

执行摘要

近年来，多模态情感分析进入快速发展期，新模型层出不穷。本文系统梳理了2019年以来的代表性模型与论文，包括基于深度学习、多模态融合和图神经网络的最新SOTA方法（如MulT、MISA、Self-MM、MMGCN、COGMEN等），并列出了其核心思路、所用数据集及指标^{1 2}。在此基础上，重点调研了COGMEN架构及强化学习（RL）方法在情感分析中的应用。COGMEN（2022）利用上下文感知的图神经网络捕捉对话中本地和全局信息，对IEMOCAP、CMU-MOSEI等数据集的情感分类任务取得了SOTA效果²。强化学习方面，已有工作如GME-LSTM(A)（2018）在多模态情感分类中应用策略梯度进行模态选择，R1-Omni（2025）和AffectGPT-R1（2025）用RL优化情绪识别任务指标^{3 4}，Li等人（2025）提出的EMO-RL通过“情感相似度加权奖励”训练大规模音频-语言模型，在IEMOCAP/MELD上达成SOTA性能⁵。同时，需要指出多模态情感分析仍面临挑战：模态对齐困难、长时序依赖、数据稀缺与标签偏差、噪声干扰、可解释性不足、跨域泛化差、实时性与隐私等问题^{1 6}。针对这些问题，本文提出多项创新研究方向：包括基于RL的模态自适应融合、对话场景下的上下文优化、结合自监督与对抗训练的预训练策略、引入心理生理信号的多源融合、跨域迁移与元学习方法等，每个方向都给出技术思路、预期改进、难点风险及实验设计框架。最后，给出了短中长期研究路线图（甘特图形式），并推荐可复现的开源工具（PyTorch、DGL、HuggingFace Transformers等）和数据集（CMU-MOSI/MOSEI、IEMOCAP/MELD、CH-SIMS、DEAP等）及其预处理方案。

多模态情感分析SOTA模型与论文

近五年内，多模态情感分析方法主要集中在有效融合文本、音频、视觉（以及生理信号）特征以提升情感分类或回归性能。下表总结了代表性模型及其核心贡献：

- **GME-LSTM(A)**（2018年，Chen等）：提出“门控多模态嵌入LSTM+时序注意力”架构，在词级进行模态融合，通过策略梯度优化模态门控，过滤噪声模态。CMU-MOSI数据集上实现了当时的分类/回归SOTA⁷。（[代码](#)）
- **MulT (Multimodal Transformer)**（2019年，Tsai等）：引入跨模态注意力Transformer，直接在非对齐时序上进行双向交叉注意力融合，无需显式对齐，显著提升了IEMOCAP、MOSI/MOSEI等数据集上的性能¹。实验表明，MulT在对齐和非对齐数据上均超越先前方法¹。（[代码](#)）
- **MISA**（2020年，Hazarika等）：提出模态不变/特定空间表示，将每种模态映射到不变子空间（捕捉通用情感成分）和特定子空间，减少跨模态分布差异。在CMU-MOSI/MOSEI上实现显著领先⁸。（[代码](#)）
- **Self-MM**（2021年，Yu等）：设计自监督标签生成模块，自动生成各模态的独立标注，并联合训练多模态与单模态任务。该自监督多任务方法在CMU-MOSI、MOSEI上达到当时SOTA⁹。（[代码](#)）
- **MMGCN**（2021年，Hu等，ACL）：基于对话情感识别，提出多模态融合图卷积网络。通过构建包含说话者信息的融合图，对话中每个语句节点使用GCN建模不同模态的依赖关系，在IEMOCAP和MELD上取得F1的显著提升^{10 11}。（[代码](#)）
- **COGMEN**（2022年，Joshi等，NAACL）：引入“上下文图神经网络”架构，同时利用局部对话上下文和全局对话语境信息。采用GraphTransformer处理说话者关系，实现对话中每句话单人情感识别。在IEMOCAP（6类和4类情感）和CMU-MOSEI（情感极性分类）任务上均刷新SOTA^{2 12}。（[代码](#)）

