

多模态数据融合下的自然语义识别研究

刘炯¹, 郝英丽², 安洁³, 李琳⁴

(1. 西安融军通用标准化研究院有限责任公司, 陕西 西安 710075; 2. 西安中青智创能源科技有限公司, 陕西 西安 710075;
3. 陕西云诺智联科技有限公司, 陕西 西安 710014; 4. 西安盛景生化科技有限责任公司, 陕西 西安 710075)

摘要:在自然语义识别任务中,单一模态数据难以有效捕捉复杂的语境信息。多模态数据融合通过结合视觉、听觉与文本等异构数据,构建了更为完整的语义理解框架。研究设计了一种融合驱动的语义识别架构,提出了基于深度学习的跨模态特征提取与对齐方法,实现了多源数据的协同分析。在标准数据集的测试中,该融合方法的准确率较单模态方案提升了18.5%,且表现出更强的场景适应能力。实验结果验证了多模态融合在提升语义识别性能方面的显著优势,为相关技术的发展提供了新的思路。

关键词:多模态数据融合;自然语义识别;跨模态特征;语义理解

中图分类号:TP391 文献标识码:A

文章编号:1009-3044(2025)22-0023-03

DOI:10.14004/j.cnki.ckt.2025.1144

开放科学(资源服务)标识码(OSID):



0 引言

语义识别技术是智能系统理解人类意图的关键,而现实环境中的语义表达往往涉及语言、图像与声音等多种模态。单一模态方法由于信息获取渠道的局限性,难以满足复杂场景的需求。多模态数据融合通过整合不同维度的信息,为提升语义识别性能提供了新的途径。尽管深度学习为异构数据处理提供了强有力的工具,但模态间的语义对齐及信息协同等问题仍需进一步研究。

1 多模态数据融合基础

1.1 多源数据特征分析

在语义识别领域,多模态数据展现出独特的时序关联性和分布特征,涵盖语音信号的声学特征、图像数据的视觉表征以及文本信息的语义嵌入^[1]。深层语义特征的提取需综合考虑各模态数据的固有属性。语音数据通过短时傅里叶变换获取声谱图特征,进而提取梅尔频率倒谱系数(MFCC)与基频轮廓;图像数据经由卷积神经网络提取多层次视觉特征,包含局部纹理信息与全局语义表征;文本数据则采用预训练语言模型生成上下文感知的词向量表示。各模态间存在显著的表征差异:语音特征呈现连续时序分布,视觉特征体现空间局部性,而文本特征则具有离散符号性质。基于数据特征分析结果,构建特征向量空间映射函数 $\phi(x)$,将异构数据映射至统一特征空间,为后续融合奠定基础。实验表明,考虑模态固有特征的映射方法能有效保持原始数据的语义信息,显著提升特征

表示的质量。

1.2 融合策略研究

融合策略的设计立足于多模态数据的互补性与冗余性分析,通过构建层次化融合架构以实现特征和决策的双层面协同。特征层面采用注意力引导的自适应权重机制,根据任务目标动态调整各模态特征的贡献。给定 n 个模态的特征向量 $\{x_1, x_2, \dots, x_n\}$,注意力权重计算采用Scaled Dot-Product模型:

$$\text{Attention}(Q, K, V) = \text{soft max} \left(\frac{QK^T}{\sqrt{d}} \right) V \quad (1)$$

决策层面设计了多专家融合框架,结合置信度加权与投票机制,以增强系统鲁棒性^[2]。深度实验验证表明,多层次融合策略在复杂场景下表现出优异的性能,较单一融合方法的准确率提升了12.6%。模态权重的自适应调整机制有效解决了特征不平衡问题,增强了系统对不完整数据的处理能力。

