

CE-DF-Scanner: 基于ConvNext和EfficientNet的AI生成人脸检测系统

摘要

随着深度学习技术的飞速发展，生成对抗网络（GAN）等人工智能模型在生成高度逼真的人脸图像方面取得了突破性进展。然而，这些技术的滥用也带来了严重的社会和安全隐患，特别是在虚假信息传播和身份盗用等方面。为了应对这一挑战，本研究提出了CE-DF-Scanner，一个基于ConvNext和EfficientNet模型的深度伪造检测系统。该系统通过巧妙结合两种先进的卷积神经网络模型，能够有效地区分真实人脸图像和AI生成的合成图像。本文详细介绍了CE-DF-Scanner的系统架构、数学模型、数据集准备、模型训练策略以及全面的评估方法。实验结果表明，该系统在检测AI生成的人脸方面表现出色，在准确率、精确率、召回率和F1分数等多个评估指标上均达到了极高的水平。

1. 引言

近年来，深度学习技术，尤其是生成对抗网络（GAN）的快速发展，使得生成高质量的虚假图像变得越来越容易。这些生成图像在视觉上与真实图像几乎无法区分，给社会带来了诸多潜在威胁，如虚假新闻传播、身份盗用、金融欺诈等。因此，开发高效、准确的深度伪造检测系统成为了一个迫切需要解决的重要研究课题。

本研究提出的CE-DF-Scanner系统，巧妙地结合了ConvNext和EfficientNet两种先进的卷积神经网络模型。ConvNext是一种新型的卷积神经网络架构，通过引入多尺度特征融合和注意力机制，具有较强的特征提取能力。EfficientNet则通过高效的模型缩放策略，在保持高精度的同时显著减少了计算资源的消耗。通过融合这两种模型的优势，CE-DF-Scanner能够在检测深度伪造图像方面表现出色，同时保持较低的计算复杂度。

2. 系统架构

CE-DF-Scanner的整体架构设计包括数据预处理模块、特征提取模块、分类模块和后处理模块。每个模块都经过精心设计，以实现最佳的检测性能。

2.1 数据预处理

数据预处理是深度学习模型训练过程中的关键步骤，对模型的最终性能有着重要影响。在CE-DF-Scanner中，我们采用了以下预处理步骤：

- 图像尺寸调整**：将输入图像统一调整为384x384像素的分辨率。这一尺寸在保证图像细节的同时，能够有效减少计算资源的消耗。
- 标准化处理**：对图像进行标准化，将像素值缩放到[0,1]范围内，并进行均值和标准差的归一化。具体使用的均值和标准差为：
 - 均值：[0.485, 0.456, 0.406]
 - 标准差：[0.229, 0.224, 0.225]

标准化的数学表达式如下：

$$x_{normalized} = \frac{x - \mu}{\sigma}$$

其中， x 是原始像素值， μ 是均值， σ 是标准差。

- 数据增强**：在训练过程中，我们采用了随机水平翻转的数据增强技术，以增加数据的多样性，有效防止模型过拟合，提高模型的泛化能力。

2.2 特征提取

特征提取模块是CE-DF-Scanner的核心部分，主要由ConvNext和EfficientNet模型组成。

2.2.1 ConvNext

在本系统中，我们使用了ConvNext-Tiny变体。ConvNext的核心思想是通过改进卷积操作来提高模型性能。其主要特点包括：

- 深度可分离卷积**：使用深度可分离卷积替代标准卷积，减少参数量和计算复杂度。深度可分离卷积的数学表达式如下：

$$Y = (X * K_d) * K_p$$

其中， X 是输入特征图， K_d 是深度卷积核， K_p 是逐点卷积核。

- 注意力机制**：引入自注意力机制，增强模型对重要特征的关注。自注意力机制可以表示为：

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

其中， Q 、 K 、 V 分别是查询、键和值矩阵， d_k 是键的维度。

2.2.2 EfficientNet

我们选用了EfficientNet-B0变体。EfficientNet的核心思想是通过复合缩放来平衡网络的深度、宽度和分辨率。其缩放公式如下：

$$\begin{aligned}depth : d &= \alpha^\phi \\width : w &= \beta^\phi \\resolution : r &= \gamma^\phi\end{aligned}$$