- **EMO-RL** (2025年, Li等, EMNLP Find.) : 针对大型音频-语言模型 (如SpeechLLM), 引入基于“情感规则”的RL训练策略 (Emotion Similarity-Weighted Reward + 结构化推理)。该方法显著提升了模型对MELD和IEMOCAP情感识别的泛化性能, 实现了SOTA水平 ⁵。
- **R1-Omni** (2025年, Zhao等) : 首次将可验证奖励的强化学习 (RLVR) 应用于多模态大模型的情感识别。利用RLVR优化Omni模型的推理能力、准确率和泛化性, 在内外部分布数据上均有明显改善, 并能够解释不同模态 (视听) 对情感识别的贡献 ⁴。
- **AffectGPT-R1** (2025年, Lian等) : 针对开放类别情感识别 (OV-MER), 使用强化学习将情绪轮 (Emotion Wheel) 评价指标作为奖励函数, 通过策略优化最大化该奖励。引入辅助奖励正则化情绪推理与预测行为, 在MER-UniBench等数据集上刷新SOTA ³。
- **RL驱动噪声增强 (Data Augmentation)** (2024年, Ranjan等, Interspeech) : 提出基于强化学习的数据增强方法, 在语音情感识别中用策略梯度挑选噪声频带生成增强样本。实验显示RL方法优于随机噪声, 对IEMOCAP的4类情感分类在噪声鲁棒性测试 (跨语料和跨语言) 中效果更佳 ¹³。

下表可为以上主要模型的概览:

模型	核心方法	数据集	性能指标与基线对比	代码链接
GME-LSTM(A) ⁷	模态门控 + LSTM+时序注意力 (策略梯度优化)	CMU-MOSI	以词级融合多模态, 去噪效果佳, 取得当时CMU-MOSI分类、回归任务的SOTA	github.com/...
MuT ¹	跨模态 Transformer (交叉注意力, 无需对齐)	CMU-MOSI/ MOSEI, IEMOCAP	无需对齐时序, 长距离依赖建模效果好。实验显示超越前序方法, 效果“远超”之前SOTA ¹	github.com/yaohungt/Multimodal-Transformer
MISA ⁸	不变/特定子空间投影	CMU-MOSI/ MOSEI	将每模态投影到共享与特定空间, 减少模态间差异。MOSI、MOSEI上实验显著优于前端方法 ⁸	github.com/snipsco/misa
Self-MM ⁹	自监督标签生成 + 多任务学习	CMU-MOSI/ MOSEI	自动生成各模态独立标签联合训练。MOSI、MOSEI分类性能超越当时SOTA ⁹	github.com/thuiar/Self-MM
MMGCN ^{10 11}	融合图卷积网络 (对话情感 ERC)	IEMOCAP, MELD	构建包含说话人关系的多模态图, 显著提升IEMOCAP与MELD对话情感F1值, 相较DialogueRNN/GCN提升明显 ^{10 11}	github.com/RUCAIBox/MMGCN
COGMEN ^{2 12}	上下文感知GNN (对话情感 ERC)	IEMOCAP, CMU-MOSEI	同时利用局部与全局对话信息的图神经网络。IEMOCAP6类、4类和MOSEI情感分类任务均刷新SOTA ^{2 12}	github.com/Exploration-Lab/COGMEN

