

中图分类号: D918.91 文献标识码: A 文章编号: 1006-8961(XXXX)XX-0001-20

论文引用格式: Guo Tianli, Li Jisong, Tang Yunqi. Advances in Forensic Science of Deepfake Facial Images Evidence[J/OL]. Journal of Image and Graphics, XXXX: 1-20. DOI: 10.11834/jig.250610. (郭甜利, 李纪松, 唐云祁. 深度伪造人脸图像证据的司法鉴定研究进展[J/OL]. 中国图象图形学报, XXXX: 1-20. DOI: 10.11834/jig.250610. ) [DOI: 10.11834/jig.250610]

## 深度伪造人脸图像证据的司法鉴定研究进展

郭甜利, 李纪松, 唐云祁\*

中国人民公安大学侦查学院, 北京 100038

**摘要:** 随着生成对抗网络(GAN)、自动编码器(AE)和扩散模型(DM)等深度伪造技术的快速发展,深度伪造人脸图像被广泛用于制作非自愿色情内容、电信诈骗及传播政治虚假信息等恶意行为。作为新型数字证据,深度伪造不仅频繁出现在当前犯罪活动中,还可能被用于虚假指控或操控监控记录,严重削弱法庭证据科学中人脸图像分析的可靠性,进而威胁刑事司法的公正性。本文系统梳理了关键数据集的发展脉络与深度伪造人脸图像的技术演进路径,深入剖析其生成机制的核心架构,涵盖主要技术路线、理论基础与功能逻辑;基于Web of Science数据库筛选出的737篇文献,采用VOSviewer进行文献计量分析,识别该领域的研究趋势,并揭示现有检测技术与司法实践中证据鉴伪实际需求之间的显著脱节。本研究提出将深度伪造人脸图像检测方法划分为技术原理导向与应用阶段适配两类互补框架,以应对不同阶段的证据鉴定需求。针对深度伪造人脸图像作为新型犯罪载体所带来的证据可信性挑战,本文主张法庭证据鉴定人应推动检测向鉴定的转化,实现检测结果与贝叶斯推理框架的融合,提升其在司法证明中的可采性与说服力。本研究衔接计算机视觉技术发展与法庭证据实践需求,明确了从技术检测到司法鉴定的逻辑路径,为法庭科学研究人员与司法实务工作者提供了具有技术价值的理论参考。

**关键词:** 司法;证据;人脸图像;深度伪造;检测;鉴定

## Advances in Forensic Science of Deepfake Facial Images Evidence

Guo Tianli, Li Jisong, Tang Yunqi\*

School of Criminal Investigation, People's Public Security University of China, Beijing 100038, China

**Abstract:** With the rapid advancement of deepfake technologies—including Generative Adversarial Networks (GANs), Autoencoders (AEs), and Diffusion Models (DMs)—deepfake facial images have been increasingly exploited for malicious purposes, such as the creation of non-consensual pornographic content, telecommunication fraud, and the propagation of political disinformation. In recent years, the continuous iteration of these core deepfake technologies has significantly lowered the technical threshold for generating high-fidelity fake facial images: GANs, with their adversarial training mechanism between generators and discriminators, enable the rapid synthesis of realistic facial features; AEs, relying on unsupervised learning and feature reconstruction capabilities, excel in simulating subtle facial expressions and skin textures; Diffusion Models, through step-by-step noise reduction and feature enhancement, have further broken through the limitations of traditional deepfake technologies in terms of image clarity and detail authenticity, making deepfake facial images

收稿日期: 2025-12-02; 修回日期: 2026-02-13

\*通信作者: 唐云祁, 通信作者, 男, 教授, 主要研究方向为电子数据检验。E-mail: tangyunqi@ppsuc.edu.cn

基金项目: 中国人民公安大学刑事科学技术双一流创新研究项目(项目编号: 2023SYL06); 中国人民公安大学研究生科研项目(项目编号: 2024yjsky032)

Supported by: Double First-Class Innovation Research Project of Criminal Science and Technology, People's Public Security University of China (Project No.: 2023SYL06); Graduate Research Project of People's Public Security University of China (Project No.: 2024yjsky032)

©中国图象图形学报版权所有

increasingly difficult to distinguish with the naked eye and even challenging preliminary manual identification in forensic investigations. As an emerging form of digital evidence, deepfakes not only pervade contemporary criminal activities but also risk being weaponized for false accusations or the tampering of surveillance records, thereby severely undermining the reliability of facial image analysis in forensic science and threatening the integrity of the criminal justice system. To address the aforementioned challenges, this paper systematically traces the developmental trajectory of key datasets and the technological evolution of deepfake facial images, while conducting an in-depth dissection of the core architecture underlying their generation mechanisms, encompassing primary technical routes, theoretical underpinnings, and functional logics. In terms of dataset development, the evolution process from early small-scale, single-scene datasets (such as CelebA and LFW) to large-scale, multi-scenario, high-diversity datasets (such as DFDC, FF++, and Celeb-DF) is sorted out, with the advantages and limitations of each type of dataset analyzed in terms of sample size, fake generation technology coverage, and scene diversity. It is further pointed out that the lack of standardized, forensic-oriented datasets has become one of the key bottlenecks restricting the practical application of deepfake detection technologies. In terms of technological evolution, the development of deepfake facial image technology is divided into three stages: the initial exploration stage dominated by GANs, the improvement stage supplemented by AEs, and the mature stage driven by Diffusion Models. The theoretical breakthroughs and technical improvements of each stage are elaborated, and the internal logic of the continuous improvement of fake image fidelity and the continuous reduction of generation costs is revealed. In the in-depth analysis of generation mechanisms, the core principles of the three mainstream technologies are focused on, the differences in their technical routes and functional characteristics are clarified, and the common key links and potential vulnerabilities in the generation process are summarized, laying a theoretical foundation for the subsequent design of targeted detection and identification methods. Based on 737 relevant articles retrieved and screened from the Web of Science database, a bibliometric analysis is performed via VOSviewer to identify research trends in the field and illuminate the substantial disconnect between existing detection technology research and the practical requirements for evidence authentication in judicial practice. The retrieval and screening process strictly followed the PRISMA guidelines: "deepfake", "facial image", "forensic science", "identification", and "detection" are adopted as key search terms, the retrieval time range is defined from the emergence of deepfake technology to the present, and irrelevant literature, review papers with low innovation, and duplicate publications are excluded through preliminary screening and full-text reading. The bibliometric analysis focuses on three core dimensions: research hotspots, institutional cooperation, and research frontiers. The results show that current research in the field is mainly concentrated on the optimization of detection algorithms (such as deep learning-based feature extraction and classification models), while research on the translation of detection results into forensic authentication, the formulation of forensic identification standards, and the adaptation of detection technologies to judicial practice scenarios is relatively scarce. This disconnect is mainly reflected in two aspects: on the one hand, most existing detection technologies focus on technical accuracy, ignoring the admissibility requirements of judicial evidence (such as objectivity, authenticity, and relevance); on the other hand, forensic examiners lack effective methods to integrate technical detection results into judicial authentication, leading to the situation where advanced detection technologies fail to be effectively applied in judicial practice. In response to the identified research-practice disconnect, this study proposes a classification of deepfake facial image detection methods into two complementary frameworks: technology-principle-oriented and application-scenario-adaptive, designed to address evidence authentication needs across distinct stages. The technology-principle-oriented framework takes the technical principles of deepfake generation as the core, classifies detection methods in accordance with the technical characteristics of different generation technologies (such as GAN-based detection, AE-based detection, and Diffusion Model-based detection), and focuses on improving the accuracy and generalization ability of detection, thus being mainly applicable to the preliminary technical screening stage of evidence in forensic investigations. The application-scenario-adaptive framework takes the actual needs of different judicial stages (such as investigation, prosecution, and trial) as guidance, classifies detection methods in light of the scenario characteristics and evidence requirements of each stage, and focuses on enhancing the adaptability and operability of detection methods, thereby being mainly applicable to the formal evidence authentication stage in judicial practice. The two frameworks are complementary and mutually supportive: the technology-principle-oriented framework provides technical support for the application-scenario-adaptive

framework, and the application-scenario-adaptive framework guides the direction of technical optimization, thus forming a complete detection system covering the entire process of forensic evidence authentication. In response to the evidentiary credibility challenges posed by deepfake facial images as a novel criminal vector, this paper argues that forensic examiners should facilitate the translation of detection outcomes into forensic authentication, integrating detection results with the Bayesian inference framework to enhance their admissibility and probative force in judicial proceedings. Specifically, a Bayesian inference-based forensic authentication model is proposed: technical detection results are taken as the prior probability, forensic examiners' professional experience and case context are integrated to determine the likelihood ratio, and the posterior probability of the authenticity of facial images is ultimately calculated via the Bayesian formula, so as to convert objective technical indicators into credible judicial authentication conclusions. Furthermore, it is emphasized that forensic examiners should continuously improve their understanding of deepfake technologies, master the working principles and limitations of detection technologies, and avoid over-reliance on technical detection results, so as to ensure the scientificity and rigor of forensic authentication. In conclusion, by bridging the advancements in computer vision technology with the practical demands of forensic evidence practice, this research articulates the logical pathway from technical detection to judicial authentication, providing a theoretically insightful and practically relevant reference for forensic science researchers and judicial practitioners alike. For forensic science researchers, this study clarifies the research direction of integrating technical research with judicial practice, identifies the key issues to be addressed in the future (such as the construction of forensic-oriented datasets, the optimization of detection algorithms adapted to judicial scenarios, and the formulation of forensic identification standards), and offers a theoretical framework for subsequent research. For judicial practitioners, this study provides a practical operational path for the forensic authentication of deepfake facial images, facilitates the better grasp of the application methods and precautions of deepfake detection technologies, and elevates the efficiency and accuracy of digital evidence authentication. Additionally, this study offers a reference for the formulation of relevant laws and regulations, which helps to improve the legal system of digital evidence authentication and promote the sound development of the criminal justice system in the era of deepfake technology.

**Key words:** forensic science; evidence; facial image; deepfake; detection; identification

## 0 引言

“深度伪造”作为“深度学习”与“伪造技术”的融合产物,特指通过深度学习模型生成或篡改图像、视频、音频等数字内容的技术过程,其中人脸替换是最典型的表现形式,此外还包括人脸合成以及视频音频伪造等衍生形态。2017年,美国社交平台Reddit上首次出现名人面部篡改内容的大规模传播,标志着深度伪造技术正式进入公众视野(Dayal S. B. 等, 2021)。经过多年迭代,该技术的影响力已渗透至社会治理、刑事司法等关键领域。2024年《自然》杂志将深度伪造检测列为年度值得关注的七大技术之一(Eisenstein M. 等, 2024),凸显了其全球范围内的技术重要性与治理紧迫性。

深度伪造技术的滥用已引发多维度、跨领域的严重危害。

在色情与隐私侵犯领域,韩国占全球深度伪造色情内容的53%,99%的受害者为女性,从教师裸照

被传播到22万成员的Telegram群组批量生成非法内容(Ji S., 2025),相关黑色产业链已形成规模化;西班牙女学生也遭遇了此类网络欺凌(Narvali A. M. 等, 2024)。

在金融领域,诈骗者通过模仿高管语音、伪造投资广告等手段实施犯罪。相关案件造成了严重损失:2019年虚假信息引发英国大都会银行挤兑潮,导致股价下跌9%(Van Der Sloot B. 等, 2022),且单起案件授权转账金额最高达2560万美元(Brown A. 等, 2025);2023-2025年间,英国储户相关诈骗案涉案金额达3500万美元(Corbett R., 2025),企业及金融部门年度损失均超数十万美元(Patishman 等, 2024),欧洲、加拿大储户及西班牙相关诈骗案的涉案金额达2000万美元(Katte 等, 2025)。

在政治与国际关系层面,伪造领导人言论、外交虚假图像、战争宣传视频等内容已成为地缘政治博弈工具,相关案例持续加剧国际或地区紧张局势:特朗普、奥巴马相关虚假言论(Pesetski A., 2020; Ramluckan T., 2024),意大利前总理伦齐、加蓬总统的伪



