



计算机科学

Computer Science

ISSN 1002-137X, CN 50-1075/TP

《计算机科学》网络首发论文

题目：基于多模态双协同 Gather Transformer 网络的虚假信息检测方法
作者：向旺，王金光，王一飞，钱胜胜
收稿日期：2023-10-10
网络首发日期：2024-03-15
引用格式：向旺，王金光，王一飞，钱胜胜. 基于多模态双协同 Gather Transformer 网络的虚假信息检测方法[J/OL]. 计算机科学.
<https://link.cnki.net/urlid/50.1075.TP.20240314.1758.006>



网络首发：在编辑部工作流程中，稿件从录用到出版要经历录用定稿、排版定稿、整期汇编定稿等阶段。录用定稿指内容已经确定，且通过同行评议、主编终审同意刊用的稿件。排版定稿指录用定稿按照期刊特定版式（包括网络呈现版式）排版后的稿件，可暂不确定出版年、卷、期和页码。整期汇编定稿指出版年、卷、期、页码均已确定的印刷或数字出版的整期汇编稿件。录用定稿网络首发稿件内容必须符合《出版管理条例》和《期刊出版管理规定》的有关规定；学术研究成果具有创新性、科学性和先进性，符合编辑部对刊文的录用要求，不存在学术不端行为及其他侵权行为；稿件内容应基本符合国家有关书刊编辑、出版的技术标准，正确使用和统一规范语言文字、符号、数字、外文字母、法定计量单位及地图标注等。为确保录用定稿网络首发的严肃性，录用定稿一经发布，不得修改论文题目、作者、机构名称和学术内容，只可基于编辑规范进行少量文字的修改。

出版确认：纸质期刊编辑部通过与《中国学术期刊（光盘版）》电子杂志社有限公司签约，在《中国学术期刊（网络版）》出版传播平台上创办与纸质期刊内容一致的网络版，以单篇或整期出版形式，在印刷出版之前刊发论文的录用定稿、排版定稿、整期汇编定稿。因为《中国学术期刊（网络版）》是国家新闻出版广电总局批准的网络连续型出版物（ISSN 2096-4188，CN 11-6037/Z），所以签约期刊的网络版上网络首发论文视为正式出版。

基于多模态双协同 Gather Transformer 网络的虚假信息检测方法

向旺¹ 王金光² 王一飞¹ 钱胜胜³

1 郑州大学 河南先进技术研究院 郑州 450000
2 合肥工业大学 计算机与信息学院 安徽 230601
3 中国科学院自动化研究所多模态人工智能系统全国重点实验室 北京 100190
(1584462772@qq.com)

摘要 社交媒体网站是人们在日常生活中分享信息、表达和交换意见的便捷平台。随着用户数量的不断增加，社交媒体网站上出现了大量的信息数据。然而，由于用户没有检查共享信息的可靠性，这些信息的真实性难以保证。这种情况导致了大量虚假信息在社交媒体上广泛传播。然而，现存的方法大多存在以下局限性：1) 现有的方法大多数通过简单提取文本与视觉特征，将其拼接后得到多模态特征来进行虚假信息判断，忽略了模态间和模态内细粒度内在联系，缺乏对关键信息的检索和筛选；2) 多模态信息间缺乏指导性的特征提取，文本和视觉等特征之间缺乏交互增强，对多模态信息理解不足。为了应对这些挑战，提出了一种新颖的基于多模态双协同 gather transformer 网络 (Multimodal Dual-Collaborative Gather Transformer Network, MDCGTN) 来解决现有方法的局限性。在 MDCGTN 模型中，通过文本-视觉编码网络对文本和视觉信息的特征表示进行提取，将获得的视觉和文本特征表示输入多模态 gather transformer 网络进行多模态信息融合，使用 gather 机制提取关键信息，充分捕捉和融合模态内和模态间细粒度关系。除此之外，设计了一个双协同机制对社交媒体帖子多模态信息进行整合，以此实现模态之间信息的交互和增强。文章在两个公开可用的基准数据集上进行了大量实验，与现有的先进基准方法相比，准确率有明显提升，证明了 MDCGTN 方法对于虚假信息检测的优越性能。

关键词： 多模态；虚假信息检测；社交媒体；gather transformer 网络

Multi-modal Dual Collaborative Gather Transformer Network for Fake News Detection