模型	核心方法	数据集	性能指标与基线对比	代码链接
EMO-RL ⁵	RL优化音频-语言模型推理	IEMOCAP, MELD	在预训练音频-语言模型基础上，引入“情感相似度加权奖励”和结构化推理，MELD/IEMOCAP上达成SOTA性能 ⁵	(待开源)
R1-Omni ⁴	可验证奖励RL (多模态大模型)	(vis+audio)	首次使用RLVR优化多模态大语言模型情感识别：提升推理能力与识别准确度，在跨域测试中表现更稳健 ⁴	(待开源)
AffectGPT-R1 ³	RL优化开放情绪识别	MER-UniBench等	将情绪轮评价指标作为奖励进行策略优化。显著提高OV-MER性能，在MER-UniBench上刷新SOTA ³	(待开源)
RL噪声增强 ¹³	RL选择噪声数据增强 (SER任务)	IEMOCAP (4类)	RL策略挑选频带噪声合成训练数据，使SER对噪声鲁棒。跨语料/跨语言测试中优于随机增强 ¹³	-

以上模型覆盖了基于Transformer、自监督、多任务学习、图神经网络与强化学习等多种技术路线，对比基线均在相关数据集上给出了显著提升²³。此外还有其他方法如LMF、TFN、MFN、DialogRNN等经典模型，此处不再赘述。总体上，多模态融合策略从早期的逐层拼接到如今的动态注意力与结构化建模不断演进，为情感分析精度提升奠定了基础。

COGMEN架构与强化学习在情感分析中的应用

COGMEN架构：COGMEN (Contextualized Graph Neural Network based Multimodal Emotion recognition) 由Joshi等人在2022年提出，用于对话中逐句情感分类²¹²。其核心模块包括：文本/语音/视觉特征编码器；对话**全局上下文编码**（通常用Transformer捕捉整段对话信息）；基于说话人关系的**图神经网络 (GraphTransformer/RGCN)** 捕捉本地上下文及说话者间依赖²。通过同时利用对话历史（全局信息）与说话者相关信息（局部信息），COGMEN显著提升了对话中每句情感的预测准确度。在IEMOCAP和MOSEI数据集上，对比以往仅使用文本或简单拼接的方法，COGMEN分别在多类情感F1上实现了领先²¹²。然而，目前关于COGMEN的研究较少，后续改进可能包括使用更强大的上下文编码器、探索图结构学习等方向。

强化学习在情感分析中的应用：强化学习 (RL) 在情感分析领域的研究仍较新颖，但近年来已有多项尝试：

- **策略梯度用于模态融合：**如GME-LSTM(A)将策略梯度RL用于训练门控模块，在多模态输入中自动选择可信特征，从而抑制噪声模态，提高了情感分类性能⁷。
- **RL用于数据增强：**Ranjan等人 (2024) 使用PPO策略选取适当的噪声频段用于语音数据增强，使得语音情感识别模型对噪声更加鲁棒，尤其在跨语料和跨语言设置下效果优于传统方法¹³。
- **RL优化大模型推理：**最新工作R1-Omni⁴、EMO-RL⁵和AffectGPT-R1³利用RL (策略梯度) 直接优化情感识别任务的评价指标。如AffectGPT-R1将情绪轮指标作为奖励进行策略优化，大幅超越基于标记训练的模型³；EMO-RL为音频-语言模型设计了带约束的策略优化，极大增强了模型的泛化和推理能力⁵。

- **对话系统中的RL**：虽非多模态融合，部分工作关注对话生成中的情感控制，如使用Actor-Critic模型调优聊天机器人回复情绪¹⁴（见上面RL评论综述）。但多模态对话情感（如CMU-MOSI式任务）中，将RL用于对话策略尚较少。

总体来看，已有方法成功将RL引入情感领域的部分环节，但也存在局限：如需要设计合理奖励函数、RL训练不稳定、样本效率低等问题^{3 5}。此外，目前多模态情感分析中绝大多数工作仍以监督学习为主，对比之下RL应用相对稀缺。未来可探索的方向包括：使用离线RL学习人类标注数据中的情感模式；对抗式训练结合RL增强鲁棒性；利用元学习快速适应新领域等。

