

Paper Reading Report

1 论文背景与核心问题

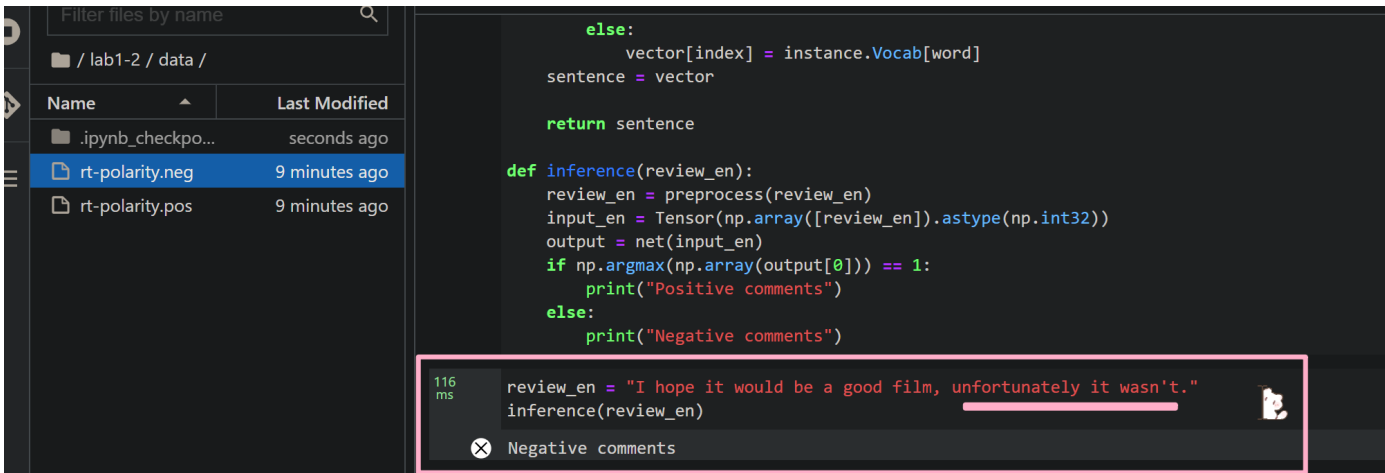
我选择阅读的论文为发表在第32届国际多媒体会议（[dblp: MM 2024](#)）上的《Dual-path Collaborative Generation Network for Emotional Video Captioning》，其文章链接为：[3664647.3681603](#)。

1.1 背景介绍

随着抖音、推特等自媒体视频平台的兴起，越来越多人选择在网络上用视频来表达自己的观点。而表达就必然会传递情感取向，因此，在线上视频数量增多的当下，基于情感的视频理解任务越来越受到关注，EVC（Emotional Video Captioning，情感视频字幕）就是其中之一。这是多媒体社区中一个新兴话题，不仅涉及到对真实视频内容的理解，还要求识别出视频中包含的复杂感情脉络，并结合情景生成字幕。

1.2 核心问题

传统自然语言处理（NLP）领域用于情感分类的模型往往忽视了情境中人物的情感动态变化，而给出一个直白的结果。比如我在NLP课程上的实验所训练出的情感分类模型（根据人物评价分析对电影的喜恶，见下图），它并不能分析出人物从期待到失望的心理变化，而仅仅能根据浮于表面的词汇得出非黑即白的评价。



传统的视频字幕同样存在这样的问题，这种方法忽视了视频内在情感的动态变化，在用于视频情感分析时，和现实场景中的真实叙事存在较大的割裂。另一方面，传统视频字幕将情感线索纳入分析的每个步骤，这放大了情感的指导作用，也在一定程度上忽略了事实内容在生成描述时的参考作用。更坏的情况下，当传统的EVC方法预测得到错误的全局情绪，在该情绪的指导下，很可能生成和视频事实内容无关的错误情感描述。

为了填补这方面的研究空白，研究团队提出了一种新的双路径协作生成网络（Dual-path Collaborative Generation Network）用于情感视频字幕。两条路径分别为动态情感感知路径和自适应字幕生成路径，前者感知每个步骤中的情绪演变，后者则在感知正确情绪的前提下自适应地生成与情绪相关的单词。它在生成情感字幕的同时动态感知情绪线索的演变，提出的两条路径通过协作学习相互促进，促进了情感字幕生成性能。

2 研究方法与创新点

2.1 研究方法

该研究提出的双路径协作生成网络（DCGN）通过两条互补路径解决“情感动态性”与“事实-情感平衡”这两大问题。

在**动态情感感知路径**中，模型首先提取视频的视觉特征（如人物动作、场景变化），并结合已生成的字幕历史信息，逐步调整情感特征：

1. 元素级筛选：根据当前生成的字幕内容（例如“流泪”），从视觉特征中过滤出相关的情感线索（如“湿润的眼睛”或“低垂的嘴角”），强化与当前情感匹配的特征。
2. 子空间分解：将情感特征拆分为不同的语义层次（如“肢体动作”“面部表情”“背景音乐”等），动态调整各部分的权重占比。例如，当视频从“紧张对话”转向“轻松拥抱”时，模型会增强肢体动作子空间的权重，捕捉肢体放松的变化。

