能。其中,与之前在原始波形上表现的方法Waveform+Raw PC-DARTS^[36]相比,所提方法在EER和min t-DCF两个指标上分别降低了1.42%和0.040。在检测精度方面,与当前最佳的频谱图特征方法SpoTNet^[41]对比后,所提方法在EER和min t-DCF两个指标上分别降低0.60%和0.034,均取得了SOTA性能。通过与现有工作相比,所提方法在特征丰富度表示上更为全面,能够削弱伪造语音检测任务对幅度和相位信息的严重依赖,伪造检测的性能也得到了显著的提升。此外,与Yue等^[40]提出的将频谱特征和频谱矩阵相融合方法相比,所提方法性能更好,其中,EER降低了0.61%、min t-DCF降低了0.019。

表 2 不同方法在 ASVspoof 数据集上的实验结果
Table 2 Experimental results of different methods on ASVspoof dataset

数据集	模型	EER/%	min t-DCF
ASVspoof 2019 LA	ProsoSpeaker ^[27]	5.39	
	CFE-ResNet ^[28]	5.23	
	ARawNet2 ^[29]	4.61	
	Stat-SE-Res2Net50 ^[30]	2.86	0.068
	Student Net-KA ^[31]	2.39	0.067
	OC-softmax ^[32]	2.19	0.059
	Spec+Multi-branch ^[33]	2.17	0.118
	Spec+LFCC+DenseNet ^[34]	1.98	0.047
	CQT+MCG-Res2Net50 ^[35]	1.78	0.052
	Waveform+Raw PC-DARTS ^[36]	1.77	0.051
	MFCC+GTCC+SMOTE-LSTM ^[37]	1.60	
	L-VQT+DenseNet ^[38]	1.54	0.045
	LFCC+OCT ^[39]	1.06	0.034
	Fusion-PA-SE-ResNet ^[40]	0.96	0.030
	SpoTNet ^[41]	0.95	0.045
	Proposed method	0.35	0.011
ASVspoof 2021 LA	CQCC-GMM ^[6]	15.62	0.479 4
	LFCC-GMM ^[42]	19.30	0.575 8
	LFCC-LCNN ^[43]	9.26	0.344 5
	RawNet2 ^[44]	9.50	0.425 7
	Proposed method	8.94	0.406 3
ASVspoof 2021 DF	CQCC-GMM ^[6]	25.56	
	LFCC-GMM ^[42]	25.25	
	LFCC-LCNN ^[43]	23.48	
	RawNet2 ^[44]	22.38	
	Proposed method	19.81	

此外,所提方法在 ASVspoof 2021 LA 和 DF 数据集上的表现同样出色,展现了较强的泛化能力。在 ASVspoof 2021 LA 数据集上 EER 上实现了 8.94%,显著优于现有的其他方法。与传统的 GMM 方法相比,所提方法的性能显著提升。而与 LFCC-LCNN 方法^[43]相比,所提方法在 EER 上降低 0.32%,在 min t-DCF 上略高,说明所提方法在特征信息挖掘方面更关注于伪造检测,在一定程度上忽视了对说话人信息的区分。与此同时,在 ASVspoof 2021 DF 数据集上,所提方法的 EER 为 19.81%,相比 CQCC-GMM^[6]、LFCC-GMM^[42]和 LFCC-LCNN 方法^[43]均表现出显著的性能提升。与 RawNet2 方法^[44]相比,本文方法也取得了 2.57% 的性能提升。结果表明,所提方法不仅在 ASVspoof 2019 LA 上表现优异,而且在其他数据集上同样表现出较好的跨数据集泛化性。

为了验证本文所提方法的泛化能力,还对 ASVspoof 2019 LA 中未知类型的欺骗攻击进行了性能测试,实验结果如表 3 所示。通过观察结果可知,在面对目前最具挑战性的 A17 攻击时,所提方法在 EER 指标的评估中仍表现出色。A17 攻击采用语音波形滤波技术来处理伪造语音,通过对原始语音进