多模态情感分析的未解决挑战

尽管已有诸多进展，多模态情感分析依然面临以下关键挑战，并严重影响模型性能和应用：

- **模态对齐问题**：不同模态的数据采样率和时序不一致导致对齐困难，如音频帧与视频帧长度往往不同¹。MuT论文指出，“数据固有的不对齐性”及模态间长距离依赖是建模难点，需特殊的交叉注意力机制才能捕捉异步信号间的关联¹。对齐误差会降低融合效果，量化指标如F1往往明显下降（如在对话情感识别中错过某些重要声音/表情线索）。
- **时序建模与长依赖**：多模态序列往往较长，单纯的RNN难以捕获长距离跨模态互动。Transformer等架构虽然有效，但计算开销大，且仍难以保证捕获长期上下文。结果之一是**记忆瓶颈**：模型容易遗忘早期信息，无法充分利用整段对话或视频。例如MuT在评测中需专门设计交叉注意力来应对这一问题¹。
- **噪声鲁棒性**：真实世界数据噪声多，如背景噪音、视频失真、传感器干扰等。数据可能存在缺失（例如人脸遮挡、声音丢帧）或质量低（低分辨率、低采样率）⁶。Ranjan等工作显示，传统方法在噪声环境下性能显著下降，而引入RL数据增强可部分缓解¹³。但总体而言，模型对噪声的敏感度仍高，需要进一步研究数据增强、鲁棒特征提取或鲁棒训练方法。
- **数据稀缺与标签偏差**：当前多模态情感数据集规模普遍偏小（数千条级别），且往往只涵盖少数语言、文化或场景¹⁵。少样本导致容易过拟合、泛化差。Yu等人（CH-SIMS）也指出，缺乏大规模多语言数据集难以训练具有强泛化性的模型¹⁵。此外，多模态数据集的**标注偏差**问题严重：情感标签本身具有主观性，且传统标注往往对三个模态给出统一一标签，忽略模态间矛盾（如文本积极但表情消极）。这可能导致模型学习到不准确的多模态关联，影响性能。
- **可解释性**：多模态融合模型结构复杂，难以解释。已有工作（如R1-Omni⁴）强调解释不同模态对预测的贡献，但大部分模型仍难以说明决策依据。缺乏透明度使得模型在敏感场景下（如情感干预、医疗）难以被信任。可视化注意力、因果推断和符号推理等是潜在解决方案，但尚未成为标配。
- **实时性与计算成本**：多模态模型往往网络复杂（如跨模态Transformer、图神经网络等）且输入维度高，推理和训练计算量大。实时情感分析系统需要低延迟和高效推理，目前此方面研究仍不足。需要对模型进行优化（轻量化、量化）或设计实时流式处理架构，但相关文献较少。
- **跨域泛化**：现有模型通常在特定数据集上训练和评测，难以直接应用于不同领域。例如在R1-Omni中发现，引入RL能使模型在**外部测试集**上表现更稳健¹⁶。但大多数情感分析系统对新的说话风格、文化背景、录音设备等泛化能力有限，需要跨域自适应技术，如领域对抗训练、元学习等。
- **隐私与伦理**：多模态数据可能包含个人敏感信息（面部影像、语音、生理信号）。如何在保护隐私的同时利用这些信号进行情感识别，是一个亟待解决的问题。目前公开文献中针对多模态情感分析的隐私保护研究很少，需要考虑联邦学习、加密计算以及伦理审查等方案。

上述挑战既有理论复杂性，也有实验限制，需要结合算法设计和数据工程共同攻关。例如，量化影响方面，已有研究发现“在无噪声条件下，情感分类准确率可达90%以上；而噪声环境下可能下降20%以上”¹³。同时，大型模型（如SpeechLLM）在情感识别任务上参数量巨大，对小数据集的过拟合风险和推理延迟必须权衡。