造视频争议(Venema A. E. 等,2020;Schiff K. J. 等,2025),外交虚假图像(Alhajjar E. 等,2022),2024年模仿拜登的选民压制语音机器人(Federal Communications Commission,2024),战争宣传视频(Bohávcek M. 等,2022),以及2025年印巴冲突中的伪造“道歉”视频(Mahajan A. S.,2025)均在其列。司法实践中,伪造音频还曾干扰监护权判决(Gabriella Swerling,2020)。

在社会层面,高中学生因伪造色情视频遭受社交与心理损害(Henry Ajder等,2019),英国记者成为深度伪造攻击目标(Kaupins G.,2025);

在军事与国家安全领域,篡改卫星图像导致军事误判(Çiftçi U. A. 等,2023);技术安全层面,78%的深度伪造内容成功欺骗微软 Azure 实验性服务 API,暴露了商业面部识别系统的脆弱性(Tariq S. 等,2021);

在儿童权益保护领域,2023-2024年人工智能生成的儿童性虐待材料数量翻两番,勒索案件同步增加(Krishna S. 等,2024)。

从法庭证据科学视角来看,洛卡德物质交换原理作为刑事侦查的核心理论,其内涵已从物理环境延伸至数字空间——任何数字篡改行为都会留下可追溯的技术痕迹(Kaur M. 等,2024;Iqbal A. 等,2020;Amundsen A. E. 等,2017)。但与指纹(Levanon L. 等,2025)、脚印、DNA、声纹、面部特征、虹膜图案、步态动态(Deng W. 等,2022)等传统生物识别证据不同,深度伪造人脸图像证据具有显著的数字特性:一是数字无形性,痕迹仅存在于数字系统中,无物理载体;二是技术异质性,生成模型与操纵算法的多样性导致伪造结果呈现高度差异化;三是传播突发性,数字内容可在短时间内实现大规模扩散,放大社会危害;四是鉴别复杂性,对抗性生成机制使人工鉴别可靠性极低,必须依赖自动化检测工具;五是刑事侦查模糊性,现有刑事侦查与司法框架中缺乏针对该类数字证据的可采性标准与证据权重评估规范。这些特性对法庭证据科学提出了全新的理论与实践挑战,亟需建立适配数字环境的证据分析与评估体系。

围绕深度伪造的生成机制与检测技术,现有研究已开展了大量探索。Seow 等(Seow J. W. 等,2022)系统梳理了深度伪造的生成技术、检测方法、数据集及研究空白,为领域发展提供了全面参考;

Ramanathan 等(Ramanaharan R. 等,2025)深入探讨检测模型的泛化能力,指出当前研究存在过度拟合数据集等局限性,并提出开发更逼真数据集的未来方向;Garg 等(Garg D. 等,2025)指出当前研究缺乏多模态方法、对混合模型探索不足的问题,进而提出多模态混合学习检测框架;Abbas 等(Abbas F. 等,2024)探讨了基于人工智能的深度伪造检测与生成技术,阐明了人工智能在该领域的应用潜力,并提出相关政策建议。然而,现有研究仍存在显著短板:其一,研究视角集中于计算机视觉领域,缺乏与法庭证据科学的跨学科整合,未能回应司法实践对证据评估的核心需求;其二,检测技术分类体系零散,未形成适配法庭证据应用阶段的结构化框架;其三,对深度伪造证据的定量评估方法研究不足,难以支撑刑事诉讼程序中的证据采信。

法庭科学作为整合生物学、物理学、化学和计算机科学知识的应用学科,核心目标是对 DNA(生物类证据)、指纹(形态类证据,含足迹、枪弹痕迹等)、电子证据(含图像、音频、视频等电子数据)等进行定量研究与概率评估,为刑事侦查和刑侦诉讼程序提供支撑。针对现有研究缺口与刑事司法实践需求,本研究以法庭证据科学核心逻辑为切入点,厘清深度伪造技术的核心原理与发展趋势,建立适配法庭证据场景的检测方法分类体系,提出科学可行的深度伪造人脸图像证据的鉴定可行性方案。

## 1 深度伪造人脸

深度伪造技术的演进路径,也是其作为犯罪载体的发展过程。

### 1.1 深度伪造人脸数据集

从时间维度看,深度伪造技术的发展呈现出多样化、精细化和向真实场景演进的趋势。详见图1。

2018年为技术奠基阶段,核心生成框架初步建立,研究聚焦低分辨率人脸合成与基础换脸技术,主要采用生成对抗网络(GAN)和自动编码器(AE)。例如,DeepFakeTIMIT数据集利用GAN结合卷积神经网络(CNN)分割与MTCNN检测的面部特征点,实现初步换脸。UADFV则通过自动编码器配合仿射变换和边界平滑减少拼接痕迹。但由于依赖单一模型,伪造结果存在边缘模糊、表情僵硬等明显视觉

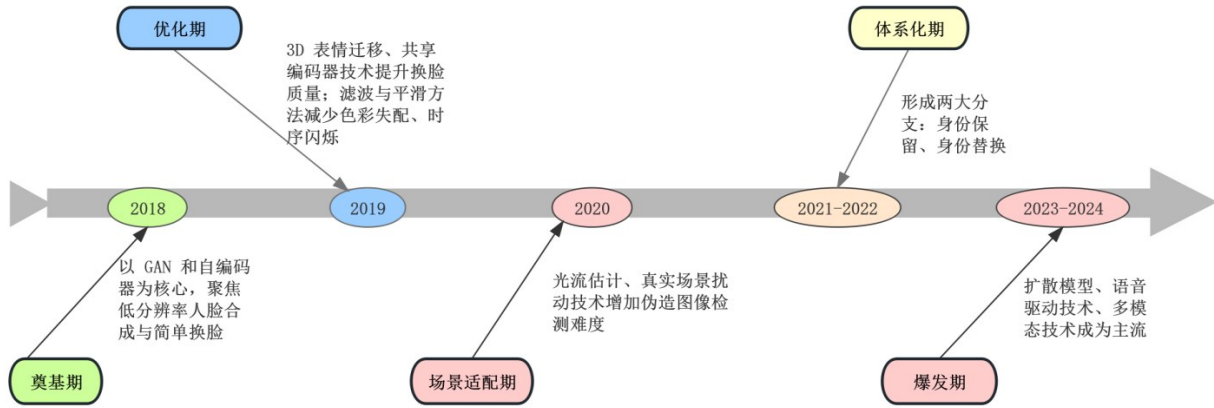


图1 数据集时间发展的鱼骨图

Fig. 1 Fishbone diagram of the temporal development of the dataset

缺陷。

2019年进入技术分化期,转向多方法融合与细节优化。FaceForensics++集成Face2Face(3D表情迁移)和DeepFakes(身份替换)等多种技术,并通过不同比特率的压缩模拟真实传播环境。Celeb-DF-v2引入色彩校正与卡尔曼滤波,有效缓解色偏和画面闪烁,显著提升视频的时间连贯性与真实感。

2020年迈向实际应用,强调合成内容与真实环境的融合。DF-1.0采用光流估计保障时间一致性,并加入JPEG压缩、噪声、运动模糊等七类现实失真,更贴近数字媒体的实际退化情况。WildDeepFake收集来自网络平台的真实深度伪造视频,涵盖图像拼接与生成合成等多种手段,增强生态效度,贴近司法实践中面临的复杂证据形态。2021至2022年,深

度伪造方法体系趋于系统化,逐渐形成“身份保留”和“身份替换”两大伪造模式。KoDF整合FaceSwap(身份替换)与FOMM(面部重演,身份保留),并引入FGSM等对抗扰动以增强抗检测能力。ForgeryNet明确分类操纵技术,区分StarGAN2类属性编辑与FaceShifter类高保真身份迁移,构建了结构化评估体系。2023至2024年,前沿伪造方法与基准集中涌现,扩散模型与多模态融合逐步成为主流。DeepFakeFace基于Stable Diffusion进行局部精细编辑,结合ControlNet实现像素级操控;DF40基准收录多种前沿伪造技术,覆盖StyleGAN3全脸生成与Sad-Talker音频驱动唇动同步。这些进展标志着深度伪造进入高保真、强适应、低残留的新阶段,对电子数据检验与数字取证工作构成前所未有的挑战。

表1 2018年至2024年深度伪造人脸的代表性数据集

Table 1 Representative datasets of deepfake faces from 2018 to 2024

| 年份   | 数据集  | 大小   | 类型                | 方法                       | 机构              | 论文                          | 链接  |
|------|--|--|-------------------|--------------------------|-----------------|-----------------------------|---|
| 2018 | DeepFakeTIMIT                              | 总计 640 段(分<br>高清、标清版<br>本)                               | 视频                | GAN 驱动                   | 艾迪<br>亚研<br>究所  | (Korshunov P.<br>等, 2018)   | <a href="https://www.idiap.ch/dataset/DeepFakeTimit">https://www.idiap.ch/dataset/DeepFakeTimit</a>   |
| 2018 | UADFV                                      | 真实 49 段, 伪<br>造 49 段                                     | 视频                | 检测框定、<br>仿射变换与<br>边界平滑   | 奥尔巴<br>尼大<br>学等 | (Yang X. 等,<br>2019)        | <a href="https://docs.google.com/forms/d/e/1FAIpQLScKP0Ov15TIZ9Mn0nGScIVgKRM9tFW0mj9eHKx57Yp-Xcnx/viewform">https://docs.google.com/forms/d/e/1FAIpQLScKP0Ov15TIZ9Mn0nGScIVgKRM9tFW0mj9eHKx57Yp-Xcnx/viewform</a> |
| 2018 | Fake Face in the<br>Wild FFW data-<br>base | 伪造图像 5.3<br>万张(150 段视<br>频), 真实图像<br>7.85 万张(150<br>段视频) | 视<br>频/<br>图<br>像 | 编<br>码<br>器、手<br>动篡改、CGI | 奥勒松<br>大学等      | (Khodabakhsh<br>A. 等, 2018) | <a href="http://ali.khodabakhsh.org/ffw/">http://ali.khodabakhsh.org/ffw/</a>   |