XIANG Wang¹, WANG Jin-guang², WANG Yi-fei¹ and QIAN Sheng-sheng³

1 Henan Institute of Advanced Technology, Zhengzhou University, Zhengzhou 450000, China
2 School of Computer Science and Information Engineering, Hefei University of Technology, Hefei 230601, China
3 State Key Laboratory of Multimodal Artificial Intelligence Systems, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China

Abstract Social media platforms are convenient platforms for people to share information, express opinions, and exchange ideas in their daily lives. With the increasing number of users, a large amount of data has emerged on social media websites. However, the reliability of the shared information is difficult to guarantee due to users' lack of verification. This situation has led to the widespread dissemination of a large amount of fake news on social media. However, existing methods suffer from the following limitations: 1) Most existing methods rely on simple text and visual feature extraction, concatenating them to obtain multimodal features for detecting fake news, while ignoring the fine-grained intrinsic connections within and between modalities, and lacking retrieval and filtering of key information; 2) There is a lack of guided feature extraction among multimodal information, with insufficient interaction and understanding between textual and visual features. To address these challenges, a novel Multimodal Dual-Collaborative Gather Transformer Network (MDCGTN) is proposed to overcome the limitations of existing methods. In the MDCGTN model, textual and visual features are extracted using a text-visual encoding network, and the obtained features are input into a multimodal gather transformer network for multimodal information fusion, using the gather mechanism to extract key information, fully capturing and fusing fine-grained relationships within and between modalities. In addition, a dual-collaborative mechanism is designed to integrate multimodal information of social media posts, thereby achieving interaction and enhancement of information between modalities. Extensive experiments are conducted on two publicly available benchmark datasets, and the results show that the proposed MDCGTN method significantly outperforms existing state-of-the-art methods in terms of accuracy, proving its superiority for fake news detection.

到稿日期：2023-10-10 返修日期：2024-3-1
基金项目：国家自然科学基金 (62276257)
This work was supported by the National Natural Science Foundation of China (62276257).
通信作者：钱胜胜 (shengsheng.qian@nlpr.ia.ac.cn)

information fusion. The gathering mechanism is used to extract key information, fully capturing and fusing fine-grained relationships within and between modalities. In addition, a dual-collaborative mechanism is designed to integrate multimodal information in social media posts, enhancing interaction and understanding between modalities. Extensive experiments are conducted on two publicly available benchmark datasets. Compared to existing state-of-the-art benchmark methods, the proposed MDCGTN method achieves significant improvement in accuracy, demonstrating its superior performance in detecting fake news.

Keywords Multi-modal, Fake news detection, Social media, Gather transformer network

1 引言

近年来,随着社交媒体平台的广泛普及和便利性的增加,越来越多的人加入了在线新闻的发布和使用行列。这一趋势的迅速发展可以归因于社交媒体的普及,以及使用便捷性的提升。由于用户群体的不断壮大,社交媒体网站已经成为各种信息的集散地。然而,由于用户往往忽略对其分享内容的可靠性审查,导致发布信息的真实性变得不可靠。这种核实上的松懈为大量虚假信息的猖獗传播铺平了道路。虚假信息恶意地扭曲和捏造事实信息,这种现象对个人和整个社会都产生了有害的影响。例如,2023年三月,“四川德阳‘中江机场建设成功’”的虚假信息广泛传播于社交媒体,造成了较大社会讨论和议论,占据了大量公众媒体资源。因此,为了避免不必要的损失和负面影响,我们需要开发方法来帮助识别虚假新闻,以使读者能够获得真实的信息。

近年来,虚假信息检测一直是研究关注的焦点。早期的方法主要依赖于用户报告和专家确认,然而这些方法非常耗时且劳动密集。因此,研究人员开始探索检测自动化的方法来解决该问题,且可以广泛分为两种类型:(1)第一种方法^{[1]-[4]}是基于手工制作特征的方法,这些方法通过从帖子内容和用户社交背景中提取特征,使用支持向量机(SVM)^{[1][4]}、决策树分类器^{[2][3]}等机器学习算法实现虚假信息检测。然而,手工制作的特征可能无法完全捕捉虚假信息复杂内容,存在一定局限性。(2)第二种基于深度学习的方法,利