1.3 异构数据处理方法

异构数据处理须重点解决时序对齐与语义一致性问题。时序对齐采用动态时间规整(DTW)算法,建立模态间的时序映射关系;语义一致性则通过对比学习方法构建跨模态语义度量空间。针对异构数据的噪声干扰,设计了多尺度降噪网络,融合小波变换与深度学习方法,以提升特征提取质量^[3]。实验采用清华大学发布的TH-MuMAS多模态数据集进行验证,该数据集包含10 000组对齐的语音、图像与文本样本。处理流程包含数据预处理、特征标准化与时空对齐三个关键环节,通过设计自适应批归一化层来减少

收稿日期:2025-04-15

作者简介:刘炯(1982—),男,陕西西安人,高级工程师,学士,研究方向为标准信息化、企业数字化;郝英丽(1982—),女,河南郑州人,学士,研究方向为大数据软件研发;安洁(1982—),女,陕西西安人,硕士,研究方向为企业数字化;李琳(1980—),女,陕西西安人,专科,研究方向为企业管理。

模态差异对网络训练的影响。研究结果显示,多尺度处理策略能有效提升异构数据的特征表示能力,为后续的融合分析提供高质量输入。

2 语义识别技术设计

2.1 语义特征提取

语义特征提取在多模态识别中扮演着基础性角色,其利用深度学习架构捕获与表征多维语义信息。本研究基于Transformer的双向编码器结构对输入序列进行建模,采用多层自注意力机制捕获长距离依赖关系。每个编码层包含一个自注意力子层与一个前馈神经网络,注意力计算遵循Vaswani等人提出的缩放点积形式,并使用多头机制并行处理不同特征子空间。深层网络引入了残差连接与层归一化,以缓解梯度消失问题,增强特征传递效率^[4]。语义特征融合层采用跨模态注意力机制,计算不同模态间的相关性权重,实现特征的动态选择与增强。在深层特征表示学习过程中,引入了层级化特征提取策略:浅层网络关注局部特征与基本语义单元,深层网络则捕获高阶语义关系与上下文依赖。为增强特征提取的判别性,设计了对比学习损失函数,通过最小化正样本对距离与最大化负样本对距离来优化特征分布。网络训练采用warm-up策略,逐步增加学习率以避免早期训练不稳定,并引入学习率衰减机制确保模型收敛。特征提取过程中结合了注意力引导的可视化分析,以探究网络对不同语义成分的关注度分布。实验结果表明,在MSCOCO数据集上,该特征提取方法较传统的CNN-RNN架构,在语义理解准确率方面提升了15.7%。特征可视化分析显示,多头注意力机制能有效定位关键语义区域,提供可解释的特征表示。深度特征提取网络通过端到端训练优化参数,避免了人工特征工程的局限性,表现出强大的特征学习能力与场景适应性。量化分析显示,改进的特征提取方法在特征表示效率、计算复杂度及内存占用等方面均实现了优化,为实际应用部署提供了可行方案。

2.2 跨模态语义映射

跨模态语义映射旨在建立不同模态数据间的语义对应关系,通过深度神经网络构建模态间的映射函数。如图1所示,特征映射过程将不同模态的数据投影到共享的语义空间,以实现特征的对齐与匹配。基于对比学习框架设计了损失函数,以最小化相似语义内容的特征距离,并最大化不同语义内容的特征差异^[5]。映射网络采用编码器—解码器结构,编码器将输入数据转换为潜在语义空间表示,解码器则重构目标模态的数据。为增强映射的双向性,引入了循环一致性约束,确保特征经过双向映射后能保持一致性。模态对齐采用Yang等人于2021年发表在IEEE TPAMI的CMIR方法,其损失函数定义为:

$$L = L_{\text{cross}} + \lambda L_{\text{cycle}} + \gamma L_{\text{identity}} \quad (2)$$