表3 面对不同攻击类型时的EER和min t-DCF

Table 3 EER and min t-DCF for different attack types

模型 攻击类型	OCSoftmax ^[32]		Inc-TDSSDNET ^[45]		Res-TDSSDNET ^[45]		Proposed method	
	EER/%	min t-DCF	EER/%	min t-DCF	EER/%	min t-DCF	EER/%	min t-DCF
A07	0.122 2	0.003 4	0.692 7	0.019 5	1.443 0	0.039 5	0.098 5	0.002 0
A08	0.180 0	0.004 2	16.562 5	0.300 4	0.750 4	0.020 6	0.203 7	0.005 4
A09	0.122 2	0.009 4	0.431 2	0.034 1	0.040 7	0.002 5	0.000 0	0.000 0
A10	1.140 8	0.024 4	0.831 9	0.024 2	1.786 0	0.049 3	0.220 7	0.006 4
A11	0.122 2	0.002 5	0.668 9	0.019 1	0.064 5	0.001 5	0.017 0	0.000 5
A12	0.465 2	0.010 7	0.709 6	0.019 9	0.180 0	0.004 3	0.098 5	0.001 8
A13	0.221 0	0.006 1	0.163 0	0.004 1	0.064 5	0.001 5	0.017 0	0.000 5
A14	0.692 7	0.018 7	0.326 0	0.008 8	0.139 2	0.003 2	0.023 8	0.000 5
A15	1.102 3	0.037 6	0.797 9	0.021 9	2.094 9	0.054 1	0.139 2	0.003 9
A16	0.326 0	0.007 4	1.164 6	0.027 8	1.246 1	0.027 6	0.098 5	0.002 4
A17	9.218 4	0.615 8	12.617 1	0.786 0	6.006 4	0.325 4	1.198 6	0.120 1
A18	0.896 4	0.037 8	1.728 2	0.111 9	1.117 1	0.064 2	0.098 5	0.009 2
A19	0.896 4	0.031 7	2.078 0	0.079 3	1.443 0	0.051 1	0.122 2	0.004 5
Total	2.178 8	0.055 8	4.043 1	0.097 6	1.642 1	0.048 1	0.353 6	0.010 5

行时域处理,实现伪造语音中音色的变化和噪声的降低,使伪造痕迹难以被提取。本文通过对语音波形的时域信息进行精细处理,成功提取了时域共性特征,显著提升了对经波形滤波操作生成的伪造语音的检测能力。与Hua等^[45]提出的Res-TDSSDNET方法对比,本文方法在A17类型攻击检测的EER指标上实现了4.807 8%的降低,展现了更高的检测精度和优越性。然而,当面对A08这类攻击时,本文所提算法的EER稍逊于现有方法。对数据集进行内容分析后发现,A08攻击是一种基于神经网络的语音伪造技术,因此在语音伪造过程中,会在语音波形上留下一些不自然的痕迹。而本文所提算法与原始语音在时域上的波形变化密切相关,这种依赖性会严重影响原始任务的检测性能^[22]。此外,本文还进行了min t-DCF的测试,实验结果如表3所示。观察结果发现,本文所提方法在应对大部分未知攻击时仍表现出较强的泛化性能。相较于Res-TDSSDNET^[45],本文所提方法在A17类型攻击检测的min t-DCF降低了0.205 3,然而在A08类型攻击下,min t-DCF的效果稍显不佳。

鉴于ASVspoof 2019 LA数据集中真实语音与伪造语音的样本不平衡问题,本文选用召回率和 F_1 得分作为评价指标,以评估算法在应对多种未知类型攻击时的泛化能力。在测试中,将欺骗语音标记为正样本,真实语音则被视为负样本,通过观察图6和图7可知,本文方法在召回率和 F_1 得分上的表现均优于基准方法,进一步表明本文算法具有较好的泛化性。