用神经网络来捕捉深层次的特征。例如, Ma^[5]等人利用循环神经网络(RNNs)从帖子中提取潜在的特征表示,而 Yu^[6]等人则利用卷积神经网络(CNNs)进行特征提取,并辨别它们在虚假信息内容中的高阶相互作用。近年来,社交媒体内容构成已经从单纯的文本扩展到包括图片和视频等多媒体格式^{[7]-[10]},这些格式为识别虚假信息提供了补充信息。此外,用户评论也包含了确定帖子内容真实性的关键线索。然而,上述大部分方法仅集中在文本内容中,忽视了具有多模态信息的社交媒体帖子。为解决该问题,基于深度神经网络的多模态方法被提出。虽然这些方法在虚假信息检测任务中发挥不错效果,但大部分方法仍无法充分考虑多模态信息间和信息内的细粒度内在联系,缺乏对关键信息的检索和筛选^{[11]-[13]}。

此外,社交媒体帖子对应的用户评论同样可以提供价值的线索和补充信息,进而提高帖子内容真实性的检测^{[9][14][15]}。例如, Lin^[16]等人提出了一种基于无向交互图的创新的声明导向分层图注意力网络,通过考虑包括声明和相关响应评论在内的全面社交上下文信息,增强了对响应帖子特征表示的学习。Shang^[18]等人则设计了名为 Duo-Generative Explainable Misinformation Detection (DGE Explain) 的模型,它能巧妙利用用户评论定位和阐明相关新闻中错误信息。然而,这些方法往往缺乏对用户评论的筛选和利用,不同模态信息之间缺乏相互指导、彼此交互增强,导致对多模态信息理解不足。

为了构建一个更为有效虚假信息检测方法, 需要解决以下挑战:

(1) 如何探索和捕捉多模态信息间和信息内细粒度关系, 并对关键信息进行有效检索和提取?

(2) 如何提高帖子中不同模态的指导性特征提取, 提高多模态信息间彼此交互增强的能力, 进而提高对多模态信息理解?

针对上述的局限性, 我们提出了一种新颖的基于多模态双协同 gather transformer 网络 (MDCGTN) 的虚假信息检测方法, 该网络模型充分融合了文本和视觉信息。(1) 针对挑战 1, 我们引入了 gather transformer 网络, 通过索引和选择检测信息中的关键内容, 提高模型多模态信息关键特征提取(2) 针对挑战 2, 我们引入双协同机制, 增强文本与视觉、评论信息间彼此指导关系, 提高特征间交互增强能力, 提高对多模态信息的理解。在得到社交媒体帖子的统一表示之后, 我们利用具有相应激活函数的全连接层来对帖子真实性进行分类。

总而言之, 我们的工作有如下几个贡献点:

1) 我们通过添加一个 gather 机制来改进传统 Transformer 网络结构, 充分探索和捕捉多模态信息间和信息内的细粒度内在联系, 有效筛选多模态信息中关键部分; 2) 我们提出了一个双协同机制, 以增加多模态信息间的彼此指导和协同, 进而提高多模态信息间的交互增强能力, 提高对多模态信息理解能力; 3) 在 MMCoVar^[17] 和 ReCOvery^[18] 两个公共基准数据集上的广泛实验结果表明, 我们所提出的模型比目前最先进的虚假信息检测方法表现更好。

2 相关工作

2.1 基于多模态的虚假信息检测

一些学者利用机器学习技术, 通过视觉特征提取器和文本特征提取器提取图像和文本的特征,

将它们拼接起来, 用于虚假信息检测任务。例如, 2019 年提出的 Spotfake 方法利用 VGG19 提取视觉信息和 BERT 提取文本信息, 将它们拼接起来并输入到分类器中进行虚假信息分类。Spotfake⁺^[19]方法则采用 VGG 和 XLNET 提取多模态特征。MSRD^[20]方法则考虑了帖子的图片中包含的文本信息, 使用 LSTM 建模文本信息以及图像中的文本信息, 并使用 VGG 建模视觉信息, 将多模态信息进行拼接, 得到最终的多模态特征表示。由于直接将多模态信息拼接的方法过于简单, 无法充分利用多模态信息。因此, 一些学者设计了辅助任务, 以帮助模型更好地理解多模态信息。例如, EANN^[21]方法采用 VGG 提取视觉特征和 Text-CNN^[22]提取文本特征, 将它们拼接起来得到帖子的特征表示, 并利用事件鉴别器将拼接的多模态信息作为输入, 输出事件类别用于辅助判断。