创新研究方向与技术思路

基于前述分析，我们提出以下至少五个具有可行性的创新研究方向，每个方向均阐述技术方案、预期改进、潜在风险与难点、所需资源和评估方案。

- 1. 基于RL的自适应模态融合：**设计一个策略学习器（RL agent），根据上下文动态分配模态权重或选择使用哪些模态。例如，在噪声环境下强化选择文本信息，在缺失视频时加强音频等。技术思路：将多模态特征融合看作决策过程，使用政策梯度（如PPO）优化情感分类的准确率作为奖励。与传统拼接或固定加权融合相比，预期能够自动忽略噪声模态，提高鲁棒性。**风险：**奖励设计困难、训练不稳定；策略偏好可能陷入局部。**资源：**所需计算量中等（可在单卡GPU上训练），数据集如CMU-MOSI/MOSEI或多模态对话（IEMOCAP/MELD）。**评估：**与固定融合基线比较分类准确率、F1；可增加不同噪声比例实验；消融实验比较有无RL模块的效果。
- 2. 对话情感中的上下文强化学习：**结合COGMEN架构，在对话情感识别中使用RL优化对话上下文的利用。技术思路：针对对话中的当前发言，将其上下文（历史发言）分为局部窗口和全局对话，让RL agent决策使用多长的上下文或强调哪些说话者的信息。例如Agent决定关注过去3句 vs 5句对话，或加权发言者贡献。**预期：**提升对话情绪分析精度，并自动适应对话长度变化。**风险：**多轮对话数据复杂，训练样本较少；状态空间大。**资源：**计算需求中等，可在对话数据集上训练，如IEMOCAP/MELD；**评估：**在对话ERC基准上比对固定上下文模型的F1；设计ablation验证不同上下文长度策略。
- 3. 多模态预训练与自监督融合：**借鉴大模型预训练的思路，构建多模态情感分析的预训练任务（如对模态置换、缺失进行预测），然后在情感任务上微调。技术思路：使用大规模无标注视频/语音数据进行多任务自监督（如对比学习、跨模态重建），学习通用情感相关表征。结合COGMEN中的图神经网络结构，在预训练时加入图预测任务。**预期：**提高模型对多模态模式的理解能力，在小数据集上更易泛化。**风险：**需要收集/使用大规模未标注数据，预训练成本高；任务设计需与情感预测相关。**资源：**需高计算（多卡GPU或TPU），利用YouTube等开放视频；**评估：**对比基线和预训练+微调模型在CMU-MOSEI、CH-SIMS等上的表现提升。
- 4. 心理生理信号的深度融合：**引入如心率(GSR)、脑电(EEG)等生理信号作为额外模态，与语音/视频/文本联合建模。技术思路：设计多模态融合网络（如多通道Transformer或混合图结构），处理标准信号流与视觉/语言序列。由于生理数据通常连续且噪声大，可借助RL自动选择有效特征频段。**预期：**生理信号提供用户真实情绪线索，可增强情感识别的精确度和鲁棒性。**风险：**获取此类数据困难且昂贵；信号噪声大（如运动伪影）；数据处理和对齐复杂。**资源：**使用公开数据集（如DEAP、AMIGOS、MuSe），低至中等计算资源；**评估：**以传统多模态（无生理）为基线，引入生理模态后对比精度；消融实验去除各模态评估贡献。
- 5. 跨域与元学习方法：**针对数据稀缺和泛化差，采用元学习或领域自适应技术。技术思路：使用元学习（Model-Agnostic Meta-Learning等）训练模型在多种环境/说话人上快速适应新域；或者基于RL的域适应策略，选择最适合目标域的源域知识。**预期：**在跨不同语言、场景或传感器条件下表现更稳健。**风险：**设计合适的元学习任务复杂；领域差异难以量化。**资源：**需要多个数据集训练（如英文与中文混合）；**评估：**跨域测试：训练在一个域（如西方访谈），测试另一域（如中文对话），对比无适应技术的基线。