为了便于观察模型的分类效果,在图8(a)和图8(c)中,分别用蓝色和红色标识真实语音(Bonafide)和欺骗语音(Spoof),可发现本文所提方法对真实语音和欺骗语音展现出更高效的聚类能力。进一步分析图8(b)和图8(d),蓝色代表真实语音,而其他不同颜色则分别对应不同类型的攻击。通过对比图8(b)和图8(d),可以清晰地看到,相较于OCSoftmax方法^[32],本文方法在聚类不同类型攻击时表现出了更高的准确性,从而验证了其在语音伪造检测领域的优越性。此外,在DET指标测评中,本文所提方法表现出更佳的性能,如图9所示。左下角的覆盖面积低于其他3种方法,但假负例率略高于Res-TDSSDNET^[45],表明存在部分欺骗语音被识别为真实语音的情况。在ROC指标测试中,如图10所示,本文所提的方法在指标AUC达到了99.99%,超越其余3种方法,进一步证明了本文方法的有效性。

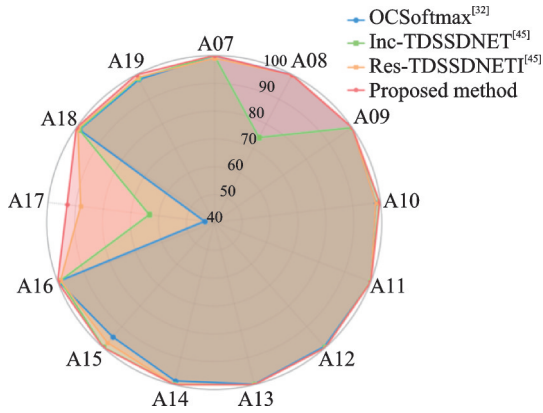


图6 对A07~A19的召回率

Fig.6 Recall rate of A07—A19

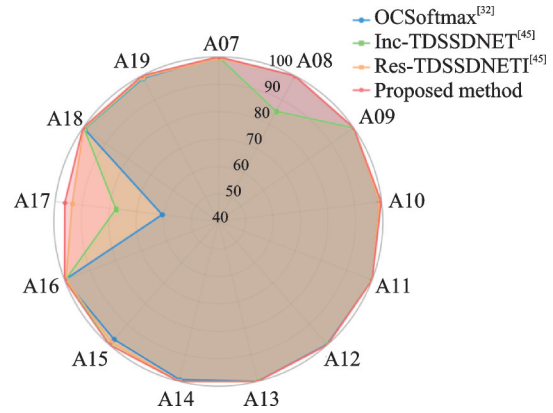
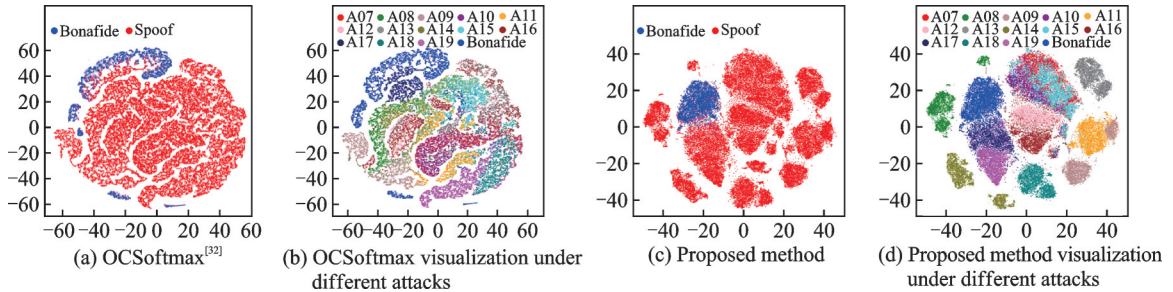
图7 对A07~A19的 F_1 得分Fig.7 F_1 -score of A07—A19