表1续表

| 年份   | 数据集                             | 大小                                  | 类型    | 方法              | 机构         | 论文                     | 链接  |
|------|---------------------------------|-------------------------------------|-------|-----------------|------------|------------------------|---|
| 2018 | 100K-Faces                      | 伪造图像 10 万张                          | 图像    | StyleGAN        | 英伟达        | (Almars A.M., 2021)    | <a href="https://generated.photos/">https://generated.photos/</a>   |
| 2018 | PGGAN                           | 伪造图像 8 万张                           | 图像    | PGGAN           | 英伟达        | (Karras T. 等, 2018)    | <a href="https://github.com/tkarras/progressive_growing_of_gans">https://github.com/tkarras/progressive_growing_of_gans</a>               |
| 2019 | FaceForensics++                 | 图像超 180 万张 (1000 段原始视频, 4000 段伪造视频) | 图像/视频 | 3D 建模、编码        | 慕尼黑工业大学等   | (Rossler A. 等, 2019)   | <a href="https://github.com/ondyari/FaceForensics">https://github.com/ondyari/FaceForensics</a>   |
| 2019 | TPDNEThis Person Does Not Exist | 图像 15 万张                            | 图像    | StyleGAN+潜空间生成  | 英伟达        | (Karras T. 等, 2019)    | <a href="https://thispersondoesnotexist.com">https://thispersondoesnotexist.com</a>   |
| 2019 | Celeb-DF-v2                     | 真实 590 段, 伪造 5639 段                 | 视频    | 编解码             | 奥尔巴尼大学等    | (Li Y. 等, 2020)        | <a href="http://www.cs.albany.edu/~lsw/celeb-DeepFakeforensics.html">http://www.cs.albany.edu/~lsw/celeb-DeepFakeforensics.html</a>       |
| 2019 | Celeb-DFv1                      | 真实视频 408 段; 深度伪造视频 795 段            | 视频    | 传统伪造            | 奥尔巴尼大学等    | (Li Y. 等, 2019)        | <a href="https://github.com/danmohaha/celeb-DeepFakeforensics">https://github.com/danmohaha/celeb-DeepFakeforensics</a>                   |
| 2019 | DFDC                            | 总计 12.815 万段                        | 视频    | 8 类伪造方法锐化, 泊松融合 | 脸书等        | (Dolhansky B. 等, 2020) | <a href="https://ai.facebook.com/datasets/dfdc">https://ai.facebook.com/datasets/dfdc</a>   |
| 2020 | WildDeepFake                    | 伪造视频 707 段, 人脸图像 118 万张             | 视频/图像 | 野生采集            | 复旦大学等      | (Zi B. 等, 2020)        | <a href="https://github.com/OpenTAI/wild-DeepFake">https://github.com/OpenTAI/wild-DeepFake</a>   |
| 2020 | DF-1.0 Deeper-Forensics-1.0     | 真实 5 万段, 伪造 1 万段                    | 视频    | 光流时间一致性、真实场景扰动  | 南洋理工大学等.   | (Jiang L. 等, 2020)     | <a href="https://liming-jiang.com/projects/DrF1/DrF1.html">https://liming-jiang.com/projects/DrF1/DrF1.html</a>                           |
| 2020 | iFakeFaceDB                     | 伪造图像 8.7 万张                         | 图像    | 卷积自编码器去除指纹      | 葡萄牙贝拉内大学   | (Neves J.C. 等, 2020)   | <a href="https://github.com/socialabusi/iFake-FaceDB">https://github.com/socialabusi/iFake-FaceDB</a>                                     |
| 2021 | OpenForensics                   | 真实 16.0676 万张, 伪造 17.3660 万张        | 图像    | 潜向量修改           | 日本国立信息研究所  | (Le T.-N. 等, 2021)     | <a href="https://sites.google.com/view/ltngghia/research/openforensics">https://sites.google.com/view/ltngghia/research/openforensics</a> |
| 2021 | KoDF                            | 403 名亚洲志愿者; 真实 6.22 万张, 伪造 17.58 万张 | 视频    | 6 类核心合成模型       | 钱脑公司       | (Kwon P. 等, 2021)      | <a href="https://moneybrain-research.github.io/kodf">https://moneybrain-research.github.io/kodf</a>                                       |
| 2021 | ForgeryNet                      | 图像 290 万张, 视频 22.1247 万段            | 视频/图像 | FOMM, StarGAN2  | 中国商汤科技研究院等 | (He Y. 等, 2021)        | <a href="https://yinanhe.github.io/projects/forgerynet.html">https://yinanhe.github.io/projects/forgerynet.html</a>                       |
| 2021 | FakeAVCeleb                     | 真实 500 段, 伪造 1.95 万                 | 视频    | 人脸相似度           | 韩国成均馆大学等   | (Khalid H. 等, 2021)    | <a href="https://sites.google.com/view/fakeavcelebdash-lab/">https://sites.google.com/view/fakeavcelebdash-lab/</a>                       |
| 2021 | DFGC-21 testing dataset         | 真实 1000 张, 伪造 1000 张                | 图像    | 仿射              | 深圳大学等      | (Peng B. 等, 2021)      | <a href="https://github.com/yuezunli/celeb-DeepFakeforensics">https://github.com/yuezunli/celeb-DeepFakeforensics</a>                     |

表1续表

| 年份   | 数据集  | 大小  | 类型            | 方法   | 机构                                     | 论文                       | 链接  |
|------|--|---|---------------|--|--|--------------------------|---|
| 2022 | Localized Audio<br>Visual DeepFake<br>LAV-DF | 真实 3.6431 万<br>段, 伪造片段<br>9.9873 万段;<br>153 名个体                 | 视频            | 唇形同步   | 澳大利<br>亚莫纳<br>什大学<br>等                 | (Cai Z. 等,<br>2022)      | <a href="https://github.com/ControlNet/LAV-DF">https://github.com/ControlNet/LAV-DF</a>   |
| 2022 | DFDMDeepFakes<br>from Different<br>Models    | 伪造 6450 段   | 视频            | 自编码器   | 美国纽<br>约州立<br>大学布<br>法罗分<br>校等         | (Jia S. 等,<br>2022)      | <a href="https://github.com/shanface33/Deep-Fake_Model_Attribution">https://github.com/shanface33/Deep-Fake_Model_Attribution</a>   |
| 2022 | DeePhy Deep-<br>Fake Phylogeny<br>dataset    | 伪造 5040 段,<br>真实 100 段  | 视频            | 1-3 次人脸<br>换脸                                      | 印度焦<br>特布尔<br>印度理<br>工学院<br>等          | (Narayan K.<br>等, 2022)  | <a href="http://iabrubric.org/deephy-database">http://iabrubric.org/deephy-database</a>   |
| 2022 | GOTCHA                                       | 伪造 409 段, 真<br>实 5.5838 万段                                      | 视频            | FSGAN v2   | 纽约大<br>学等                              | (Mittal G. 等,<br>2024)   | <a href="https://github.com/mittalgovind/GOTCHA-DeepFakes">https://github.com/mittalgovind/GOTCHA-DeepFakes</a>   |
| 2023 | DeepFakeFace<br>DF                           | 真实 3 张, 伪造<br>9 万张  | 图像            | 扩散模型   | 清华大<br>学等                              | (Song H. 等,<br>2023)     | <a href="https://github.com/OpenRL-Lab/Deep-FakeFace/">https://github.com/OpenRL-Lab/Deep-FakeFace/</a> <a href="https://huggingface.co/datasets/OpenRL/DeepFakeFace">https://huggingface.co/datasets/OpenRL/DeepFakeFace</a> |
| 2023 | StyleGAN2-<br>FFHQ                           | 图像 100 张(深<br>度伪造人脸 50<br>张; 真实人脸<br>50 张)                      | 图像            | StyleGAN2  | 伦敦大<br>学等                              | (Bray S.D. 等,<br>2023)   | <a href="https://osf.io/tfn7v/">https://osf.io/tfn7v/</a>   |
| 2023 | Div-DF                                       | 真实 150 张, 伪<br>造 250 张(换脸<br>100 张, 重演<br>100 张; 唇形同<br>步 50 张) | 视频            | FSGAN  | 德里技<br>术大学<br>等                        | (Dagar D. 等,<br>2023)    | <a href="https://forms.gle/WhWDeVoM3wHTo-BYi9">https://forms.gle/WhWDeVoM3wHTo-BYi9</a>   |
| 2023 | DFMD   | 真实 1000 段,<br>伪造 1000 段   | 视频            | FOMM   | 伊玛目<br>阿卜杜<br>勒拉赫<br>曼·本<br>·费萨<br>尔大学 | (Alnaim N.M.<br>等, 2023) | <a href="https://creativecommons.org/licenses/by-nc-nd/4.0/">https://creativecommons.org/licenses/by-nc-nd/4.0/</a>   |
| 2024 | DF40   | 视频片段超 10<br>万段, 图像超<br>100 万张                                   | 视<br>频/<br>图像 | 10 类换脸,<br>12 类人脸重<br>演, 12 类全<br>脸合成, 5 类<br>人脸编辑 | 北京大<br>学等                              | (Yan Z. 等,<br>2024)      | <a href="https://github.com/YZY-stack/DF40">https://github.com/YZY-stack/DF40</a>   |

## 1.2 深度伪造人脸图像中的信息

人脸图像包含丰富信息。如图 2 所示, 原始人脸图像可提取多种特征, 包括语义掩码、RGB 颜色、深度、频域和空域特征等。这些内在属性被生成模型用于生成深度伪造人脸。

尽管不同类型信息的具体表示与特征提取方法较为复杂, 如图 2 中的各子图仍能直观反映不同模态的基本分布与结构。

其中, 语义掩码通过像素级标注实现结构化面部分割, 明确眼、鼻、口等关键面部部位的边缘轮廓与边界, 完成空间定位与语义区分。

颜色与深度信息互补。颜色通道反映表面外观(如色彩分布与纹理), 深度数据则通过空间距离刻画 3D 面部几何特征(如鼻梁凸起、下颌凹陷)。二者共同提供多模态的面部外观描述。

在频域分析中, 高频成分体现皮肤纹理、睫毛等



细节;中频成分连接局部结构与整体形态;低频成分主导全局轮廓与宏观布局(如脸型、五官排列),反映粗粒度结构特征。

空域特征则随维度变化而呈现不同表征能力。在二维空间中,像素梯度与密度可以刻画局部复杂

性,例如眼部区域因纹理更为丰富而表现出更剧烈的梯度变化;在三维空间中,不同坐标轴用于捕捉特定几何属性,z轴聚焦于深度方向的形变,x/y轴描述平面内的轮廓与纹理,从而形成多视角的空间表达。

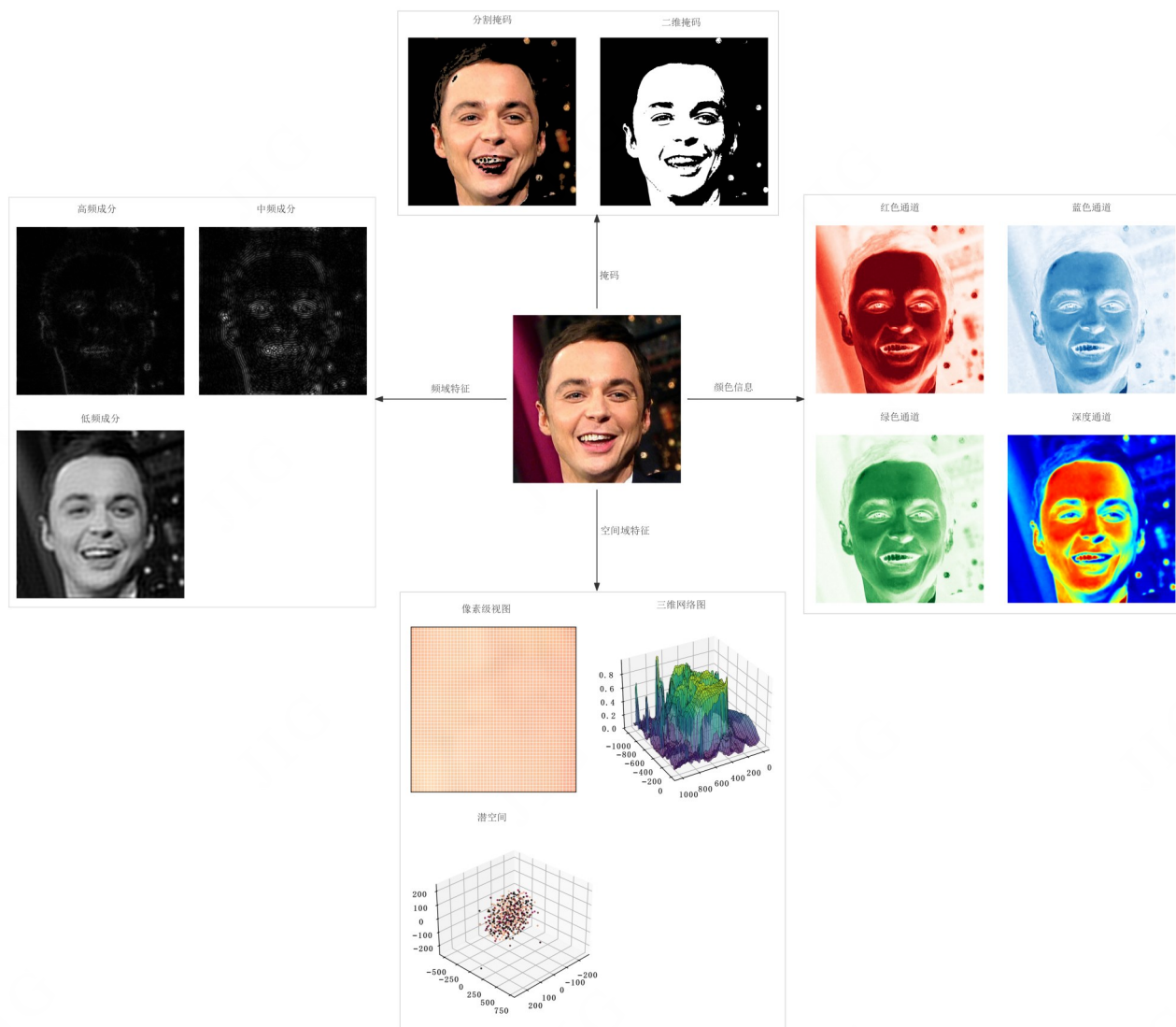


图2 从人脸图像中提取的关键信息

Fig. 2 Key information extracted from facial images

另外,深度神经网络学习到的潜在特征通常是高维、抽象且不可直接可视化的。尽管难以用直观方式加以解读,这些潜在特征却编码了人脸生成所需的全局统计规律与结构先验,是深度伪造人脸合成与编辑的核心表征基础。

