

# 多模态驱动的3D艺术头像生成技术——融合文本描述、参考图与面部动作编码

张龙庆

杭州思代尔阿特科技有限公司 浙江 杭州 310000

**摘要:** 随着数字化发展,3D艺术头像生成技术备受关注。本文聚焦多模态驱动的3D艺术头像生成技术,融合文本描述、参考图与面部动作编码展开研究。先阐述研究背景,点明传统生成技术局限,强调多模态融合在弥补不足、提升头像生成质量与丰富度上的重要意义。通过多模态数据获取与预处理,为生成模型提供优质数据;探索不同融合机制,优化数据利用。构建融合数据驱动的3D头像生成模型,并针对面部动作编码优化动态头像生成。研究成果有望革新3D艺术头像生成方式,在虚拟社交、游戏等领域展现广阔应用前景。

**关键词:** 多模态;3D艺术头像;文本描述;参考图;面部动作编码

## 1 引言

在数字化进程飞速迈进的当下,虚拟形象在社交、娱乐、游戏等众多领域的重要性与日俱增,其中3D艺术头像作为用户虚拟身份的直观代表,其生成技术成为研究热点。传统的3D头像生成技术多依赖单一数据模态,存在诸多局限,生成的头像在细节丰富度、个性化表现和动态展示上难以满足日益增长的多样化需求。多模态驱动技术的出现为这一困境带来突破曙光。通过融合文本描述、参考图与面部动作编码等多源数据,能为3D头像生成注入更全面、精准的信息。文本描述可赋予头像独特的风格与属性,参考图提供直观视觉特征,面部动作编码实现头像的动态生动表达。对多模态驱动的3D艺术头像生成技术展开深入研究,将极大推动其在各领域的应用与发展。

## 2 多模态数据获取与预处理

### 2.1 文本描述的数据收集与整理

在构建多模态驱动的3D艺术头像生成技术时,文本描述数据的收集处理尤为关键。我们广泛收集各类数据源,在社交媒体平台挖掘用户对头像的个性化需求,如“想要蓝色长发、梦幻精灵耳朵的头像”,直观反映大众对独特头像的追求;在艺术创作论坛,收集独特角色形象的文字设定,汲取创意灵感;从文学作品中摘取人物外貌的细致刻画,丰富数据来源<sup>[1]</sup>。

数据收集完成后,马上开展清洗工作,排查并去除乱码、重复内容和无关特殊符号,净化数据。接着,利用自然语言处理技术,按关键元素对文本细致分类,将包含“卷发”“直发”的文本归为发型类别,对描述面部特征、配饰等元素的文本也如此处理,精准提取核心词汇,

搭建完善的特征库。这能为后续3D头像生成提供清晰准确的语义指令,确保生成头像与文本描述高度相符。

### 2.2 参考图的采集与筛选

为给3D艺术头像生成提供优质视觉参考,参考图采集范围广泛。来源包括专业摄影作品,其高分辨率与拍摄技巧能清晰呈现面部结构和光影细节,为写实细节塑造提供关键依据;还有艺术插画,卡通风、写实风、赛博朋克风等不同风格都是创意源泉;网络上的多样素材图片也被纳入其中。

采集时,先按图像分辨率和内容相关性初步筛选,剔除模糊及与头像主题无关的图片。再从构图、色彩搭配、细节完整性等维度评估图片质量,选出最具代表性和参考价值的图。比如生成古风头像时,会着重挑选有古典服饰、发型、背景元素,且画面精致、细节丰富的参考图,作为直观视觉样本,辅助模型快速学习古风元素特征,生成风格与特征相似的3D头像。

### 2.3 面部动作编码的数据记录与预处理

借助专业的面部动作捕捉设备,在可控环境下对不同人物进行面部动作采集。让受试者做出各种表情和动作,如微笑、愤怒、眨眼、转头等,设备精准记录面部肌肉运动和关键点位移信息,并转化为面部动作编码。采集完成后,对编码数据进行预处理。先通过滤波算法去除因设备噪声或外界干扰产生的异常值,再对数据进行标准化处理,使不同受试者、不同采集时段的数据处于统一尺度。比如将所有表情动作的强度值映射到0-1的区间,方便后续在3D头像生成中,准确地将面部动作编码映射到头像模型上,实现自然、流畅的动态表情展示<sup>[2]</sup>。