- 6. **鲁棒对抗训练**：借鉴对抗学习思想，引入对抗样本训练增强模型鲁棒性。技术思路：在训练过程中生成对抗的模态输入（如对抗噪声扰动），并用强化学习评估对抗强度。**预期**：提升模型对攻击和异常情境的稳定性。**风险**：对抗训练易导致收敛困难；不当策略可能损坏正常准确率；**资源**：中等GPU计算，使用已有人脸/语音对抗生成器；**评估**：在对抗测试集上比较鲁棒性，并与未对抗训练模型进行对比。

以上方向均结合COGMEN和RL思想，如使用RL进行融合和结构优化，同时也可独立尝试其他新范式（如元学习、自监督、融合生理信号）。每个方向的**评估方案**通常包括：使用公开多模态情感基准数据（CMU-MOSI/MOSEI、IEMOCAP/MELD、CH-SIMS等），与当前SOTA基线比较（准确率、F1、MAE/相关性等指标），进行消融实验以验证新模块的有效性，并在不同环境（噪声、跨域）下检验鲁棒性。

实验方案设计

为了验证上述创新，我们设计了两个可实施的实验方案如下，并用流程图说明训练与评估流程。

方案1：RL驱动模态融合实验

- **模型架构**：输入层提取文本、视觉、语音特征（如BERT编码文本，ResNet/CNN处理视频帧，LSTM处理音频）。接着进入融合网络（如多模态Transformer或拼接+全连接）。在此基础上并联一个**RL决策模块**：该模块根据当前批次特征或先前预测的信心（State）输出动作，即各模态的融合权重或是否使用某模态（Action）。最终融合后的特征送入情感分类器输出。
- **训练流程**：如图所示（下图），采用交替训练：常规地用标注数据训练情感分类器；同时根据分类准确率等指标计算奖励，更新RL策略（例如使用PPO算法）。损失函数包含分类交叉熵和策略梯度损失。
- **奖励设计**：可以设定实时奖励，如一个批次正确率的平均值；或更精细的情感极性相关指标。也可引入次要奖励以鼓励使用更多模态信息或惩罚过度依赖单一模态。
- **超参数**：学习率（文本0.0001、RL 0.0005）、策略探索率（如 ϵ -greedy概率0.1）、奖励系数权重平衡等。批次大小依据设备设定（如32或64）。
- **消融实验**：比较添加RL模块前后性能；单一模态输入与多模态输入；不同奖励策略的效果。

下面给出简化的mermaid流程图：

```

flowchart TB
    subgraph "多模态情感模型"
        T[文本特征编码] --> M[Fusion 层];
        A[语音特征编码] --> M;
        V[视觉特征编码] --> M;
        M --> C[情感分类器];
    end
    subgraph "强化学习模块"
        C --> R[奖励计算];
        R --> Agent[RL策略网络];
        Agent --> |输出模态权重| M;
    end
  
```

方案2：对话情感中的RL图结构优化

- **模型架构**：以COGMEN为基础，输入为对话中一系列多模态语句。首先用Transformer或BiLSTM提取每句话的多模态嵌入。然后构建对话图：每个发言作为节点，节点间初始连接依据时间顺序或说话人关系。引入**可学习边**（Edge）或**可调路径**：RL agent决定连接哪些上下文节点，或者对局部/全局上下文进行权重分配。图神经网络（RGCN/GraphTransformer）汇聚后送入分类器。
- **训练流程**：固定节点表示提取器，RL agent输入对话状态（如节点特征差距、前一轮预测准确度等），输出图结构调整（例如保留/删除某条边）。用ERC任务的F1作为奖励信号优化RL策略。训练同时更新图神经网络参数和RL策略。
- **评估**：在对话情感基准（IEMOCAP、MELD）上对比固定图结构模型。分析RL学习到的连接模式是否符合对话逻辑。