### 1.3 深度伪造人脸的技术本质

#### 1.3.1 技术路线

深度伪造人脸的生成流程包含四个核心阶段。

准备阶段。利用深度神经网络从待伪造人脸中提取高维抽象特征。这些特征不同于传统法庭科学关注的几何比例、皮肤纹理等可观察形态,而是对颜色分布、空间结构、频域特性等多维信息的隐式编码,兼具高维度与强抽象性。同时提供待迁移特征的原始图像和伪造结果,并通过面部关键点检测与姿态对齐建立源与目标间的特征映射,为后续建模提供关键约束。



伪造建模阶段。使用生成模型,学习源与目标特征的关联,实现特征的定向迁移或全新生成。例如,将源数据的身份特征融入目标的姿态框架,或在目标约束下生成符合特定属性的新面部特征。

重新渲染阶段。将生成的抽象特征转换为可视图像,通过调整光照一致性、像素级重建、补充皮肤纹理等细节填充,赋予图像符合视觉感知的颜色、纹理与空间结构。

后处理阶段。通常进行扰动增强,如微调高频成分、添加噪声,提升图像真实感或增强逃避反深度伪造人脸检测的能力,最终输出符合预期的伪造结果。

### 1.3.2 技术基础

生成对抗网络及其变体。生成对抗网络(GAN)由生成器和判别器组成,两者在对抗训练中不断提升各自能力,最终使伪造图像在视觉效果上接近真实图像,以至于判别器几乎接近随机猜测。这一机制如图4所示。

自动编码器。自编码器通过编码器从人脸图像中提取潜在特征,并利用解码器将其重构为连贯的人脸图像。在面部替换中,采用两个共享编码器权重的对称编码器-解码器结构。共享编码器确保身份特征的一致性编码,各解码器则在其对应数据集上独立训练。解码路径相互解耦,便于独立优化以提升生成真实性。面部替换效果依赖于姿势、表情和头部形态的一致性。编码器通过分离身份与非身份属性实现特征的精确解缠,双解码器结构则在保持源图像外观上下文的同时重建目标身份。该过程如图5所示。

扩散模型。扩散模型包含两个阶段,正向扩散和反向去噪。在正向过程中,噪声逐步添加到原始人脸图像,经过多步迭代后图像完全变为噪声。该过程是固定的,由预定义的方差调度控制每一步的噪声强度。反向过程中,神经网络从纯噪声出发,逐步重建原始图像。该网络经训练可估计并去除每步噪声,从而实现正向过程的有效逆转。其架构与原理如图6所示。

### 1.3.3 功能逻辑

从功能分类来看,深度伪造人脸图像主要包括面部重演、面部编辑、面部替换与面部合成四类,其核心差异体现在对身份信息的处理逻辑上。面部重演聚焦面部动态迁移,将源个体的头部姿势、表情、

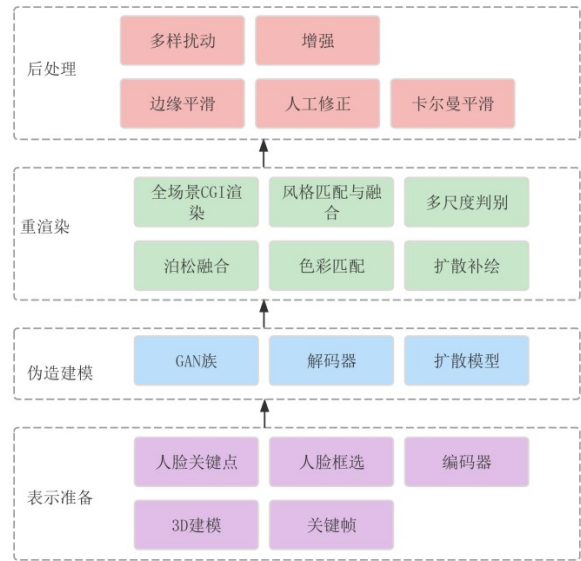


图3 生成深度伪造人脸图像的流程

Fig. 3 Deepfake Image Generation Process

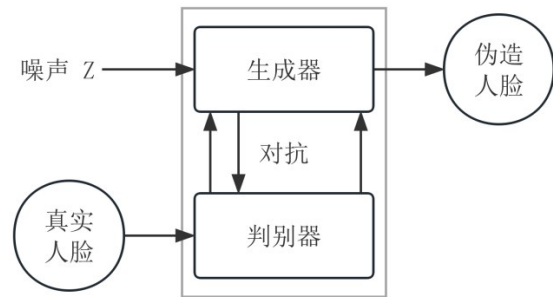


图4 基于生成对抗网络的伪造人脸图像生成机制示意图

Fig. 4 Schematic of GAN-based forged facial image generation

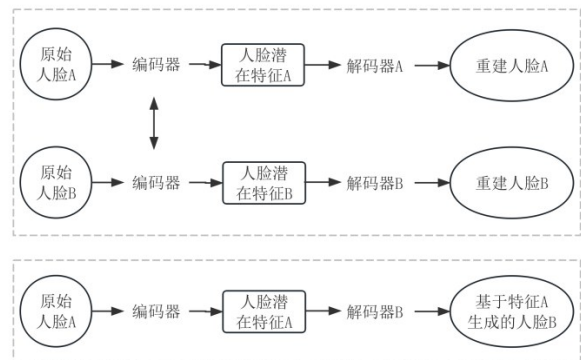


图5 基于自编码器的伪造人脸图像生成机制示意图

Fig. 5 Schematic of autoencoder-based forged facial image generation

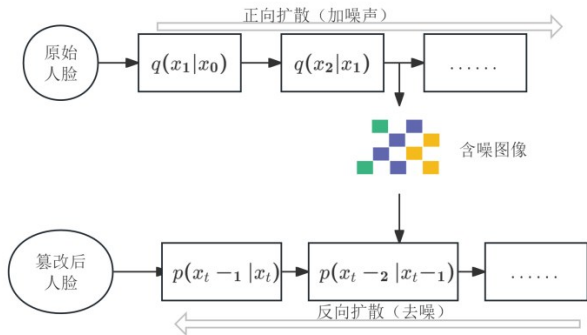


图6 基于扩散模型的伪造人脸图像生成机制示意图

Fig. 6 Schematic of diffusion-based forged facial image generation

唇部动作等特征传递至目标个体,同时完整保留目标主体的身份属性;面部编辑则针对特定面部属性进行局部修改,如调整年龄、发型或表情等,不改变个体的核心身份信息;面部替换通过特征替换实现身份重构,以源个体的面部特征替代目标个体特征,进而改变输出结果中感知到的身份归属;面部合成则基于生成模型从零构建逼真人脸图像,生成的内容不对应任何真实存在的个体。

#### 1.3.4 后处理

基于68或81个关键点的面部遮罩已从早期方法的边缘模糊,发展为自适应高斯模糊,并结合泊松融合,使伪造区域与背景梯度对齐。伪造人像通过模拟自然阴影,使传统边缘检测失效,显著增加伪影检测难度。

深度伪造技术在技术与功能层面持续演进。图像鉴定聚焦功能逻辑,核心是判断图像是否涉及换脸、表情操控等伪造行为,还是真实内容。技术上需识别其来源、方法与工具,不仅用于确认伪造,还需阐明生成方式与动因,以满足法庭证据的可采性

要求。

## 2 深度伪造人脸检测

深度伪造技术的演进路径,也是其作为犯罪载体的发展过程。

### 2.1 研究热点

本研究以“(deepfake OR "face forgery" OR "face manipulation" OR "face swap" OR "facial deepfake") AND (face OR facial)”为Topic,在Web of Science核心合集(Core Collection)检索近五年文献,共获取845篇文档;经筛选排除与深度伪造人脸不相关的内容及综述类文章后,剩余737篇有效文献。保留文献的题目与摘要信息并导入VOSViewer,以Binary Counting方式计算关键词频数,选取频数大于20的关键词共196个,展示前60%(共118个);将标签最大长度设置为10、线长设置为1000,最终生成网络可视化图谱。

图7展示了2022至2026年深度伪造人脸图像检测领域文献的发表趋势。2022年与2023年发文量分别为101篇和108篇,呈现稳步增长态势,标志着该领域处于技术积累与认知拓展阶段——随着深度伪造技术引发的安全风险逐步显现,研究重点集中于检测方法的构建,致力于建立技术框架与理论体系。2024年发文量跃升至265篇,2025年维持在260篇的高位水平,进入显著增长期。这一爆发式增长源于双重驱动机制:其一,学术界对该议题的关注持续升温,高精度检测模型与大规模标注数据集的相继推出为技术迭代提供了坚实支撑;其二,现实中深度伪造引发的虚假信息传播、身份冒用等安全事件频发,对社会信任体系与个体权益构成严峻挑

表2 深度伪造人脸图像的分类

Table 2 Classification of Deepfake Facial Images

| 类别   | 身份变化   | 技术实现                             | 潜在风险               |
|------|--------|----------------------------------|--------------------|
| 面部重演 | 保留目标身份 | 迁移源面部表情与头部姿态至目标,复刻动作             | 伪造公众人物言论或行为        |
| 面部编辑 | 保留目标身份 | 基于GAN或编解码器模型,修改年龄、发型、表情等属性,不改变身份 | 篡改数字证据、生成误导性生物特征数据 |
| 面部替换 | 替换为源身份 | 通过面部对齐、特征提取与图像融合,迁移源面部结构与纹理至目标   | 身份欺诈、深度伪造色情内容、冒充攻击 |
| 面部合成 | 生成全新身份 | 基于生成模型或3D建模,从噪声生成逼真人脸            | 传播虚假信息、伪造网络身份实施诈骗  |

战,催生了迫切的实际防控需求。上述因素共同推动研究范式从技术可行性验证向实际应用阶段适配演进。2026年截至统计时点(2026年1月14日)仅收录3篇文献,因尚处年度初期,反映该年度研究工作正处于起步布局阶段。

从图8的网络聚类可见,红色与橙色节点构成以“deepfake”为核心的技术聚类,涵盖伪造类型(如“faceswap”“fake video”“fake image”)和技术路径(如“deep neural network”“feature extraction”“convolution”),表明计算机视觉领域已围绕深度伪造图像构建起系统化的研究框架,并持续聚焦于特征提取、卷积模型优化等关键技术方向,展现出较强的技术积累与演进能力,具备支撑实际鉴伪应用的基础条件。同时,该聚类关联“precision”“accuracy”“f1 score”等评估指标,进一步凸显研究对检测性能精细化评估的高度重视,反映出该领域正由方法探索向性能验证深化,是其研究成熟度提升的重要体现。

绿色节点形成以“detection”为中心的实践聚类,延伸至“benchmark”“generalization”“artifact”等关键术语,不仅体现了计算机视觉领域当前的研究重心,也恰好契合司法实践中对证据鉴伪在可验证性、稳定性与可解释性方面的需求。然而值得注意的是,现有研究仍主要集中于技术性能的内部优化,尚未主动构建计算机视觉方法与法庭证据标准之间的系统性对接机制。尽管学界对“artifact”(伪造痕迹)、“frequency”(频率特征)等底层信号的分析已在事实上触及证据真实性的判断依据,而司法体系对证据可靠性与程序公正的要求亦亟需技术手段的支持,但目前技术供给与司法需求之间仍呈现并行发展、缺乏深度融合的局面,二者在术语体系、验证逻辑与应用阶段上的衔接尚待加强。

综上所述,从深度伪造人脸图像作为司法证据的视角出发,该网络图谱不仅反映出计算机视觉领域在深度伪造人脸检测研究方面的技术成熟度与发展深度,同时也揭示了当前技术研究与司法实践中证据鉴伪需求之间存在明显脱节的现状。未来的研究应主动构建技术与司法之间的跨学科衔接机制,推动计算机视觉成果在刑事侦查中更高效、精准地应用于证据真实性检验,实现技术能力与司法需求之间的双向对接。