## 3 多模态数据融合机制探索

### 3.1 基于特征层的融合策略

在3D艺术头像生成中，基于特征层的融合策略至关重要。此策略是在完成对文本描述、参考图以及面部动作编码的数据预处理和特征提取后，将这些不同模态数据所对应的特征向量进行直接融合。比如，把文本中提取的关于发型、五官特征的语义特征，参考图中经卷积神经网络提取的视觉图像特征，以及面部动作编码所蕴含的肌肉运动特征，按特定规则拼接在一起。通过这种融合方式，能让后续的生成模型从一开始就接触到多模态的综合信息，利于挖掘不同模态间潜在的关联，从而生成更具细节和真实感的3D头像，为头像生成提供丰富的特征基础。

### 3.2 模型中间层融合技术研究

模型中间层融合技术旨在生成模型的中间阶段融入多模态数据。当模型在处理单模态数据进行初步特征学习后，在中间隐藏层将其他模态数据的特征巧妙嵌入。例如，先由参考图数据经过前期网络层学习到图像的基础轮廓、色彩等特征，此时将文本描述对应的语义特征以及面部动作编码对应的动态特征，通过特定的融合模块，如注意力机制模块，依据不同模态特征的重要程度进行加权融合。这样做可以让模型在训练过程中不断调整对各模态数据的关注度，更好地整合多模态信息，避免早期融合可能出现的特征混淆，提升生成头像在风格、表情等方面的精准度和表现力<sup>[3]</sup>。

### 3.3 决策层融合的创新实践

决策层融合是在生成模型输出阶段进行多模态数据融合的创新尝试。此时，各个模态的数据分别经过独立的模型分支进行处理和分析。比如文本描述经自然语言处理模型转化为语义理解结果，参考图经图像分析模型提取关键视觉要素，面部动作编码经专门的动作分析模块得出动态信息。然后，在最终决策环节，根据各模态输出结果的置信度、相关性等因素，通过融合算法来决定3D头像最终的生成细节。例如，当生成一个微笑表情的头像时，面部动作编码提供表情动态信息，参考图确定面部肌肉纹理细节，文本描述补充表情风格特点，三者决策层融合，生成符合多模态信息的生动3D头像，提高生成结果的准确性和可靠性。

## 4 融合数据驱动的3D头像生成模型构建

### 4.1 生成模型架构选型与分析

在3D头像生成领域，模型架构的选择至关重要。常见的生成对抗网络（GAN）通过生成器与判别器的对抗博弈，能生成高分辨率且逼真的图像，但其训练过程不稳定，易出现模式坍塌。变分自编码器（VAE）则基于概率模型，能学习数据分布并生成具有多样性的样本，不过生成图像的细节可能不够清晰。考虑到多模态数据融合的复杂性，需综合权衡。例如，在处理参考图与文本描述时，GAN能更好捕捉图像视觉特征，而VAE在融合语义信息时展现优势，需根据实际需求分析各模型在多模态3D头像生成中的适用性。如图一所示：

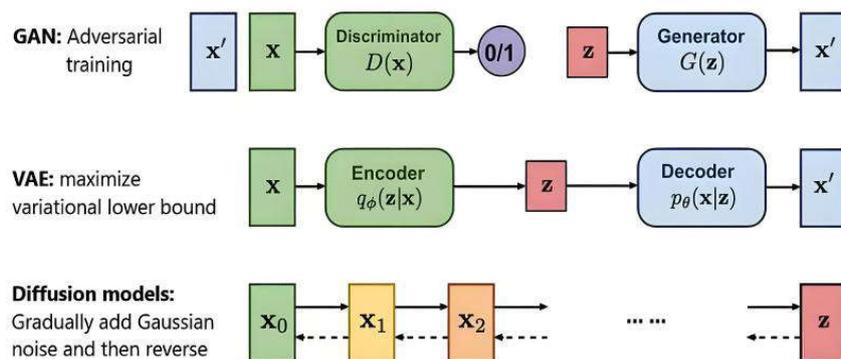


图1 生成式模型对比图

### 4.2 针对多模态融合的模型改进

为适配多模态数据，需对基础模型进行改造。首先，在模型输入层，设计多模态数据接入模块，使文本描述、参考图和面部动作编码能有效整合。如将文本特征向量、图像特征矩阵和动作编码序列进行拼接或按特

定规则融合。其次，优化模型中间层，增加跨模态交互层，让不同模态数据在模型内部充分交流。可利用注意力机制，使模型聚焦关键信息。最后，在输出层，调整网络结构以输出符合3D头像生成要求的参数，如顶点坐标、纹理映射等，从而生成满足多模态融合需求的高质