这两个实验方案可在标准硬件（单机多GPU）上实施。所用数据集包括公开的多模态对话情感集（IEMOCAP/MELD）和视频情感集（CMU-MOSI/MOSEI、CH-SIMS）。评估指标以准确率、加权F1、MAE/相关性等为主。基线包括无RL的原始模型、简单拼接融合等。消融实验则检验RL模块对不同模态或上下文长度的影响。

研究路线图

下图展示了本课题的**短期（1-2年）**、**中期（3-5年）**与**长期（5年以上）**研究计划的甘特图示例。我们建议逐步推进，从重现已有基线与小规模验证开始，再扩展到大规模预训练和跨域应用。

```
gantt
dateFormat YYYY-MM-DD
title 多模态情感分析研究路线图
section 短期（1-2年）
数据准备与基线验证           :done,    des1, 2026-06-01, 2026-12-31
RL融合模型初步实验         :active,  des2, 2026-07-01, 2027-06-30
COGMEN结构复现与优化:       des3, 2026-10-01, 2027-09-30
section 中期（3-5年）
大规模预训练与微调         :         des4, 2027-01-01, 2029-12-31
跨域/多语种适应           :         des5, 2028-01-01, 2030-12-31
鲁棒性对抗与解释研究       :         des6, 2028-06-01, 2031-12-31
section 长期（5年以上）
实时系统与硬件部署         :         des7, 2031-01-01, 2034-12-31
行业应用验证               :         des8, 2032-01-01, 2035-12-31
```

- **短期任务**：使用CMU-MOSI/MOSEI、IEMOCAP/MELD、CH-SIMS等数据集验证基线与提出模型，搭建实验平台。复制经典SOTA结果，并初步集成RL模块（如基于GME-LSTM或COGMEN的简单RL融合）。
- **中期任务**：扩展到大规模未标注预训练（使用YouTube视频数据等），以及跨数据集/语言测试。研究对抗训练、元学习提升泛化，探索算法加速与模型压缩。
- **长期任务**：探索实时部署与应用场景，如情感对话机器人或健康监护系统。评估系统级性能，确保隐私与伦理合规。

开源工具与推荐数据集：

- 深度学习框架：PyTorch、TensorFlow/Keras。图神经网络库：DGL、PyG。RL库：OpenAI Gym、Ray RLlib。

- 多模态工具：HuggingFace Transformers（文本）、Librosa（音频特征）、OpenFace/MediaPipe（视频关键点）、OpenSMILE（音频特征）。
- 推荐数据集：CMU-MOSI/MOSEI、IEMOCAP、MELD、CH-SIMS、DEAP（情感脑电/生理）、MuSe（多模态情绪）、PELD（多传感器）、MuSe-Text（跨媒体情感）等。
- 数据预处理：统一音频采样率（如16kHz）、静音截断；视频对齐关键帧、图像归一化；文本分词（包括中英文混合）、去除停用词；生理信号滤波去噪（带通滤波器）、同步对齐。

以上资源均支持可复现研究。例如，[CH-SIMS](#)项目提供公开的多模态中文情感数据与基线代码¹⁷。国内外多模态情感分析工具箱（如Multimodal MultiSense）和RL平台（如Stable Baselines3）也可辅助实现。