表3 按技术原理对深度伪造人脸检测方法的分类

Table 3 Classification of Deepfake Detection Methods by Technical Principle

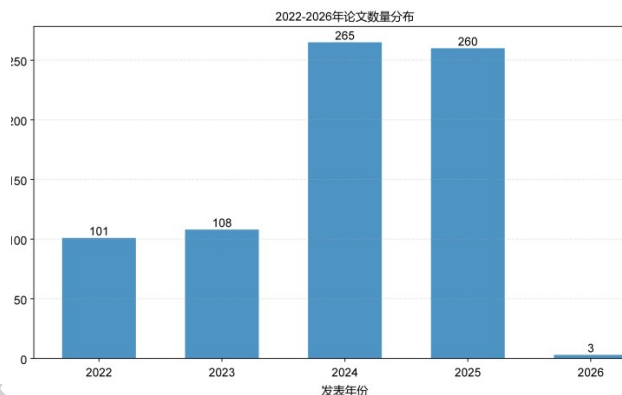


图7 深度伪造人脸图像检测文献近五年发表数量

Fig. 7 The number of published papers on deepfake face image detection in the past five years

## 2.2 检测方法分类

随着深度伪造人脸图像技术的演进,检测技术也随之不断发展(丁峰等,2024)。深度伪造人脸检测通常被视为二分类任务,即通过深度模型区分真实与伪造的人脸图像。为了便于与刑事司法需求对接,本文从两个维度对深度伪造人脸取证方法进行归类(姚文达等,2025)。其一是按技术原理划分,聚焦算法框架、特征提取和模型构建等底层逻辑,回答“如何判定图像为伪造”的核心问题;其二是按应用阶段划分,根据技术在伪造流程中的介入时点,将方法分为伪造前的源头预防、伪造中的实时监测与拦截、伪造后的溯源追责三个关键环节。

### 2.2.1 按技术原理分类

按技术原理划分,深度伪造人脸检测方法大致可分为三类:特征分析、模型能力和对抗策略。

特征分析是法庭场景下深度伪造人脸图像证据鉴伪的基础,通过时空分析、频率分析、多模态分析、生物信号分析、深度特征分析(杨少聪等,2022)等方法,识别图像的物理或生物特性,为刑事侦查提供依据。时空分析检测面部运动或表情不一致,排查伪造痕迹;频率分析揭示频率域中的压缩失真与边缘模糊,甄别虚假特征;多模态分析交叉验证图像、音频与文本,强化证据链判断;生物信号分析利用心率、瞳孔运动等难以伪造的生理特征,提供本质性鉴伪依据;深度特征分析借助深度学习发现伪造区域的细微纹理异常,辅助精准判定。



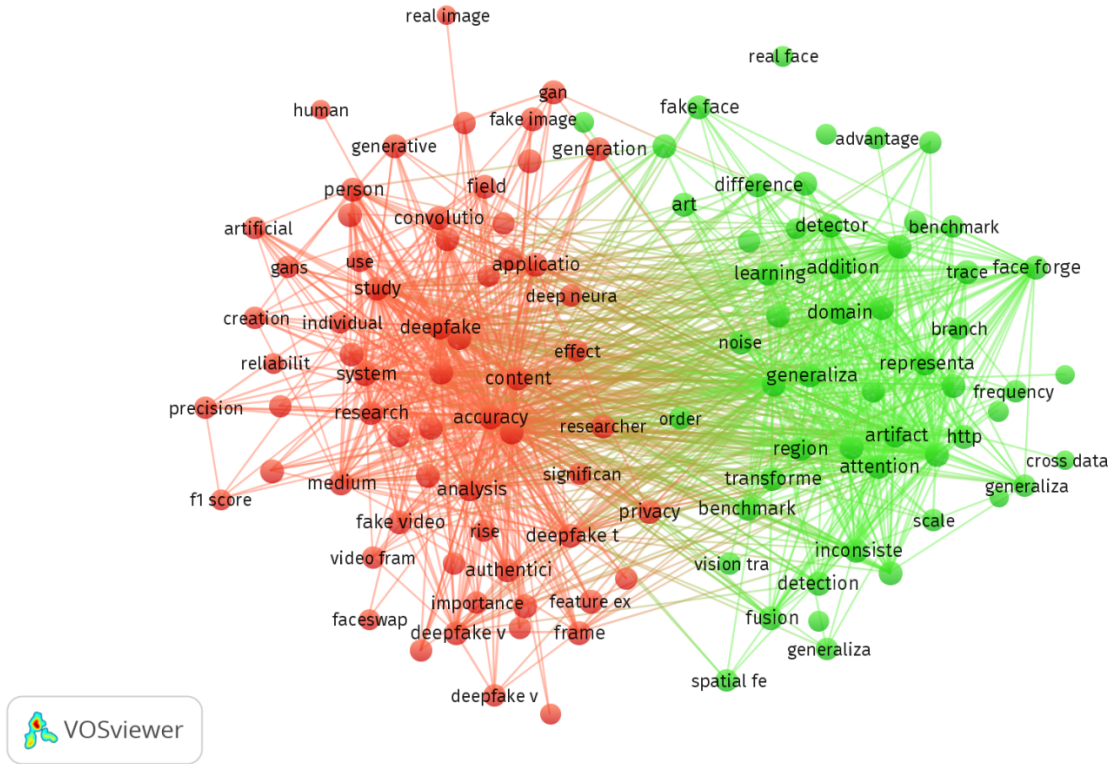


图8 使用VOSviewer的关键词共现可视化

Fig. 8 Keyword co-occurrence visualization of eligible studies using VOSviewer

| 技术原理 | 方法    | 解决问题                         | 主流方法   | 技术优势                            | 局限性                             |
|------|-------|------------------------------|--|---------------------------------|---------------------------------|
|      | 时空法   | 视频中因人脸伪造引发的帧间逻辑矛盾            | 帧间运动一致性检测、光照与姿态的时序连贯性分析 (Shahzad 等, 2025 ; Xu 等, 2024 ; Guan 等, 2024 ) | 可捕捉人眼难以察觉的伪造痕迹                  | 对视频压缩和帧率变化敏感,难以有效检测无明显帧间异常的伪造视频 |
|      | 频率法   | 挖掘伪造人脸图像技术在频域留下的信息           | 傅里叶变换频谱分析、小波变换异常检测 (Jin 等, 2024)                                       | 对图像压缩和噪声扰动具有鲁棒性                 | 对扩散模型等新型伪造技术响应较弱                |
| 特征分析 | 多模态法  | 整合多维度人脸图像证据信息,解决单一模态鉴伪不全面的问题 | 音视频同步性验证、基于CLIP等大模型的视觉-文本语义一致性分析 (周成祖等, 2025 ; Edwin Joel 等, 2025)     | 鲁棒性强,能识别复杂伪造人脸图像,降低误判           | 依赖高质量多模态数据,模型结构复杂且成本较高          |
|      | 生物信号法 | 鉴定人脸证据的生物真实性                 | 人脸生理信号(心跳、呼吸)检测与眨眼频率一致性分析(Ni 等, 2024 ; Jin 等, 2021)                    | 泛化能力强,生物特征难伪造                   | 对视频质量、光照和头部姿态要求高,实用性受限          |
|      | 深度特征法 | 提取人脸图像证据中的深层伪造痕迹             | CNN/Transformer特征提取,基于注意力机制的细微伪影检测(Liang 等, 2023 ; Wang 等, 2024 )      | 特征表达能力强,检测精度高,可适配多种伪造类型及高质量证据鉴伪 | 依赖训练数据,易发生过拟合,跨数据集泛化能力弱         |



表 3 续表

| 技术原理 | 方法     | 解决问题                                 | 主流方法   | 技术优势                             | 局限性                                  |
|------|--------|--------------------------------------|--|----------------------------------|--------------------------------------|
| 模型能力 | 自监督学习  | 减少对标注数据的依赖,提升模型对多样伪造证据的适配性           | 基于对比学习与掩码重建的特征学习(Qiao 等, 2024; Shao 等, 2025)                             | 无需大规模标注数据,泛化能力强,适配司法场景数据稀缺现状     | 模型训练难度高,部分场景检测精度低于监督学习,性能验证需更多司法案例支持 |
|      | 多分类学习  | 区分深度伪造人脸图像证据的伪造类型,为法庭提供溯源依据,辅助案件事实认定 | 基于类别标签的伪造技术分类模型(如区分 GAN 类与扩散模型类伪造)(Deng 等, 2025; Bunluesakdikul 等, 2025) | 分类结果明确,可为案件定性提供技术支持,证据指向性强       | 对新型伪造类型的识别能力弱,需持续更新训练数据中的伪造类型标签      |
|      | 元学习    | 快速适配新型伪造证据鉴伪需求,解决司法实践中伪造技术迭代快的问题     | 采用少样本/零样本检测模型与元知识迁移学习(Huang 等, 2024; Liu 等, 2024)                        | 样本需求小,可快速响应新型伪造手段,满足司法场景时效性要求    | 模型稳定性有待提升,小样本训练易致误判,需结合其他方法交叉验证      |
|      | 异常检测   | 识别偏离正常模式的伪造证据,降低误用异常证据为有效证据的风险       | 基于正常人脸数据的分布建模与重构误差异常判定(Rosca 等, 2025; Wang 等, 2024)                      | 无需标注伪造样本,可适配未知伪造类型,应用阶段灵活        | 易受数据分布差异影响,异常阈值设定主观,可能导致合理证据误判       |
|      | 鲁棒性优化  | 提升模型在复杂场景下的鉴伪可靠性                     | 对抗训练、跨域自适应优化与噪声鲁棒性增强设计(Krasilnikov 等, 2025)                              | 模型抗干扰能力强,输出稳定,证据可信度高             | 优化过程复杂,可能牺牲检测精度,鲁棒性验证需覆盖多样化的司法证据场景   |
|      | 可解释性增强 | 破解深度模型“黑箱”问题,满足法庭证据可采性要求             | 梯度可视化、注意力权重分析、特征归因方法(Peng 等, 2022; Yang 等, 2024)                         | 呈现鉴伪依据,降低法庭解释难度                  | 可解释性与检测精度存在权衡,部分解释方法结果主观性强           |
| 对抗策略 | 伪造定位   | 定位证据中的伪造区域,为法庭明确证据中真实部分与伪造部分的边界      | 伪造区域分割模型、像素级异常检测(Waseem 等, 2023; Zeng 等, 2025)                           | 能区分证据中“真”“伪”成分,服务于鉴定人            | 定位精度受图像质量影响大,边缘区域定位准确性不足,需结合其他方法验证   |
|      | 对抗防御   | 抵御伪造技术对鉴伪模型的攻击,确保鉴伪结果稳定可靠            | 对抗样本训练、防御性蒸馏、模型集成策略(Lei 等, 2025; Dong 等, 2023)                           | 提升模型抗攻击能力,减少恶意伪造技术对鉴伪结果的干扰       | 增加模型复杂度与计算成本,部分防御策略可能影响正常的检测效率       |
|      | 源模型溯源  | 追溯伪造证据的生成模型,为法庭提供伪造行为溯源线索,辅助追查相关责任人  | 生成模型指纹提取、模型输出特征匹配溯源(Sun 等, 2025)   | 能为案件侦破提供技术线索,强化伪造行为与责任人的关联性      | 溯源精度受模型迭代、参数微调影响大,跨平台溯源难度高,技术成熟度不足   |
|      | 联邦学习   | 在保护数据隐私的前提下开展多机构鉴伪模型协作               | 面向跨办案机构的联邦训练框架、模型参数加密传输(Zhang 等, 2025)                                   | 保护涉案数据隐私,实现司法资源共享,提升模型泛化         | 系统部署与协调成本高,模型训练通信开销大,需顶层协调           |
|      | 知识蒸馏   | 构建轻量化鉴伪模型,满足司法一线快速鉴伪需求,降低设备部署与使用门槛   | 教师-学生模型蒸馏、关键特征迁移学习(Wang 等, 2025; Khan 等, 2024)                           | 模型体积小,推理速度快,适配司法一线现场终端鉴伪场景,部署成本低 | 蒸馏过程可能损失部分检测精度,复杂伪造证据鉴伪能力弱于原始大模型     |