其中， ϕ 是缩放系数， α 、 β 、 γ 是通过网格搜索确定的常数。

2.3 分类

分类模块通过全连接层和Sigmoid激活函数，将提取的特征进行分类，输出图像为真实或伪造的概率。具体实现步骤如下：

- 特征融合**：将ConvNext和EfficientNet提取的特征进行拼接，形成一个高维特征向量。
- 特征降维**：通过一个全连接层将高维特征向量映射到一个标量输出。全连接层的数学表达式为：

$$y = Wx + b$$

其中， W 是权重矩阵， x 是输入特征， b 是偏置项。

- Sigmoid激活**：使用Sigmoid函数将输出映射到[0,1]区间，表示图像为伪造的概率。Sigmoid函数的数学表达式为：

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

在训练过程中，我们采用二元交叉熵损失函数（BCEWithLogitsLoss）来衡量模型的分​​类误差。其数学表达式为：

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\sigma(x_i)) + (1 - y_i) \log(1 - \sigma(x_i))]$$

其中， N 是样本数量， y_i 是真实标签， x_i 是模型输出。

2.4 后处理

后处理模块对分类结果进行进一步处理，以提高结果的可解释性和实用性。主要包括以下步骤：

- 阈值判断**：对分类概率进行阈值判断，如果某个图像的伪造概率超过预设阈值（例如0.5），则判定其为伪造图像。
- 结果输出**：将判断结果保存为CSV文件，包含文件名和对应的分类结果（0表示真实，1表示AI生成）。

3. 数据集准备

在本研究中，我们使用了一个包含真实和AI生成人脸图像的数据集。数据集的组织结构如下：

- 数据目录
 - 真实人脸子目录
 - AI生成人脸子目录

我们的FaceDataset类负责加载和预处理这些图像。主要的数据集准备步骤包括：

- 数据加载：**从指定目录加载真实和AI生成的人脸图像。
- 标签分配：**为真实人脸分配标签0，为AI生成人脸分配标签1。
- 数据划分：**将数据集按照8:2的比例随机划分为训练集和验证集。
- 数据增强：**对训练集应用随机水平翻转的数据增强技术。

4. 模型训练

在模型训练过程中，我们采用了一系列先进的优化策略和技术，以提高训练效率和模型性能。

4.1 优化器选择

我们选择了Adam优化器来优化模型参数，初始学习率设置为0.0005。Adam优化器的更新规则如下：

$$\begin{aligned}m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\\hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\\hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \\\theta_t &= \theta_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}\end{aligned}$$

其中， m_t 和 v_t 分别是梯度的一阶矩和二阶矩估计， β_1 和 β_2 是衰减率， α 是学习率， ϵ 是一个小常数。

4.2 学习率调度

我们采用了余弦退火学习率调度策略（CosineAnnealingLR），该策略能够在训练初期保持较高的学习率以快速收敛，在训练后期降低学习率以微调模型参数。学习率的计算公式如下：

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{t}{T}\pi))$$

其中， η_t 是第 t 个epoch的学习率， η_{min} 和 η_{max} 分别是最小和最大学习率， T 是总的训练epoch数。

4.3 训练过程

模型训练持续20个epoch，每个epoch后都会在验证集上评估模型性能，并保存模型检查点。我们还实现了混合精度训练，以加速训练过程并减少内存占用。混合精度训练使用FP16（半精度浮点数）进行部分计算，同时保持FP32（单精度浮点数）的模型权重，从而在不损失精度的情况下提高训练速度。

5. 实验结果与分析

我们在384x384像素的图像数据集上进行了全面的实验评估。以下是详细的实验结果：

- 损失 (Loss) : 0.0173
- 准确率 (Accuracy) : 0.9941
- 精确率 (Precision) : 0.9985
- 召回率 (Recall) : 0.9942
- F1分数 (F1 Score) : 0.9964

这些评估指标的计算公式如下：

- 准确率: $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$
- 精确率: $Precision = \frac{TP}{TP+FP}$
- 召回率: $Recall = \frac{TP}{TP+FN}$
- F1分数: $F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision+Recall}$

其中，TP是真正例，TN是真负例，FP是假正例，FN是假负例。

这些结果表明，CE-DF-Scanner在检测AI生成的人脸方面表现出色，具有极高的准确性和可靠性。特别是在精确率方面的出色表现，意味着该系统在实际应用中能够最大限度地减少误报，这对于维护社会信任和防止错误指控至关重要。

6. 结论

本文提出了CE-DF-Scanner，一个基于ConvNext和EfficientNet模型的深度伪造检测系统。通过巧妙结合两种先进的卷积神经网络模型，CE-DF-Scanner在检测AI生成的人脸方面取得了显著效

果。实验结果表明，该系统在准确率、精确率、召回率和F1分数等多个评估指标上均达到了极高的水平。