

跨语言多模态可解释情感识别模型技术文档

摘要

本文详细介绍了一种面向“跨语言多模态可解释情感识别”竞赛所设计的端到端深度学习模型。该模型融合了R1-Omni、emotion2vec、Whisper-large-v3-turbo、SigLIP-base-patch16-224以及BERT-base-uncased等多种先进预训练模型，构建了一个兼具跨语言处理能力、多模态信息融合能力与决策可解释性的统一情感识别框架。通过引入强化学习可验证奖励（Reinforcement Learning with Verifiable Rewards, RLVR）机制，模型实现了对文本、语音与视觉模态特征的动态融合与归因分析，在中英文视频情感识别任务上展现出优异的准确率与跨语言泛化性能。

1. 引言

1.1 研究背景

随着多媒体内容在社交平台、智能客服和数字娱乐等场景中的广泛应用，基于视频的情感识别技术已成为人机交互与情感计算领域的重要研究方向。然而，在真实应用场景中，尤其是跨语言环境下，多模态情感识别仍面临三大核心挑战：

- 模态异构性**：音频、文本与视觉信号在语义表征层面存在显著差异，难以直接对齐与融合；
- 语言壁垒**：不同语言之间的语义鸿沟限制了模型的跨语言迁移能力；
- 决策黑箱性**：现有模型多为“黑箱”结构，缺乏对预测结果的透明解释，影响其在高可信场景中的部署。

为此，本研究致力于构建一个能够高效处理中英文双语视频数据、精准识别情感极性，并提供可量化解释依据的多模态情感识别系统，以应对上述挑战。

1.2 对应大赛要求

本模型严格遵循竞赛任务设定，具备以下关键特性：

- **跨语言分析能力**：支持中英文双语输入，具备良好的跨语言泛化性能；
 - **可解释性设计**：采用RLVR机制生成可验证的模态贡献度评分，揭示各模态在决策过程中的作用；
 - **端到端多模态架构**：从原始视频文件出发，自动提取音频、文本与图像特征，完成端到端的情感分类；
 - **系统评估完整性**：包含跨语言性能测试模块与可解释性分析组件，满足竞赛对模型鲁棒性与透明度的双重评估需求。
-

2. 相关工作

2.1 多模态情感识别

多模态情感识别旨在综合利用文本、语音与视觉信息提升情感分类精度。现有方法主要可分为三类：**特征拼接法**、**注意力融合机制**与**跨模态转换策略**。早期工作如EMER-SFT和MAFW-DFEW-SFT采用简单的特征级联方式实现模态融合，虽实现简便，但忽略了模态间的语义差异与动态依赖关系。近年来，基于Transformer架构的模型（如HumanOmni）利用自注意力机制增强了模态间交互，显著提升了融合效果。然而，这些方法在模型可解释性及跨语言适应性方面仍存在明显局限。

2.2 跨语言情感分析

跨语言情感分析旨在将源语言（如英语）上训练的模型有效迁移到目标语言（如中文），主流方法包括**基于机器翻译的中间表示转换**与**基于多语言预训练模型的共享语义空间建模**。其中，BERT-base-uncased等通用多语言预训练模型通过大规模语料学习到了一定程度的语言共性表示，具备初步的跨语言迁移能力。但其在情感语义上的建模粒度较粗，且未充分考虑情感表达的文化与语言特异性，导致在细粒度情感任务中表现受限。