模型能力指通过自监督学习、多类别学习、元学习、异常检测、鲁棒性优化、可解释性增强和伪造定位,提升证据取证的可靠性与刑事侦查适配性。其中,自监督学习利用图像修复、时间预测等无标签数

据模式,降低对标注数据的依赖;多类别学习区分真实内容与特定伪造类型,实现细粒度分类;元学习基于少量样本快速适应新型伪造,应对新形态虚假证据的挑战;异常检测通过建模正常数据分布标记异常,辅助识别篡改内容;鲁棒性优化增强对噪声、压缩和篡改的抵抗能力,保障复杂场景下的稳定性;可解释性增强通过可视化可疑区域提高结果可信度,满足法庭对过程可追溯的要求;伪造定位精准识别被篡改的面部或合成背景,明确证据篡改范围。

对抗策略通过反制与溯源技术,既实现抗规避检测,又可追踪证据来源,为刑事追责提供支撑。具体包括对抗防御、源模型溯源、联邦学习与知识蒸馏。对抗防御抵御微小扰动等规避攻击,提升模型韧性,确保鉴伪结果可靠;源模型溯源通过模型特有噪声等独特痕迹识别生成模型,助力追查伪造源头;

联邦学习支持隐私保护下的分布式训练,在保障案件数据安全的同时提升泛化性;知识蒸馏将复杂模型知识迁移至轻量模型,提升移动端与边缘设备部署效率,适配现场取证与庭审快速鉴伪需求。

文献统计显示,特征分析是法庭证据鉴伪的基础。其中频率特征占比最高(31.4%),其次为多模态特征(13.1%)、深度特征(4.1%)、时空特征(3.2%)、生物特征(3.1%),这些特征是鉴伪判断的核心依据。

模型能力研究占文献的53.1%,是刑事侦查适配性优化的核心。鲁棒性优化占39.4%,伪造定位占18.0%,异常检测占13.6%、自监督学习占9.8%、可解释性占9.0%、多类别学习占8.9%、元学习占7.9%。以上覆盖模型全周期,契合刑事侦查鉴伪的严苛要求。

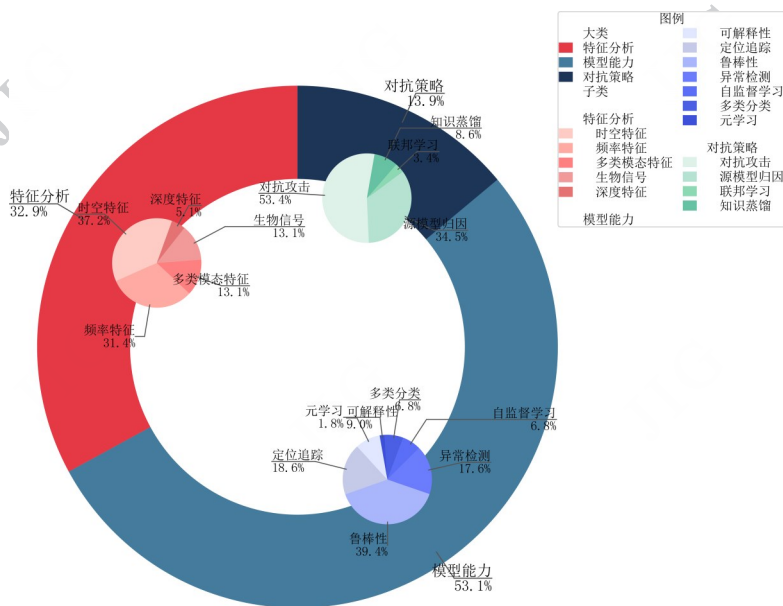


图9 按技术原理对深度伪造人脸检测方法的分类

Fig. 9 Classification of deepfake detection methods based on technical principles

对抗策略研究占13.9%,聚焦攻防对抗与证据溯源。对抗防御占53.4%、模型溯源占34.5%,辅以知识蒸馏占8.9%、联邦学习占4.9%,这恰满足刑事侦查实践中多样化的鉴伪与追责需求。

### 2.2.2 按应用阶段分类

深度伪造人脸的预防与控制分为事前预防、事中监测和事后追溯,实现从源头阻断到案件响应的端到端刑事司法过程全覆盖,适配刑事侦查与法庭证据。

事前预防依托指纹水印与身份关联,可服务于

侦查源头防控。指纹水印在图像创建时嵌入不可见标识符,通过验证其完整性可快速排查篡改。身份关联将伪造内容与伪造人绑定,减少伪造图像进入刑事诉讼程序的风险,降低取证干扰。

事中监测采用特征异常检测与对抗防御,支撑侦查取证。特征异常检测利用时空、频率和多模态方法分析传输的人脸图像,识别并拦截伪造内容,为固定真实证据提供支持;对抗防御应对小噪声干扰和对抗样本,确保模型在规避攻击下仍能稳定输出,满足证据有效性要求。

表 4 按应用阶段划分的深度伪造人脸检测方法分类  
Table 4 Classification of Deepfake Detection Methods by Application Stage

| 检测阶段 | 方法    | 解决问题  | 主流方法  | 技术优势   | 局限性   |
|------|-------|---|---|--|---|
| 事前预防 | 指纹水印  | 在原始人脸图像中植入不可见标识,从源头防范伪造人脸图像作为犯罪载体                     | 嵌入鲁棒性数字水印(如DWT域水印)、区块链存证水印(Ge等, 2025; Yang等, 2023)          | 源头把控人脸图像真实性,水印可作为直接溯源依据,说服力强;不影响图像视觉效果               | 依赖原始数据采集环节介入,无法覆盖已流通的无水印图像;水印易受图像压缩、裁剪等操作破坏,影响验证有效性       |
|      | 身份关联法 | 建立人脸图像与真实身份的唯一关联,防范伪造他人身份                             | 人脸生物特征绑定(如虹膜/指纹与人脸多模态关联)、身份信息加密嵌入、权威数据库身份校验(Zhang等, 2025)   | 身份关联性强,能快速排查伪造身份的的证据;验证逻辑直观,易被法庭理解采纳;适配身份类案件证据鉴伪需求   | 依赖权威身份数据库支撑,跨区域/跨机构数据库协同难度大;对无身份信息的匿名图像无效,应用阶段受限          |
| 事中监控 | 异常检测  | 在深度伪造人脸图像实时识别异常篡改/伪造行为,避免伪造证据进入司法程序,降低错案风险            | 实时帧间一致性监控、图像哈希值动态校验、异常特征实时预警模型(Stamnas等, 2025)              | 实时性强;无需标注伪造样本,适配未知伪造类型;操作便捷,适配司法一线流转场景               | 易受环境干扰导致误预警;对轻微伪造行为敏感度不足;对实时监控算法和算力要求高                    |
|      | 对抗防御  | 抵御伪造技术对证据鉴伪系统的攻击,确保鉴伪过程稳定可靠,保障鉴伪结果的证据效力               | 对抗样本训练增强鲁棒性、防御性蒸馏模型、多模型集成防御策略(Galdi等, 2024)                 | 抗攻击能力强,能应对恶意伪造技术的针对性干扰;鉴伪结果稳定性高                      | 模型训练复杂度高,计算成本大;防御范围有限,难以覆盖所有新型攻击手段;可能牺牲鉴伪效率               |
| 事后追溯 | 人工判定  | 对疑似伪造的证据进行最终人工鉴定,为法庭提供权威鉴伪意见                          | 法医图像专家人工审核、多专家交叉验证、结合司法案例经验判定(Bharati等, 2026)               | 和其他如指纹、枪弹痕迹等证据鉴定流程一致;可结合案件背景综合判定,灵活性强                | 专家需结合检测结果与可视化输出,进行图结构分析和案件事实判断,依赖算法且主观性强                  |
|      | 伪造定位  | 精准定位证据中伪造区域的范围与边界,明确证据“真”“伪”成分                        | 基于深度网络的伪造区域分割、像素级异常特征定位、伪造边缘轮廓提取(Zhang等, 2025; Peng等, 2023) | 鉴伪结果精细化;定位依据客观,可通过可视化呈现增强说服力;更符合法院、检察院、公安、案件相关人的常识   | 定位精度受图像质量影响大,模糊图像边缘定位准确性与精确性不足;对轻微篡改的定位敏感度低;需结合其他鉴伪方法交叉验证 |
|      | 公平性优化 | 解决不同人群(如不同肤色、性别)人脸伪造鉴伪的偏差问题,保障鉴伪结果的公平性,避免因数据集偏差影响司法公正 | 公平性约束损失函数优化、多人群均衡数据集训练、跨人群迁移学习(Xu等, 2024; Ju等, 2024)        | 鉴伪结果公平性强,避免对特定人群的歧视性误判;符合司法公正的核心要求;提升模型在多样化人群案件中的适配性 | 公平性与检测精度存在权衡,可能牺牲部分单一人群的鉴伪性能;多人群均衡数据集构建难度大,数据获取成本高        |
|      | 模型交互法 | 通过多模型协同交互验证鉴伪结果,解决单一模型鉴伪的局限性,提升复杂伪造证据鉴伪的可靠性与说服力       | 多模型投票决策、专家系统与深度学习模型交互、跨模态模型协同验证(Lin等, 2024; Sun等, 2025)     | 鉴伪结果可靠性高,误判率低;能应对复杂多样的伪造场景;可通过多模型结果互补增强法庭说服力         | 系统架构复杂,部署与维护成本高;多模型协同存在时延,适配实时鉴伪场景受限;鉴定人的水平参差增加不确定性       |

事后追溯包含四项核心任务,适配法庭证据审查与侦查需求。一是人为决策支持,借助可解释AI输出伪造区域标签和热图,辅助人工审查,确保证据

符合质证标准;二是伪造定位,精确标记图像中的篡改区域,为认定篡改范围提供直观依据;三是模型迭代,基于已破案件通过元学习、联邦学习和知识蒸馏



更新模型,提升对新型伪造的识别能力;四是公平性验证,评估结果以防止因年龄、性别等属性导致偏差。

图 10 展示伪造检测系统的双层环结构阶段划分及文献方法分布,各阶段占比契合刑事司法资源配置逻辑。

事前阶段占 19.7%,聚焦主动风险控制。其中指纹水印占 38.0%、身份关联占 62.0%,通过水印验证或身份认证实现早期拦截,减少后续取证障碍。

事中阶段占 42.2%,构成实时检测核心(李启运等,2019)。异常检测占 54.7%、抗攻击防御占 45.3% 共同保障侦查中图像检测的稳健性与适应性。

事后阶段占 38.1%,支撑法庭回顾性分析。决策支持占 23.9%、伪造定位占 26.6%、公平性评估占 42.4%、互动机制占 6.7%,可助力物证鉴定人鉴定伪造图像。

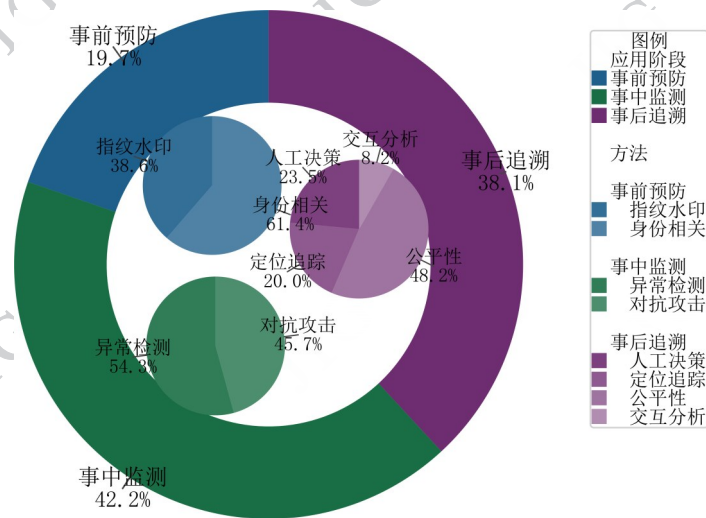


图 10 按应用阶段对深度伪造人脸检测方法的分类

Fig. 10 Classification of deepfake detection methods by application scenario

综上所述,将深度伪造人脸图像检测方法按技术原理和应用阶段分类,为新型深度伪造犯罪案件的法庭证据鉴定提供系统且精准的支持。技术原理上,特征分析、模型能力与对抗策略三类方法各有适用场景与优劣,明确了鉴伪依据,并为图像、视频、多模态等不同类型证据的技术选择提供指导,解决技术上检测伪造人脸图像的核心问题。应用阶段上,从事前防范、事中控制到事后追溯的全流程划分,契合社会治理与公共安全需求,既能减少伪造证据进入司法程序,保障取证真实,也可通过伪造定位、公平性验证等手段,帮助法庭界定证据真伪,提升质证说服力。以上分类体系有助于司法人员从技术上开展对该新型犯罪载体的研究。