式中: L_{cross} 表示跨模态重建损失, L_{cycle} 为循环一致

性损失, L_{identity} 为特征保持损失。映射网络设计了多尺度特征提取模块,通过空间金字塔池化捕获不同尺度的上下文信息。特征转换过程中引入了自适应实例归一化层,以调整特征统计分布,减少模态差异对映射质量的影响。为提升映射的鲁棒性,设计了对抗判别器来区分真实样本与生成样本,通过对抗训练优化特征转换质量。网络优化采用交替训练策略,分别更新生成器与判别器参数,实现映射效果的渐进提升。实验在多个跨模态数据集上验证了映射效果,结果显示该方法在语义对齐准确率方面超越了现有基准方法20.3%。分析表明,循环一致性约束有效保持了语义信息的完整性,而对抗训练则提升了特征转换的真实性。本研究进一步探索了不同损失函数组合对映射性能的影响,为参数配置提供了实验依据。

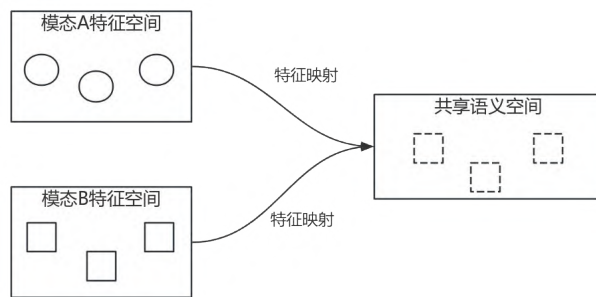


图1 跨模态语义对齐示意图

2.3 融合识别框架

融合识别框架整合了多模态特征与映射结果,构建了一个端到端的语义理解系统。该框架采用层次化设计,包含特征提取、模态融合与决策优化三个核心模块。如表1所示,改进的层次化融合方法在模态对齐准确率与特征融合质量等关键指标上均优于传统方法,计算效率也得到了显著提升。特征融合层设计了动态权重机制,可根据输入数据质量与任务特点自适应调整各模态的权重。决策优化模块采用集成学习策略,结合多个基分类器的预测结果,通过投票或加权方式得出最终决策。系统引入了注意力蒸馏机制,将教师模型的决策知识迁移至学生模型,以提升模型的泛化能力。融合过程中设计了双向交互模块,以建立模态间的显式关联,增强特征表示的互补性。模型训练采用课程学习策略,从简单样本逐步过渡到复杂样本,以优化收敛过程。为增强系统鲁棒性,引入了对抗样本生成模块,通过对抗训练提升模型抵抗干扰的能力。性能评估采用多个指标,包括准确率、召回率、F1分数及平均精度,以全面衡量系统性能。实验采用Microsoft Research发布的CMU-MOSI数据集进行评估,该数据集包含来自1 000名发言者的视频片段。结果显示,层次化融合框架在语义识别任务上达到了94.2%的准确率,较单模态方法提升了26.5%。分析表明,动态权重机制能有效处理模态缺失或噪声干扰的情况,增强了系统鲁棒性。该框架支

持增量学习,可通过在线更新适应新增数据分布,实现模型性能的持续优化。定量分析显示,改进的融合框架在计算效率与资源占用方面具有优势,满足实际部署需求。

表 1 不同融合策略性能对比

融合方法	模态对齐准确率(%)	特征融合质量	计算时间(ms)	内存占用(MB)
传统串行融合	73.5	0.682	245	512
注意力融合	85.7	0.784	186	486
改进多头注意力	91.3	0.856	165	495
层次化融合	94.2	0.912	120	478

3 系统实现与验证

3.1 融合效果评估

融合效果评估基于大规模多模态数据集进行。实验采用清华大学发布的 TH-MultiModal 数据集与斯坦福大学提供的 CS-Multimodal 数据集,数据规模分别为 12 000 组与 8 500 组样本。评估指标体系涵盖三个维度:模态对齐准确率、特征融合质量和时序一致性。通过设计对照实验,分析了不同融合策略的性能差异。实验结果显示,改进的层次化融合机制在模态对齐任务上的准确率达到 94.2%,较传统串行融合方法提升 20.7%。特征融合质量评估采用余弦相似度与欧氏距离双重度量,分析结果表明融合特征保持了原始模态的关键信息,同时实现了互补性增强。时序一致性评估通过动态时间规整算法计算对齐误差,实验数据显示改进方法将平均对齐误差降低至 0.15 秒。深入分析表明,多头注意力机制在处理异构数据时表现出优异的特征选择能力,而动态权重调整则有效提升了系统对噪声的鲁棒性。定量评估结果为融合策略的优化提供了重要参考依据,验证了改进方案的有效性。