图8 在ASVspoof 2019 LA数据集上的T-SNE可视化结果

Fig.8 T-SNE visualization results on the ASVspoof 2019 LA dataset

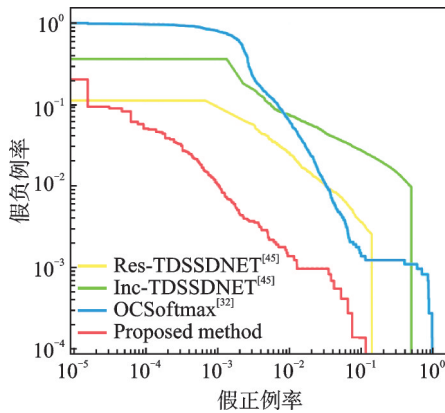


图9 DET曲线

Fig.9 DET curves

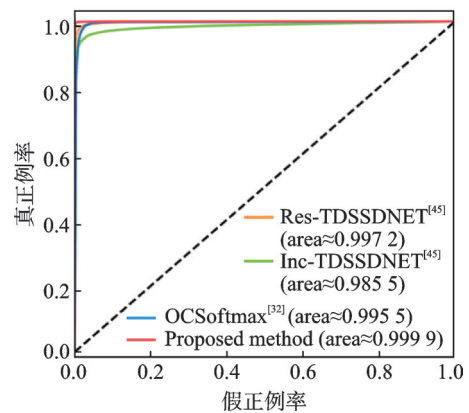


图10 ROC曲线

Fig.10 ROC curves

为了验证不同参数对算法性能的影响,还进行了消融实验,实验结果如表4所示。相较于单频谱图,双频谱图在EER方面表现出1.82%的较低数值,同时在min t-DCF指标上取得了更佳的效果。因此,采用双频谱图作为验证TCAM模块是有效的。当进行双频谱图的对比时,引入了TCAM模块后,EER的值下降了0.45%。对于引入TCAM模块的多特征融合方法,能够从时频角度挖掘语音中的伪

表 4 消融实验结果

Table 4 Ablation experimental results

模型	EER/%	min t-DCF
Signal-spectrum-branch (CQCC)	3.41	0.103
Signal-spectrum-branch (LFCC)	3.70	0.097
Signal-waveform-branch w/o Dual-spectrum-branch	2.14	0.062
Dual-spectrum-branch w/o Waveform-branch, w/o TCAM	1.82	0.062
Dual-spectrum-branch w/o Waveform-branch, w/ TCAM	1.37	0.040
Multi-feature fusion w/o TCAM, w/o SAL	0.70	0.022
Multi-feature fusion w/o TCAM, w/ SAL	0.58	0.018
Multi-feature fusion w/ TCAM, w/o SAL	0.50	0.015
Proposed method (Multi-feature fusion w/ TCAM, w/ SAL)	0.35	0.011

造痕迹, EER 达到了 0.50%。同时, 将 SAL 模块引入多特征融合方法后, 提高了数据多样性, 加强了模型对未知伪造语音能量分布的学习, 使得 EER 达到 0.58%。将 TCAM 和 SAL 引入多特征融合方法中后, EER 为 0.35%, min t-DCF 为 0.011, 再次证明本文所提算法的高效性。

4 结束语

语音识别作为一种身份认证技术, 目前已被广泛应用于各类身份认证场景。然而, 随着深度伪造技术的发展, 伪造语音被用来尝试欺骗各类语音识别系统, 引发了一系列严重的隐私泄露、非授权访问等安全纠纷, 因此对待认证语音的真伪进行鉴别至关重要。现有的伪造语音检测方法, 往往基于单一特征, 难以捕捉伪造语音中的语义连贯性, 且过度依赖时域的相位信息, 致使检测精度和泛化能力受限。因此, 本文提出了一种多分支特征融合网络, 分别从频率、共振峰以及时域等角度挖掘伪造语音中的能量异常痕迹及语义不连贯性细微差别, 以更好地表征真伪语音的音高变化、峰值差异和波形变化。然后, 进行时频伪造定位, 在时域上聚焦局部伪造特性, 且在频域上捕捉伪造语音中的偏移或异常频率, 并对时域和频域上的空间信息进行联合编码, 以便更好地关注伪造区域。实验结果表明, 本文所提方法在伪造语音检测和未知攻击识别方面均取得了显著的优势。

在接下来的工作中, 我们还将探索更高效的特征融合和伪造语音片段的定位方法, 以进一步解决伪造语音检测任务中面临的可解释难的问题。

参考文献:

[1] 李云峰, 闫祖龙, 高天, 等. 基于多任务学习的语音情感识别[J]. 数据采集与处理, 2024, 39(2): 424-432.
LI Yunfeng, YAN Zulong, GAO Tian, et al. Speech emotion recognition with multi-task learning[J]. Journal of Data Acquisition and Processing, 2024, 39(2): 424-432.