#### 4 检测与鉴定

深度伪造人脸图像作为一种新型犯罪载体,其生成机制基于生成器与判别器之间的对抗博弈模

型,导致伪造样本在特征空间中呈现出复杂的分布模式与高度隐蔽的伪影特征。此类技术的持续迭代,不仅表明传统判别器已难以精准识别特征空间中的伪造痕迹,更说明对该类图像的鉴定无法仅依赖专家对表层视觉特征的人工判读完成。如图 11 所示,唯有依托专业的检测模型,通过对高维特征空间的深度挖掘与系统性特征建模,才能实现科学、可靠的鉴定结论。

在“鉴定依赖检测”这一前提下,需进一步探讨当前检测技术为何尚难达到司法鉴定标准。若暂不考虑深度伪造图像的鉴定标准与技术规范问题,转而从可操作性角度分析,其根本原因在于:指纹、枪弹痕迹(Song 等,2020)、声纹(Morrison,2011)、足迹(Ma 等,2024)、DNA(Agudo 等,2024)以及人脸识别(Macarulla Rodríguez 等,2024)等物证均已纳入贝叶斯推理框架,以似然比作为标准化的证据表达形式(Van Lierop 等,2024)。该框架基于控方假设与辩方假设进行推断。



$$LR = \frac{f(\text{score}|H_p, I)}{f(\text{score}|H_d, I)} \quad (\text{公式 1})$$

例如,在犯罪现场提取一枚指纹并怀疑其属于犯罪嫌疑人X时,存在两种对立假设:控方假设为“该指纹来源于嫌疑人X”,辩方假设为“该指纹来源于某一随机个体”。通过提取该指纹的特征,可计算其与嫌疑人指纹之间的相似度。依托已有数据库,可预先建立两类概率密度函数:一是两枚指纹同源(来自同一人)时相似度的分布函数,二是异源(来自不同人)时相似度的分布函数。将实际观测到的相似度分别代入这两个函数,所得同源概率密度与异源概率密度之比,即为该指纹证据的似然比,用以量化其证据强度。如公式(1)所示其中LR(likelihood ratio)代表似然比,score是现场提取的指纹检材与嫌疑人指纹样本之间的相似度,Hp(Prosecution Hypothesis)代表控方假设。Hd(Defence Hypothesis)代表辩方假设。如 $LR=10^3$ ,意味着该指纹相似度的观测结果,在“指纹属于嫌疑人X”的前提下发生的可能性,是“指纹属于随机个体”前提下的1000倍,属于强证据,能够有力支撑控方的主张。目前,在鉴定科学领域,所有证据普遍采用似然比作为统一的评估指标。

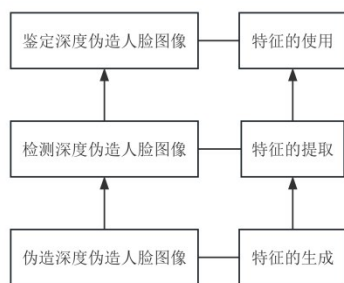


图 11 深度伪造人脸图像的鉴定前提

Fig. 11 Premises for the authentication of deepfake facial images

深度伪造人脸检测模型的典型输出为一组二分类概率值 $(p, 1-p)$ ,其中 $p$ 代表输入图像为深度伪造内容的预测概率, $1-p$ 则对应图像为真实拍摄内容的预测概率。在工程化应用中,模型常以0.5为判定阈值,直接输出“伪造”或“真实”的离散结论。需明确的是,此类基于经验风险最小化的概率输出与阈值判定方式,并未直接纳入贝叶斯推理的完整分析框架——其本质是模型基于训练数据集的特

征学习与概率拟合,未结合司法场景下的先验信息完成后验概率的推演。

类比指纹鉴定的法庭科学统计分析范式,深度伪造图像作为新型电子证据的鉴定工作,需构建一组相互对立的检验假设,以支撑证据证明力的量化评估。针对不同案件的待证事实,核心假设可分为两类:

#### (1)基础属性假设

控方假设:送检的涉案人脸图像为深度伪造生成;

辩方假设:送检的涉案人脸图像为真实拍摄所得。

#### (2)关联属性假设

控方假设:送检的、与嫌疑人X面部特征高度相似的人脸图像为深度伪造生成;

辩方假设:送检的、与嫌疑人X面部特征高度相似的人脸图像为X本人的真实人脸图像。

上述两类假设的量化分析,均需依托深度伪造人脸检测技术提取图像的鉴别特征或者分数,并基于特征或者分数构建概率密度函数,进而通过似然比等指标实现证据强度的客观量化。

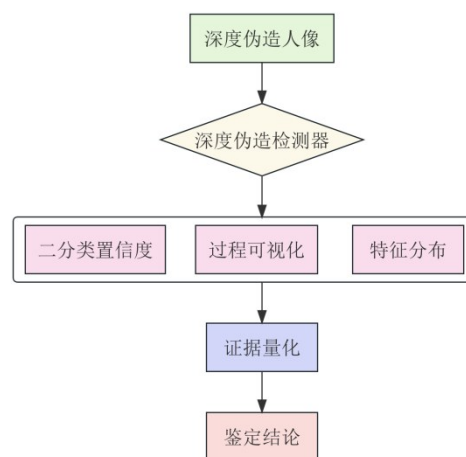


图 12 检测与鉴定的关系

Fig. 12 Detection and Identification Relationship

上述计算过程可视为深度伪造图像证据的量化分析环节。鉴于司法实务人员对技术证据在理解与采信方面的需求,可通过检测流程可视化、特征分布可视化等方式,将抽象的量化结果转化为直观的视觉呈现,如图12所示。

尽管深度伪造人脸检测技术能够生成伪造区域  
© 中国图象图形学报版权所有

热图、异常特征标注等可视化输出,其实质仍局限于技术特征的识别与量化,无法自主判断这些特征所承载的法律意义。例如,在诈骗案件中,检测到的人脸篡改痕迹是否指向被害人身份信息的伪造;在诬告案件中,图像的篡改区域是否影响关键时间戳的真实性——此类问题均需鉴定专家结合刑事侦查知识与具体案情进行专业解读。此外,涉案图像元数据的合法性、检测结果与待证事实之间的关联性,以及检测模型是否存在年龄、性别等属性上的系统性偏见,亦需由专家开展综合评估。唯有经过上述审查程序,方可将技术层面的量化结果转化为符合刑事侦查与司法规范要求的鉴定意见。

需要强调的是,人类对深度伪造内容的主观辨识能力本就有限,且随着伪造技术的不断迭代而持续下降,这进一步凸显了技术检测与专家审查相结合的重要性。实现从技术特征检测向司法证据审查的全流程转化。

## 5 结 语

本文所探讨的深度伪造人脸图像生成技术,已从生成对抗网络等早期技术演进至扩散模型与多模态驱动阶段,并初步构建起完整技术流程。该流程可实现篡改检测、工具追踪与来源溯源,充分满足刑侦证据的可追溯性需求,且配套检测技术兼具高精度与可解释性,能生成直观篡改标记并定位伪造范围。同时,文中提出的端到端框架覆盖事前水印管控、事中异常检测与事后溯源全环节,有效提升了鉴定工作的科学性。最后,本文基于鉴定领域对证据形式所要求的贝叶斯框架,提出基于检测技术的深度伪造图像鉴定可行性方法。综上,本文通过相关技术与框架的研究,推动了数字取证与法庭证据鉴定的深度融合,为应对深度伪造技术对刑事证据构成的挑战提供方法参考和实践启示。

## 参考文献(Reference)

Agudo M M, Fantinato C, Roseth A, Aanes H, Gill P, Fonnelløp A E, et al. 2024. A comparison of likelihood ratios calculated from surface DNA mixtures using MPS and CE technologies [J]. *Forensic Science International: Genetics*, 73: 103111 [DOI: 10.1016/j.fsi-gen.2024.103111]

Bharati N, Wong P, Mostéfaoui S K, Kbaier D, Collie J. 2025. Explainable deepfake detection: a multi-model framework with human-interpretable rationales for legal investigation purposes [J]. *Machine Learning with Applications*, 23: 100819 [DOI: 10.1016/j.mlwa.2025.100819]

Bunluesakdikul P, Mahanan W, Sungunnasil P, Sangamuang S. 2025. Deepfake video detection: a novel approach via NLP-based classification [J]. *International Journal of Computational Intelligence and Applications*, 24 (2) : 2550001 [DOI: 10.1142/S1469026825500014]

Deng L, Wu B, Wang J. 2025. A multi-label classification method combined with texture enhancement for deepfake face detection [J]. *Multimedia Systems*, 31 (6) : 410 [DOI: 10.1007/s00530-025-01996-y]

Dong J, Wang Y, Lai J, Xie X. 2023. Restricted black-box adversarial attack against DeepFake face swapping [J]. *IEEE Transactions on Information Forensics and Security*, 18: 2596 - 2608 [DOI: 10.1109/TIFS.2023.3266702]

Galdi C, Panariello M, Todisco M, Evans N. 2024. 2D-malafide: adversarial attacks against face deepfake detection systems [C]//2024 International Conference of the Biometrics Special Interest Group (BIOSIG). Darmstadt, Germany: IEEE: 1 - 7 [2026-02-05] [DOI: 10.1109/BIOSIG61931.2024.10786754]

Ge J W, Cao J X, Zhao Z X, Liu B. 2025. FSD-GAN: generative adversarial training for face swap detection via the latent noise fingerprint [J]. *Journal of Computer Science and Technology*, 40(2) : 397 - 412 [DOI: 10.1007/s11390-024-3337-8]

Guan W, Wang W, Peng B, Dong J, Tan T. 2025. ST-SBV: spatial-temporal self-blended videos for deepfake detection [M]//Lin Z, Cheng M M, He R, et al. *Pattern recognition and computer vision*. Singapore: Springer Nature Singapore: 274 - 288 [2026-02-05] [DOI: 10.1007/978-981-97-8620-6\_19]

Huang D, Zhang Y. 2024. Learning meta model for strong generalization deepfake detection [C]//Proceedings of the International Joint Conference on Neural Networks. Yokohama, Japan: IEEE: 1 - 8 [2026-02-05] [DOI: 10.1109/IJCNN60899.2024.10651482]

Joel J E, Kumar P K, Silar M M, Rohan J. 2025. Flexible multi-modality deepfake identification with dynamic learning and cross-modality inconsistency [C]//2025 3rd International Conference on Advancements in Electrical, Electronics, Communication, Computing and Automation, ICAECA 2025. Coimbatore, India: IEEE: 1 - 6 [2026-02-05] [DOI: 10.1109/ICAECA63854.2025.11012276]

Jin X, Wu N, Jiang Q, Kou Y, Duan H, Wang P, et al. 2024. A dual descriptor combined with frequency domain reconstruction learning for face forgery detection in deepfake videos [J]. *Forensic Science International: Digital Investigation*, 49: 301747 [DOI: 10.1016/j.fsidi.2024.301747]