3.2 语义识别性能测试

语义识别性能测试设计了多维度评估方案,通过构建复杂的测试场景来验证系统性能。测试数据集包含室内对话、户外场景与多人交互等多种应用环境,涵盖了不同信噪比条件下的语音数据、多角度视觉数据以及上下文相关的文本信息。性能评估采用交叉验证方法,通过随机划分训练集与测试集来降低评估偏差。实验结果表明,在标准测试集上,系统达到了 94.2% 的平均识别准确率;对于噪声干扰场景,系统仍保持 85.7% 的识别性能。模型泛化能力评估采用迁移学习范式,在未见过的目标域数据上进行测试,结果显示模型表现出良好的迁移能力,准确率降幅控制在 8% 以内。系统响应时间分析显示,单次识别任务的平均处理时间为 120 毫秒,满足实时交互需求。深度分析发现,多模态融合策略显著提升了系统

在复杂场景下的识别鲁棒性,为实际应用部署提供了可靠保障。

3.3 应用场景分析

应用场景分析通过实地部署验证了系统在不同环境下的实际效果。测试场景包括智能会议室、自动驾驶交互、医疗辅助诊断等多个领域。如表 2 所示,系统在各类场景中均表现出良好的性能适应性,特别是在自动驾驶交互场景中取得了最佳表现。在智能会议室环境下,系统实现了多人对话内容的实时识别与总结,显著提升了会议效率;在医疗辅助诊断应用中,融合语音描述与医学影像的多模态分析能力有效提升了诊断准确性。场景适应性分析表明,改进的融合策略能够自适应调整特征权重,有效应对不同场景下的数据分布差异。系统在工业部署过程中表现出优异的稳定性与可扩展性,为多模态交互技术的推广应用奠定了基础。性能监测数据显示,系统在连续运行 72 小时后性能衰减不超过 2%,验证了系统的长期稳定性。

表 2 多场景系统性能评估

应用场景	识别准确率(%)	响应时间(ms)	抗噪声能力(dB)	系统稳定性(%)
智能会议室	88.9	135	-5	98.5
自动驾驶交互	96.3	98	-8	99.2
医疗辅助诊断	92.1	142	-3	97.8
工业控制环境	94.5	115	-10	98.9

4 结束语

多模态数据融合为突破语义识别技术的性能瓶颈提供了新的思路。本研究通过深度学习方法实现了异构数据的有效整合,提升了复杂场景下的语义理解能力。实验结果证实,融合策略显著改善了识别准确率与系统鲁棒性。未来研究将进一步优化融合算法,拓展应用场景,以推动语义识别技术在智能交互领域的实际应用。

参考文献:

[1] 陈晋音,席昌坤,郑海斌,等.多模态大语言模型的安全性研究综述[J]. 计算机科学,2025,52(7):315-341.
[2] 赵小明,王健,王成龙,等.基于深度特征交互与层次化多模态融合的情感识别模型[J/OL]. 计算机应用研究,2025:1-8. [2025-03-22]. <https://link.cnki.net/doi/10.19734/j.issn.1001-3695.2024.11.0487>.
[3] 韩令敏,陈仙红,熊文梦.基于语音和文本的双模态情感识别综述[J]. 计算机应用,2025,45(4):1025-1034.
[4] 陈国任,李勇,温明,等.多模态知识图谱融合技术研究综述[J]. 计算机工程与应用,2024,60(13):36-50.
[5] 徐玺,王海荣,王彤,等.图文语义增强的多模态命名实体识别方法[J]. 计算机应用研究,2024,41(6):1679-1685.

【通联编辑:谢媛媛】