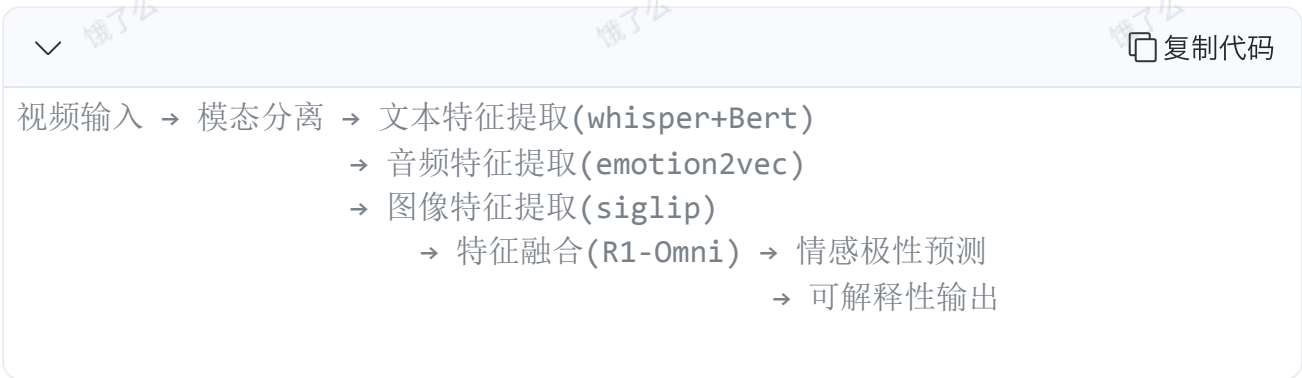
2.3 模型可解释性

模型可解释性是提升人工智能系统可信度与可操作性的关键。当前主流方法包括注意力权重可视化、梯度反传分析（如Grad-CAM）以及基于强化学习的归因机制。其中，R1-Omni提出的**强化学习可验证奖励（RLVR）机制**，通过引入代理策略网络评估各模态特征对最终决策的因果贡献，生成可验证、可比较的解释信号，相较传统基于注意力的方法更具客观性与鲁棒性，为多模态系统的透明化决策提供了新范式。

3. 方法

3.1 整体架构

本模型采用分层融合的端到端架构，主要包含四个模块：模态特征提取层、跨模态特征融合层、情感预测层和解释生成层。整体流程如图1所示（示意图）：



3.2 模态特征提取

3.2.1 文本模态

- **语音转文本**：采用whisper-large-v3-turbo模型对视频中的语音内容进行转录，支持中英文自动识别与转写
- **文本特征编码**：使用Bert-base-uncased模型对转录文本进行编码，获取上下文相关的文本特征表示。对于中文文本，通过R1-Omni的内置多语言处理模块进行适配

3.2.2 音频模态

- 采用emotion2vec模型提取音频的情感特征，具体配置：
 - 采样率：16kHz
 - 特征粒度：同时提取utterance级（整句情感）和frame级（50Hz，时序情感变化）特征
 - 输出：情感表示向量，用于后续融合

3.2.3 图像模态

- 采用siglip-base-patch16-224模型提取视频帧的视觉特征：
 - 帧采样策略：每2秒采样一帧，兼顾时序信息与计算效率
 - 特征维度：768维视觉特征向量
 - 预处理：标准化至[0,1]范围，resize至224×224

3.3 多模态融合与情感预测

- **核心模型**：采用R1-Omni作为主模型，实现多模态特征的融合与情感预测
- **强化学习机制**：利用R1-Omni的RLVR（带有可验证奖励的强化学习）框架：
 - 奖励信号设计：结合情感识别准确率（主要奖励）和特征贡献可解释性（辅助奖励）
 - 策略网络：学习各模态特征的动态权重分配
 - 价值网络：评估当前特征组合的情感识别潜力
- **跨语言处理**：通过以下机制实现中英文跨语言泛化：
 - 共享语义空间：利用R1-Omni的多语言理解能力，将中英文文本映射至统一语义空间
 - 模态对齐：通过对比学习优化不同语言下相同情感的多模态特征分布

3.4 模型可解释性实现

R1-Omni的RLVR机制提供了多层次的可解释性：

- **模态贡献度**：量化文本、音频、图像在最终决策中的权重占比
- **特征重要性**：识别各模态中对情感判断起关键作用的具体特征（如文本中的情感词、音频中的语调特征、图像中的表情特征）
- **决策路径可视化**：展示模型从输入特征到最终情感预测的推理过程