在**自适应字幕生成路径**中，模型通过对比已生成字幕与视频内容的关联性，动态决定每一步的情感介入程度：

1. 情感强度估计：计算当前生成状态（如已生成的“颤抖的手”）与视频特征的匹配度。若匹配度高（如视频中确实存在“手部颤抖”），则提高情感权重，生成更强烈的情绪词（如“恐惧地”）；若匹配度低，则降低情感干预，优先描述事实（如“拿着杯子”）。
2. 特征融合控制：将动态调整后的情感特征与原始视觉特征按权重融合。例如，生成“她笑着说”时，情感特征占主导；生成“站在讲台前”时，更多依赖物体识别等事实特征。

两条路径通过共享中间状态（如当前生成词的特征）实现协作：情感路径提供细粒度情感演化信号，生成路径反馈实际用词效果，二者共同优化最终输出，实现双路径协作生成情感视频字幕。

2.2 创新点

自适应字幕生成路径的**核心在于动态调节情感与事实的融合比例**。模型通过对比已生成文本与视频特征的关联性，计算每一步的情感强度权重。例如，生成形容词“激动地”时，情感权重较高，模型更依赖情感特征；而生成名词“奖杯”时，情感权重降低，更多参考事实性视觉特征。两条路径通过共享中间状态（如当前生成词的表征）实现协同——情感感知路径提供动态调整后的情感特征，生成路径则根据这些特征调整输出，同时将生成结果反馈至情感路径以修正后续步骤。这种设计**避免了传统方法中情感与事实的“硬性绑定”**，使模型能灵活适应不同场景需求。

3 实验结果与局限性

3.1 实验结果

实验在EVC-MSVD、EVC-VE和EVC-Combined等三个公开数据集上验证了模型效果，分别得到如下结果：

情感准确性显著提升。以EVC-VE数据集为例，情感词准确率（ Acc_{sv} ）从基线模型的63.8%提升至71.0%，情感句准确率（ Acc_c ）从62.3%提升至69.4%。这表明动态情感感知路径能更精准捕捉情绪变化，例如在“运动员带伤完赛”视频中，模型能识别从“痛苦”到“释然”的过渡，而非简单标注为“励志”。

语义质量明显改善。CIDEr分数（衡量生成字幕与人工标注的相似度）在EVC-VE上提升19%，尤其在情感转折明显的场景（如“争吵后和解”），生成描述更贴近真实情感演变。

消融实验验证了双路径协作的必要性。单独使用动态情感路径时CIDEr仅提升4.9%，单独使用自适应生成路径提升6.8%，而两者联合后提升幅度扩大至19%，说明协同效应远大于单一模块的贡献。

总体来看，双路径协作生成网络（DCGN）在情感准确性（ Acc ）和语义质量（如CIDEr）上均优于基线模型，且双路径协作机制是该网络的关键点之一，协同效应显著。

3.2 局限性

尽管模型性能优于传统方法，但仍存有明显的不足之处：

首先是模型**高度依赖CLIP等预训练视觉模型**。当视频质量较差（如模糊、低光照）时，视觉特征提取误差会传导至情感感知路径，导致误判。例如，昏暗场景中的“疲惫神情”可能被错误关联到“平静”，进而影响对情感的感知判断。

再者，模型**对混合情感的处理能力有限**。在“笑着流泪”“尴尬沉默”等多情感特征复合的复杂场景中，模型往往只能捕捉单一情绪（如“悲伤”或“快乐”），而无法表达矛盾情感的交织，这极大地减弱了模型的情感分析能力。

最后，双路径方式下情感视频字幕的**生成效率较低**。由于动态情感路径需多步迭代调整，单句生成时间比传统方法增加35%，难以满足直播字幕等实时性要求高的场景需求。

4 批判性思考

从技术改进的角度看，该研究仍存有一定的优化空间。其一，情感感知的输入信息过于依赖视觉模态，而忽略了其他模态的情感表达。例如，在某段用于在婚礼现场播放的剪辑视频中，背景音乐的快慢、歌词内容往往与情感起伏直接相关，但该论文提出的模型并未融合音频特征，便存在将“舒缓音乐下的离别场景”误判为“温馨场景”的可能。在这一点上，未来可以考虑**引入多模态融合机制**，结合音频、文本评论等多维度信号进行视频情感的判断，提升情感判别的鲁棒性。另外则是，动态情感路径的计算复杂度较高，限制了模型的实际应用（这一点在上述模型的局限性分析里也有提到）。可以尝试采用稀疏注意力机制或分层迭代策略，仅对关键帧进行细粒度情感分析，有可能在不损失精度的情况下**减少计算耗时，降低计算复杂度**。

而在考虑将该双路径协作生成网络的模型应用到实际中时，不难看出该技术可延伸至多个领域。例如，在心理健康医疗领域，可以用其分析患者自述视频的情感波动（比如从“焦虑”逐渐变得“平静”），辅助医生评估治疗进展的同时减少人力成本。而在市场庞大的影视创作领域，该模型或能自动为未剪辑素材添加情感标注（如“冲突爆发-情绪高潮-和解缓和”），提升后期剪辑效率；同时，对于已形成并上传到社交媒体上的视频，该模型也能够依据情感对其进行分类，方便社交媒体的推流等操作。

总体而言，该论文提出的双路径协作生成网络（DCGN）为情感视频理解与字幕标注提供了新思路，但在实际落地与正式应用中仍需解决效率与泛化性等问题。