主要参考文献

- Chen M., Wang S. 等. Multimodal Sentiment Analysis with Word-Level Fusion and Reinforcement Learning. ICMI 2017⁷ .
- Tsai Y.-H., Bai S. 等. Multimodal Transformer for Unaligned Multimodal Language Sequences. ACL 2019¹ .
- Hazarika D., Zimmermann R. 等. MISA: Modality-Invariant and -Specific Representations for Multimodal Sentiment Analysis. ACM MM 2020⁸ .
- Yu W., Xu H. 等. Learning Modality-Specific Representations with Self-Supervised Multi-Task Learning for Multimodal Sentiment Analysis. AACL 2021⁹ .
- Hu J., Liu Y. 等. MMGCN: Multimodal Fusion via Deep Graph Convolution Network for Emotion Recognition in Conversation. ACL 2021¹⁰¹¹ .
- Joshi B. K., Zhong Q. 等. COGMEN: COntextualized GNN based Multimodal Emotion recognition. NAACL 2022²¹² .
- Ranjan S., Chakraborty R. 等. Reinforcement Learning based Data Augmentation for Noise Robust Speech Emotion Recognition. Interspeech 2024¹³ .
- Zhao J., Wei X. 等. R1-Omni: Explainable Omni-Multimodal Emotion Recognition with Reinforcement Learning. arXiv 2025⁴ .
- Lian Z., Zhang F. 等. AffectGPT-R1: Leveraging Reinforcement Learning for Open-Vocabulary Multimodal Emotion Recognition. arXiv 2025³ .
- Li P., Zhao B. 等. EMO-RL: Emotion-Rule-Based Reinforcement Learning Enhanced Audio-Language Model for Generalized Speech Emotion Recognition. EMNLP Findings 2025⁵ .
- Lai S., Hu X. 等. Multimodal Sentiment Analysis: A Survey. arXiv 2023¹⁸⁶ .
- Yu W., Xu H. 等. CH-SIMS: A Chinese Multimodal Sentiment Analysis Dataset with Fine-grained Annotation of Modality. ACL 2020¹⁷ .
- Ranjan S., Chakraborty R. 等. Reinforcement Learning based Data Augmentation for Noise Robust Speech Emotion Recognition. Interspeech 2024¹³ .

(以上文献均为原始论文或官方资料，图表和代码均引用相应来源。)

¹ Multimodal Transformer for Unaligned Multimodal Language Sequences - ACL Anthology

<https://aclanthology.org/P19-1656/>

² ¹² aclanthology.org

<https://aclanthology.org/2022.naacl-main.306.pdf>

³ [2508.01318] AffectGPT-R1: Leveraging Reinforcement Learning for Open-Vocabulary Multimodal Emotion Recognition

<https://arxiv.org/abs/2508.01318>

- 4 16 [2503.05379] R1-Omni: Explainable Omni-Multimodal Emotion Recognition with Reinforcement Learning
<https://arxiv.org/abs/2503.05379>
- 5 [2509.15654] EMO-RL: Emotion-Rule-Based Reinforcement Learning Enhanced Audio-Language Model for Generalized Speech Emotion Recognition
<https://arxiv.org/abs/2509.15654>
- 6 15 18 [2305.07611] Multimodal Sentiment Analysis: A Survey
<https://arxiv.labs.arxiv.org/html/2305.07611>
- 7 [1802.00924] Multimodal Sentiment Analysis with Word-Level Fusion and Reinforcement Learning
<https://arxiv.org/abs/1802.00924>
- 8 [2005.03545] MISA: Modality-Invariant and -Specific Representations for Multimodal Sentiment Analysis
<https://arxiv.org/abs/2005.03545>
- 9 [2102.04830] Learning Modality-Specific Representations with Self-Supervised Multi-Task Learning for Multimodal Sentiment Analysis
<https://arxiv.org/abs/2102.04830>
- 10 11 MMGCN: Multimodal Fusion via Deep Graph Convolution Network for Emotion Recognition in Conversation
<https://aclanthology.org/2021.acl-long.440.pdf>
- 13 [isca-archive.org](https://www.isca-archive.org)
https://www.isca-archive.org/interspeech_2024/ranjan24_interspeech.pdf
- 14 Reinforcement learning in sentiment analysis: a review and future directions | Artificial Intelligence Review | Springer Nature Link
<https://link.springer.com/article/10.1007/s10462-024-10967-0>
- 17 CH-SIMS: A Chinese Multimodal Sentiment Analysis Dataset with Fine-grained Annotation of Modality - ACL Anthology
<https://aclanthology.org/2020.acl-main.343/>