[2] 康瑶, 康坊, 杨飞然. 利用互子带滤波器和稀疏特性的多通道线性预测语音去混响方法[J]. 数据采集与处理, 2024, 39(5): 1135-1146.
KANG Yao, KANG Fang, YANG Feiran. Multi-channel linear prediction for speech dereverberation using cross-band filters and sparse priors[J]. Journal of Data Acquisition and Processing, 2024, 39 (5): 1135-1146.

[3] PATEL T B, PATIL H A. Combining evidences from Mel cepstral, cochlear filter cepstral and instantaneous frequency features for detection of natural vs. spoofed speech[C]//Proceedings of the 16th Annual Conference of the International Speech Communication Association. Dresden: International Speech Communication Association, 2015: 2062-2066.

- [4] KORSHUNOV P, MARCEL S, MUCKENHIRN H, et al. Overview of BTAS 2016 speaker anti-spoofing competition[C]//Proceedings of 2016 IEEE 8th International Conference on Biometrics Theory, Applications and Systems. Niagara Falls: IEEE, 2016: 1-6.
- [5] SAHIDULLAH M, KINNUNEN T, HANILÇI C. A comparison of features for synthetic speech detection[C]//Proceedings of the 16th Annual Conference of the International Speech Communication Association. Dresden: International Speech Communication Association, 2015: 2087-2091.
- [6] TODISCO M, DELGADO H, EVANS N. Constant Q cepstral coefficients: A spoofing countermeasure for automatic speaker verification[J]. Computer Speech & Language, 2017, 45: 516-535.
- [7] QIAN Yanmin, CHEN Nanxin, YU Kai. Deep features for automatic spoofing detection[J]. Speech Communication, 2016, 85: 43-52.
- [8] LAVRENTYEVA G, NOVOSELOV S, MALYKH E, et al. Audio replay attack detection with deep learning frameworks [C]//Proceedings of the 18th Annual Conference of the International Speech Communication Association. Stockholm: International Speech Communication Association, 2017: 82-86.
- [9] LAI C T, CHEN N, VILLALBA J, et al. ASSERT: Anti-spoofing with squeeze-excitation and residual networks[EB/OL]. (2019-04-01). <https://arxiv.org/pdf/1904.01120>.
- [10] GS A S, GANESHKUMAR V, THIRUMAVALAVAN V C, et al. Synthetic speech classification using bidirectional LSTM networks[C]//Proceedings of 2022 IEEE 3rd Global Conference for Advancement in Technology. Bangalore: IEEE, 2022: 1-6.
- [11] GOMEZ-ALANIS A, PEINADO A M, GONZALEZ J A et al. A light convolutional GRU-RNN deep feature extractor for ASV spoofing detection[C]//Proceedings of the 20th Annual Conference of the International Speech Communication Association. Graz: International Speech Communication Association, 2019: 1068-1072.
- [12] JUNG J W, KIM S B, SHIM H J, et al. Improved RawNet with feature map scaling for text-independent speaker verification using raw waveforms[C]//Proceedings of the 21st Annual Conference of the International Speech Communication Association. Shanghai: International Speech Communication Association, 2020: 1496-1500.
- [13] TAK H, JUNG J, PATINO J, et al. Graph attention networks for anti-spoofing[EB/OL]. (2021-04-08). <https://arxiv.org/pdf/2104.03654>.
- [14] WANG Z Y, HANSEN J H L. Audio anti-spoofing using a simple attention module and joint optimization based on additive angular margin loss and meta-learning[EB/OL]. (2022-11-17). <https://arxiv.org/pdf/2211.09898>.
- [15] WANG Chenglong, YI Jiangyan, TAO Jianhua, et al. TO-RawNet: Improving RawNet with TCN and orthogonal regularization for fake audio detection[EB/OL]. (2023-05-23). <https://arxiv.org/pdf/2305.13701>.
- [16] XUE Junxiao, ZHOU Hao. Physiological-physical feature fusion for automatic voice spoofing detection[J]. Frontiers of Computer Science, 2023, 17(2): 172318.
- [17] WANG Chenglong, YI Jiangyan, TAO Jianhua, et al. Detection of cross-dataset fake audio based on prosodic and pronunciation features[EB/OL]. (2023-05-23). <https://arxiv.org/pdf/2305.13700>.
- [18] LIU Rui, ZHANG Jinhua, GAO Guanglai, et al. Betray Oneself: A novel audio deepfake detection model via mono-to-stereo conversion[EB/OL]. (2023-05-25). <https://arxiv.org/pdf/2305.16353v1>.
- [19] HOU Qibin, ZHOU Daquan, FENG Jiashi. Coordinate attention for efficient mobile network design[C]//Proceedings of the IEEE/CVF Conference on Computer vision and Pattern Recognition. Nashville: IEEE, 2021: 13713-13722.
- [20] PyTorch Contributors. Multi margin loss[EB/OL]. (2023-12-01). <https://pytorch.org/docs/stable/generated/torch.nn.MultiMarginLoss.html>.
- [21] PyTorch Contributors. Cross entropy loss[EB/OL]. (2023-12-01). <https://pytorch.org/docs/stable/generated/torch.nn.CrossEntropyLoss.html>.
- [22] WANG X, YAMAGISHI J, TODISCO M, et al. ASVspoof 2019: A large-scale public database of synthesized, converted and replayed speech[J]. Computer Speech & Language, 2020, 64: 101114.