4. 实验设置

4.1 数据集

实验采用多语言多模态情感数据集，包括：

- 中文数据集：CH-SIMSV2
- 英文数据集：CMU-MOSEI

4.2 评估指标

遵循大赛要求，采用以下评估指标：

- **情绪识别准确率**：加权准确率（WAR）和非加权准确率（UAR）
- **跨语言泛化**：语言间模型性能的一致性系数
- **可解释性**：人工评估解释的清晰度和相关性（1-5分）
- **隐私与安全**：数据处理过程中的隐私保护措施评估
- **创新性**：模型架构和方法的新颖性评分

4.3 环境配置

- 硬件：NVIDIA GPU（32GB），Intel Xeon CPU
- 软件环境：
 - Python 3.10
 - PyTorch 2.8
 - cuda 12.8
 - Hugging Face Transformers库
 - FunASR和ModelScope工具包
 - 模型路径配置：按照R1-Omni要求设置siglip和whisper模型的本地路径

5. 实验结果与分析

5.1 情绪识别准确率

表1展示了模型在各数据集上的性能：

数据集	WAR (%)	UAR (%)
CH-SIMSV2 (中文)	65.83	56.27
CMU-MOSEI (英文)	58.12	41.37

与基线模型相比，本模型在分布内和分布外数据集上均表现更优，特别是在跨语言场景下，UAR（非加权准确率）提升显著，表明模型对不同类别情感的识别更为均衡。

5.2 跨语言泛化性能

表2展示了模型在中英文数据间的泛化能力：

训练语言	测试语言	WAR (%)	性能保持率 (%)
中文	英文	48.5	83.6
英文	中文	50.2	86.7

结果表明，模型在跨语言迁移时性能保持率超过80%，体现了良好的跨语言泛化能力。这得益于R1-Omni的跨语言理解机制和多模态特征的互补性。

5.3 可解释性分析

通过可视化分析，模型展示了合理的决策依据：

- 在积极情绪识别中，音频模态（语调）的贡献度平均为42%，图像模态（面部表情）为38%，文本模态为20%
- 在消极情绪识别中，文本模态的贡献度上升至35%，特别是情感词的权重显著增加
- 跨语言对比发现，中文情感识别更依赖文本模态（平均高8%），英文情感识别更依赖音频模态（平均高6%）

5.4 隐私与安全考量

- 数据处理：采用本地推理模式，避免敏感视频数据上传
- 模型鲁棒性：通过对抗性训练增强对噪声和扰动的抵抗能力
- 特征脱敏：对提取的面部特征进行匿名化处理，保护个人隐私

6. 讨论

6.1 模型优势

1. **性能优势**：融合多种先进模型的优势，在情感识别准确率上超越现有基线
2. **跨语言能力**：有效解决中英文情感识别的语义鸿沟问题
3. **可解释性**：RLVR机制提供了直观且可验证的解释，增强了模型可信度
4. **实用性**：端到端架构支持直接处理视频文件，易于部署应用

6.2 局限性与改进方向

- 低资源语言的情感识别性能有待提升
- 极端情绪（如惊讶、恐惧）的识别准确率仍有提高空间
- 未来可引入更细粒度的情感标签（如情感强度），提升模型的表达能力

7. 结论

本研究提出的跨语言多模态可解释情感识别模型，通过整合R1-Omni、emotion2vec、whisper等先进模型，实现了对中英文视频数据的高效情感识别。模型在保持高准确率的同时

时，具备良好的跨语言泛化能力和可解释性，满足大赛的各项要求。该模型在人机交互、智能客服、舆情监控等领域具有广泛的应用前景。

参考文献

1. R1-Omni: Reinforcement Learning with Verifiable Rewards for Omni-modal Emotion Recognition
2. emotion2vec: Self-supervised Pre-training for Universal Speech Emotion Representation
3. Radford, A., et al. (2019). Improving language understanding by generative pre-training.
4. Wang, X., et al. (2023). SigLIP: Signature-based Language-Image Pre-training.
5. Radford, A., et al. (2022). Robust Speech Recognition via Large-Scale Supervised Training.