量3D头像<sup>[4]</sup>。

#### 4.3 模型训练与性能评估

模型训练采用大量标注好的多模态数据集，包含丰富文本描述、多样参考图及对应面部动作编码。训练过程中，设置合理超参数，如学习率、迭代次数等。利用随机梯度下降等优化算法，不断更新模型参数，使生成头像与真实样本在视觉效果、语义表达和动作表现上尽可能相似。性能评估采用多指标衡量，包括峰值信噪比（PSNR）评估图像清晰度，结构相似性指数（SSIM）衡量图像结构相似度，以及引入语义匹配度指标评估生成头像与文本描述的契合度，从多维度全面评估模型生成3D头像的质量与准确性，确保模型性能达到预期。

### 5 基于面部动作编码的动态头像生成优化

#### 5.1 面部动作与3D头像动态映射关系建立

准确建立面部动作与3D头像动态的映射关系，是实现自然生动动态头像生成的关键。借助面部动作编码系统，细致分析面部肌肉运动规律，将不同动作单元与3D头像的骨骼、关节点相连接。例如，嘴角上扬动作对应头像面部肌肉的拉伸，精确到肌肉变形程度和骨骼旋转角度。通过大量动作数据的采集与分析，构建全面且精准的映射表，确保每个面部动作都能在3D头像上真实还原，为动态头像的生成奠定坚实基础，让头像能细腻展现各种表情与动作变化<sup>[5]</sup>。

#### 5.2 动态头像生成的实时性优化策略

在追求动态头像高质量生成的同时，实时性至关重要。一方面，优化算法架构，采用轻量级模型与并行计算技术，减少计算量并提升计算效率。例如，运用剪枝算法精简神经网络模型，去除冗余连接，加速运算。另一方面，利用硬件加速技术，如GPU并行计算，充分发挥硬件性能优势。同时，优化数据传输与处理流程，采用缓存机制减少数据读取延迟，确保在有限硬件资源下，动态头像能快速响应输入的面部动作编码，实现流畅、实时的动态展示，满足用户即时交互需求。

#### 5.3 动态头像生成的稳定性改进

动态头像生成的稳定性直接影响用户体验。为解决生成过程中的闪烁、卡顿等问题，从多方面入手改进。在

数据处理环节，加强对噪声数据的过滤与修正，保证面部动作编码数据的准确性和连续性。在模型训练上，增加稳定性约束条件，使模型在不同输入情况下都能生成稳定的头像动态。优化渲染流程，合理分配计算资源，避免因渲染压力过大导致的画面异常。通过这些措施，大幅提升动态头像生成的稳定性，让头像在长时间动态展示中保持流畅、稳定，为用户提供优质的视觉效果。

### 6 结语

本研究围绕多模态驱动的3D艺术头像生成技术，融合文本描述、参考图与面部动作编码展开深入探索，取得了一系列成果。通过搭建多模态数据融合框架，为3D头像生成引入了丰富且精准的信息，有效解决了传统生成技术存在的单一性问题。构建的生成模型，经反复优化，能够生成高逼真度、高度个性化的3D艺术头像，且在动态生成方面，借助面部动作编码实现了生动自然的表情与动作变化。该技术未来在虚拟社交、游戏、影视等领域有着巨大的应用潜力，有望重塑用户的虚拟交互体验。然而，研究仍存在数据隐私保护、模型泛化能力有待提升等问题，后续可从加密技术应用、拓展数据多样性等方向持续改进，推动多模态驱动的3D艺术头像生成技术迈向新高度。

### 参考文献

- [1]周治国,马文浩.一种多层多模态融合3D目标检测方法[J].电子学报,2024,52(3):696-708.
- [2]王彩玲,闫晶晶,张智栋.基于多模态数据的人体行为识别方法研究综述[J].计算机工程与应用,2024,60(9):1-18.
- [3]冯霞,梁宇龙,卢敏,左海超.基于NNC-EPNet的多模态融合3D目标检测[J].北京交通大学学报,2024,48(5):78-87.
- [4]张青青,曾冉,浦奔放,张学军,宋冬雷,吴曦.3D Slicer多模态融合及三维重建技术在前庭神经鞘瘤手术中保护面神经功能的应用探索[J].临床神经外科杂志,2024,21(6):641-647.
- [5]陈娜.基于深度卷积网络的3D人脸重构算法[J].激光与红外,2022,52(6):923-930.