- [23] LIU X C, WANG X, SAHIDULLAH M, et al. ASVspoof 2021: Towards spoofed and deepfake speech detection in the wild [J]. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023, 31: 2507-2522.
- [24] CAI T T, MA R. Theoretical foundations of t-SNE for visualizing high-dimensional clustered data[J]. *Journal of Machine Learning Research*, 2022, 23(301):1-54.
- [25] ADLER A, SCHUCKERS M E. Calculation of a composite DET curve[C]//*Proceedings of Audio-and Video-Based Biometric Person Authentication: 5th International Conference*. Berlin: Springer, 2005: 860-868.
- [26] FAWCETT T. An introduction to ROC analysis[J]. *Pattern Recognition Letters*, 2006, 27(8): 861-874.
- [27] ATTORRESI L, SALVI D, BORRELLI C, et al. Combining automatic speaker verification and prosody analysis for synthetic speech detection[C]//*Proceedings of International Conference on Pattern Recognition*. Montréal Québec: Springer, 2022: 247-263.
- [28] ZHANG Jinghong, YI Xiaowei, ZHAO Xianfeng. A compressed synthetic speech detection method with compression feature embedding[C]//*Proceedings of the 24th Annual Conference of the International Speech Communication Association*. Dublin: International Speech Communication Association, 2023: 5376-5380.
- [29] LI Jing, LONG Yanhua, LI Yijie, et al. Advanced RawNet2 with attention-based channel masking for synthetic speech detection[C]//*Proceedings of the 24th Annual Conference of the International Speech Communication Association*. Dublin: International Speech Communication Association, 2023: 2788-2792.
- [30] LI Xu, LI Na, WENG Chao, et al. Replay and synthetic speech detection with Res2Net architecture[C]//*Proceedings of 2021 IEEE International Conference on Acoustics, Speech and Signal Processing*. Toronto: IEEE, 2021: 6354-6358.
- [31] REN Yeqing, PENG Haipeng, LI Lixiang, et al. Generalized voice spoofing detection via integral knowledge amalgamation [J]. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023, 31: 2461-2475.
- [32] ZHANG You, JIANG Fei, DUAN Zhiyao. One-class learning towards synthetic voice spoofing detection[J]. *IEEE Signal Processing Letters*, 2021, 28: 937-941.
- [33] WANG Ruoyu, DU Jun, WANG Chang. Multi-branch network with circle loss using voice conversion and channel robust data augmentation for synthetic speech detection[C]//*Proceedings of Chinese Conference on Biometric Recognition*. Beijing: Springer, 2022: 613-620.
- [34] WANG Zheng, CUI Sanshuai, KANG Xiangui, et al. Densely connected convolutional network for audio spoofing detection [C]//*Proceedings of 2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*. Auckland: IEEE, 2020: 1352-1360.
- [35] LI Xu, WU Xixin, LU Hui, et al. Channel-wise gated Res2Net: Towards robust detection of synthetic speech attacks[C]// *Proceedings of the 22th Annual Conference of International Speech Communication Association*. Brno: International Speech Communication Association, 2021: 4314-4318.
- [36] GE W, PATINO J, TODISCO M, et al. Raw differentiable architecture search for speech deepfake and spoofing detection [EB/OL]. (2021-10-06). <https://arxiv.org/pdf/2107.12212>.
- [37] CHAKRAVARTY N, DUA M. Data augmentation and hybrid feature amalgamation to detect audio deep fake attacks[J]. *Physica Scripta*, 2023, 98(9): 096001.
- [38] LI Jialong, WANG Hongxia, HE Peisong, et al. Long-term variable Q transform: A novel time-frequency transform algorithm for synthetic speech detection[J]. *Digital Signal Processing*, 2022, 120: 103256.
- [39] LI Changtao, YANG Feiran. YANG Jun. The role of long-term dependency in synthetic speech detection[J]. *IEEE Signal Processing Letters*, 2022, 29: 1142-1146.
- [40] YUE Feng, CHEN Jiale, SU Zhaopin, et al. Audio spoofing detection using constant-q spectral sketches and parallel-attention SE-ResNet[C]//*Proceedings of Computer Security—ESORICS 2022: 27th European Symposium on Research in Computer Security*. Copenhagen: Springer, 2022: 756-762.
- [41] KHAN A, MALIK K M. SpoTNet: A spoofing-aware transformer network for effective synthetic speech detection[C]// *Proceedings of the 2nd ACM International Workshop on Multimedia AI against Disinformation*. Thessaloniki: Association for