Jin X, Ye D, Chen C. 2021. Countering spoof: towards detecting deep-

- fake with multidimensional biological signals[J]. *Security and Communication Networks*, 2021:1 - 8 [DOI: 10.1155/2021/6626974]
- Ju Y, Hu S, Jia S, Chen G H, Lyu S. 2024. Improving fairness in deepfake detection[C]//*Proceedings - 2024 IEEE Winter Conference on Applications of Computer Vision, WACV 2024, Waikoloa, HI, USA: IEEE: 4643 - 4653 [2026-02-05] [DOI: 10.1109/WACV57701.2024.00459]*
- Khan S S, Hossain R, Bishal S H, Khan R. 2024. Smartphone-based deepfake detection through transfer learning and lightweight knowledge distillation technique[C]//*2024 27th International Conference on Computer and Information Technology, ICCIT 2024 - Proceedings. Cox's Bazar, Bangladesh: IEEE: 675 - 680 [2026-02-05] [DOI: 10.1109/ICCIT64611.2024.11022522]*
- Khormali A, Yuan J S. 2024. Self-supervised graph transformer for deepfake detection [J]. *IEEE Access*, 12: 58114 - 58127 [DOI: 10.1109/ACCESS.2024.3392512]
- Krasilnikov M, Nikitin M, Konushin A. 2025. VCF: a real-world video conference deepfake benchmark for face-swap detection and robustness evaluation[J]. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, XLVIII-2/W9-2025: 169 - 174 [DOI: 10.5194/isprs-archives-XLVIII-2-W9-2025-169-2025]*
- Lei S, Song J, Feng F, Yan Z, Wang A. 2025. Deepfake face detection and adversarial attack defense method based on multi-feature decision fusion [J]. *Applied Sciences*, 15 (12) : 6588 [DOI: 10.3390/app15126588]
- Liang B, Wang Z, Huang B, Zou Q, Wang Q, Liang J. 2023. Depth map guided triplet network for deepfake face detection [J]. *Neural Networks*, 159:34 - 42 [DOI: 10.1016/j.neunet.2022.11.031]
- Lin T T C. 2025. AI deepfake interaction, authentication and correction in taiwan: examining the roles of echo chamber and conspiracy mentality[C]//*Proceedings of the Annual Hawaii International Conference on System Sciences, Hawaii, USA: IEEE [2026-02-05] [DOI: 10.24251/HICSS.2025.285]*
- Liu X, Song P, Lu P, Wang Y. 2024. Meta-learning with relation embedding for few-shot deepfake detection [J]. *IEEE Access*, 12: 180135 - 180145 [DOI: 10.1109/ACCESS.2024.3499353]
- Ma X, Luo Y, Yu Y, Liu Y, Li S, Li H. 2024. Similarity quantification of bare footprint based on linear measurement and shape context contour combined score-based likelihood ratio evaluation [J]. *Forensic Science International*, 356: 111967 [DOI: 10.1016/j.forsciint.2024.111967]
- Macarulla Rodriguez A, Geradts Z, Worring M, Unzueta L. 2024. Improved likelihood ratios for face recognition in surveillance video by multimodal feature pairing [J]. *Forensic Science International: Synergy*, 8:100458 [DOI: 10.1016/j.fsisyn.2024.100458]
- Morrison G S. 2011. A comparison of procedures for the calculation of forensic likelihood ratios from acoustic - phonetic data: multivariate kernel density (MVKD) versus gaussian mixture model - universal background model (GMM - UBM) [J]. *Speech Communication*, 53(2):242 - 256 [DOI: 10.1016/j.specom.2010.09.005]
- Ni Y, Zeng W, Xia P, Yang G S, Tan R. 2024. A deepfake detection algorithm based on fourier transform of biological signal [J]. *Computers, Materials & Continua*, 79 (3) : 5295 - 5312 [DOI: 10.32604/cmc.2024.049911]
- Peng B, Lyu S, Wang W, Dong J. 2022. Counterfactual image enhancement for explanation of face swap deepfakes [M]//Yu S, Zhang Z, Yuen P C, et al. *Pattern recognition and computer vision. Cham: Springer Nature Switzerland: 492 - 508 [2026-02-05] [DOI: 10.1007/978-3-031-18910-4\_40]*
- Peng F, Zhang X, Long M. 2023. F2DLNet: a face forgery detection and localization network based on SSIM error maps [M]//Xu Y, Yan H, Teng H, et al. *Machine learning for cyber security. Cham: Springer Nature Switzerland: 355 - 369 [2026-02-05] [DOI: 10.1007/978-3-031-20099-1\_30]*
- Qiao T, Xie S, Chen Y, Retraint F, Luo X. 2024. Fully unsupervised deepfake video detection via enhanced contrastive learning [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(7):4654 - 4668 [DOI: 10.1109/TPAMI.2024.3356814]
- Rosca C M, Stancu A. 2025. AI anomaly-based deepfake detection using customized mahalanobis distance and head pose with facial landmarks [J]. *Applied Sciences*, 15 (17) : 9574 [DOI: 10.3390/app15179574]
- Shahzad S A, Hashmi A, Peng Y T, Tsao Y, Wang H M. 2025. AV-lip-sync+: leveraging AV-HuBERT to exploit multimodal inconsistency for deepfake detection of frontal face videos [J]. *IEEE Transactions on Human-Machine Systems*, 55(6):973 - 982 [DOI: 10.1109/THMS.2025.3618409]
- Shao C, Zhang F, Wang J, Zhang B. 2025. Sample based contrastive learning for DeepFake detection [M]//Yang Z, Sun G. *Proceedings of the 2nd international conference on networks, communications and intelligent computing (NCIC 2024). Singapore: Springer Nature Singapore: 993 - 1000 [2026-02-05] [DOI: 10.1007/978-981-96-5006-4\_86]*
- Song J, Chen Z, Vorburger T V, Soons J A. 2021. Evaluating likelihood ratio (LR) for firearm evidence identifications in forensic science based on the congruent matching cells (CMC) method [J]. *Forensic Science International*, 317: 110502 [DOI: 10.1016/j.forsciint.2020.110502]
- Stamnas S, Sanchez V. 2025. DiffFake: exposing deepfakes using differential anomaly detection [C]//*Proceedings - 2025 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops, WACVW 2025, Tucson, AZ, USA: IEEE: 647 - 657 [2026-02-05] [DOI: 10.1109/WACVW65960.2025.00079]*
- Sun C, Li W. 2025. A two-stage interaction approach for enhancing generalization of deepfake detection [J]. *Multimedia Systems*, 31(5) : 359 [DOI: 10.1007/s00530-025-01942-y]
- Sun Z, Chen S, Yao T, Yi R, Ding S, Ma L. 2025. Rethinking open-

- world DeepFake attribution with multi-perspective sensory learning [J]. *International Journal of Computer Vision*, 133(2): 628 - 651 [DOI: 10.1007/s11263-024-02184-7]
- Van Lierop S, Ramos D, Sjerps M, Ypma R. 2024. An overview of log likelihood ratio cost in forensic science - where is it used and what values can we expect? [J]. *Forensic Science International: Synergy*, 8: 100466 [DOI: 10.1016/j.fsisyn.2024.100466]
- Wang C, Meng L, Xia Z, Ren N, Ma B. 2025. Cross-domain deepfake detection based on latent domain knowledge distillation [J]. *IEEE Signal Processing Letters*, 32: 896 - 900 [DOI: 10.1109/LSP.2025.3540941]
- Wang H, Li S, He J, Qian Z, Zhang X, Fan S. 2024. Exploring depth information for detecting manipulated face videos [EB/OL]. [2026-02-05]. <http://dx.doi.org/10.48550/arXiv.2411.18572>
- Wang Y, Liao G. 2024. Deepfake video detection based on image source anomaly [C]//2024 IEEE 2nd International Conference on Image Processing and Computer Applications, ICIPCA 2024. Shenyang, China: IEEE: 397 - 401 [2026-02-05] [DOI: 10.1109/ICIPCA61593.2024.10709022]
- Waseem S, Abu-Bakar S A R S, Omar Z, Ahmed B A, Baloch S, Hafeezallah A. 2023. Multi-attention-based approach for deepfake face and expression swap detection and localization [J]. *EURASIP Journal on Image and Video Processing*, 2023(1): 14 [DOI: 10.1186/s13640-023-00614-z]
- Xu Y, Terhörst P, Pedersen M, Raja K. 2024. Analyzing fairness in deepfake detection with massively annotated databases [J]. *IEEE Transactions on Technology and Society*, 5(1): 93 - 106 [DOI: 10.1109/TTS.2024.3365421]
- Xu Y, Liang J, Sheng L, Zhang X Y. 2024. Learning spatiotemporal inconsistency via thumbnail layout for face deepfake detection [J]. *International Journal of Computer Vision*, 132(12): 5663 - 5680 [DOI: 10.1007/s11263-024-02054-2]
- Yang J, Sun Y, Mao M, Bai L, Zhang S, Wang F. 2023. Model-agnostic method: exposing deepfake using pixel-wise spatial and temporal fingerprints [J]. *IEEE Transactions on Big Data*, 9(6): 1496 - 1509 [DOI: 10.1109/TBDATA.2023.3284272]
- Yang Y, Joukovsky B, Oramas Mogrovejo J, Tuytelaars T, Deligiannis N. 2024. SNIPPET: a framework for subjective evaluation of visual explanations applied to DeepFake detection [J]. *ACM Transactions on Multimedia Computing Communications and Applications*, 20(8): 1 - 29 [DOI: 10.1145/3665248]
- Zeng S, Yi J, Tao J, He J, Lian Z, Liang S, et al. 2025. Adversarial training and gradient optimization for partially deepfake audio localization [C]//ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings. Hyderabad, India: IEEE: 1 - 5 [2026-02-05] [DOI: 10.1109/ICASSP49660.2025.10890470]
- Zhang B, Yin Q, Lu W, Luo X. 2025. Deepfake detection and localization using multi-view inconsistency measurement [J]. *IEEE Transactions on Dependable and Secure Computing*, 22(2): 1796 - 1809 [DOI: 10.1109/TDSC.2024.3472064]
- Zhang C. 2025. Federated learning-based cross-domain face forgery detection for consumer IOT systems [C]//Mohd Zain A B, Chen L. *Proceedings of SPIE - The International Society for Optical Engineering*. Kuala Lumpur, Malaysia: SPIE: 238 [2026-02-05] [DOI: 10.1117/12.3067515]
- Zhang Z, Zhang J, Zhou W, Zhou X, Guo Q, Zhang W, et al. 2025. FaceTracer: unveiling source identities from swapped face images and videos for fraud prevention [J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(12): 12021 - 12037 [DOI: 10.1109/TPAMI.2025.3601141]
- Zhou C Z, Zhang D, Chen Z, Zhao J Q, Zhang G B, Wei C, et al. 2025. A speaking face deepfake detection method, medium and device based on cross-modal learning [P]. China, CN120126223
- 周成祖, 章鼎, 陈忠, 赵建强, 张光斌, 魏超, 等. 2025. 基于跨模态学习的说话人脸深度伪造检测方法、介质及装置 [P]. 中国, CN120126223
- Yang S C, Wang J, Sun Y L, Tang J H. 2022. Multi-level features global consistency for human facial deepfake detection [J]. *Journal of Image and Graphics*, 27(9): 2708-2720 [DOI: 10.11834/jig.211254]
- 杨少聪, 王健, 孙运莲, 唐金辉. 2022. 多级特征全局一致性的伪造人脸检测 [J]. *中国图象图形学报*, 27(9): 2708-2720 [DOI: 10.11834/jig.211254]
- Ding F, Kuang R S, Zhou Y, Sun L, Zhu X G, Zhu G P. 2024. A survey of Deepfake and related digital forensics [J]. *Journal of Image and Graphics*, 29(2): 295-317 [DOI: 10.11834/jig.230088]
- 丁峰, 匡仁盛, 周越, 孙珑, 朱小刚, 朱国普. 2024. 深度伪造及其取证技术综述 [J]. *中国图象图形学报*, 29(2): 295-317 [DOI: 10.11834/jig.230088]
- Yao W D, Li P C, Zhao Y, Wu H C. 2025. Review of research on face deepfake detection methods [J]. *Journal of Image and Graphics*, 30(7): 2343-2363 [DOI: 10.11834/jig.240586]
- 姚文达, 李盼池, 赵娅, 吴洪超. 2025. 人脸深度伪造检测方法研究综述 [J]. *中国图象图形学报*, 30(7): 2343-2363 [DOI: 10.11834/jig.240586]
- Li Q Y, Ji Q G, Hong S D. 2019. FastFace: a real-time robust algorithm for face detection [J]. *Journal of Image and Graphics*, 24(10): 1761-1771 [DOI: 10.11834/jig.180662]
- 李启运, 纪庆革, 洪赛丁. 2019. FastFace: 实时鲁棒的人脸检测算法 [J]. *中国图象图形学报*, 24(10): 1761-1771 [DOI: 10.11834/jig.180662]