Computing Machinery, 2023: 10-18.

- [42] TAK H, PATINO J, NAUTSCH A, et al. Spoofing attack detection using the non-linear fusion of sub-band classifiers[EB/OL]. (2020-05-20). <https://arxiv.org/pdf/2005.10393>.
- [43] WANG X, YAMAGISHI J. A comparative study on recent neural spoofing countermeasures for synthetic speech detection [EB/OL]. (2021-06-13). <https://arxiv.org/pdf/2103.11326>.
- [44] TAK H, PATINO J, TODISCO M, et al. End-to-end anti-spoofing with RawNet2[C]//Proceedings of 2021 IEEE International Conference on Acoustics, Speech and Signal Processing. Toronto: IEEE, 2021: 6369-6373.
- [45] HUA G, TEOH A B J, ZHANG H J. Towards end-to-end synthetic speech detection[J]. IEEE Signal Processing Letters, 2021, 28: 1265-1269.

作者简介:



袁程胜(1989-),男,博士,副教授,研究方向:信息隐藏、多媒体内容安全等, E-mail: yuancs@nuist.edu.cn。



张雪原(2000-),男,硕士研究生,研究方向:人工智能安全。



周志立(1982-),男,博士,教授,研究方向:信息安全、区块链安全等。



李欣亭(1991-),通信作者,女,博士,副教授,研究方向:数据挖掘、情报分析等, E-mail: lixt@tju.edu.cn。



付章杰(1983-),男,博士,教授,研究方向:信息安全、区块链安全等。

(编辑:夏道家)