Khattar^[14]等人提出的 MVAE 同样采用 VGG 提取图像特征和双向 LSTM 提取文本特征, 并设计了信息重构的辅助任务, 以提高模型利用多模态信息的效率。另一些学者认为, 社交媒体帖子的图片内容和文本内容是否相符, 是判断该帖子是否为虚假信息的重要指标。基于该假设, 学者们提出了另一种多模态检测方法, 该方法可以检测帖子图文相符性, 进而判断帖子的真实性。具体方法是将社交媒体帖子的图片信息和文本信息分别编码, 并通过计算它们的相似度来判断它们是否匹配。如果相似度很高, 则说明该帖子的文本信息和视觉信息匹配, 即为真实信息; 如果相似度很低, 则说明该帖子的文本信息和视觉信息不匹配, 即为虚假信息。Zhou^[23]等人提出了一种基于以上方式的多模态检测方法, 该方法利用 image2text^[24]模型将图片信息转化为文本信息, 并通过全连接层将文本信息和图像信息映射到同一向量空间中。最终, 该方法通过比较两者之间

的相似度来判断信息的真实性。Xue^[25]等人则提出了另一种方法,该方法分别使用BERT和ResNet来提取文本和图像的特征,并计算两者之间的相似度,最终也是通过判断相似度进而判断帖子的真实性。

然而,上述方法仅简单将文本与图像特征进行提取和拼接,忽略了模态间和模态内的复杂关系。多模态信息间缺乏彼此指导和交互,导致模型对模态信息理解不足。

2.2 Attention 机制

Attention 机制已经广泛被应用在各种任务中,如图像字幕生成^[26]、机器翻译^[27]和推荐系统^{[28][29]}等。最初由Bahdanau等人^[27]引入,Attention 机制被应用于机器翻译任务中,能在预测目标词时关注句子中相关部分。随后,Transformer^[30]模型被设计出来,它利用Attention 机制作为LSTM的替代品来解决Sequence-to-Sequence的问题,并在机器翻译任务中取得前所未有的效果。Baeovski和Auli^[31]提出了自适应单词输入表示,使模型能够为常见单词分配更多空间,同时限制不常见单词的空间。Dai^[32]等人提出了Transformer-XL模型,它具有相对位置编码和缓存内容,增强了模型扩展上下文信息建模能力。在此基础上,Rae^[33]等人将Transformer-XL内存段扩展到更为复杂的压缩内存中,进一步扩展上下文长度,并在WikiText-103基础测试中达到了较高效果。尽

管以上方法都表明更长的上下文信息更有利于语言建模任务,但利用密集点数据生成增强上下文信息表示仍未得到充分探索。同时,经过大型文本语料库训练的大规模神经语言模型经过各种实验证明其显著优势。OpenAI的GPT^[34]和BERT^[35]模型分别采用自回归语言建模任务和掩码语言建模任务进行训练。最近,一些学者已开始将Attention 机制应用于虚假信息检测任务中。例如,Chen等人^[36]建议使用深度注意力模型,基于循环神经网络(RNN)来有选择地学习序列帖子的时间隐藏表示,从而帮助识别虚假信息。此外,还引入了一种复杂的分层多模态注意力网络^[15],它采用两个Transformer单元来联合建模多模态上下文数据,用于虚假信息检测。

受Attention 机制成功应用的启发,我们引入了一种多模态Gather Transformer网络,以合并多模态特征,全面捕捉模态之间细粒度关系,同时筛选和提取模态中关键信息。

3 MDCGTN 模型方法

3.1 任务描述

虚假信息检测任务可以被定义为一个二元分类问题,主要关注社交媒体上的帖子是否为虚假信息。给定社交媒体上由文本信息和相应图片组成的多模式帖子 P ,该模型将输出 $Y = \{0,1\}$ 来表示帖子的标签,其中 $Y = 0$ 和 $Y = 1$ 分别表示帖子是真实信息和虚假信息。

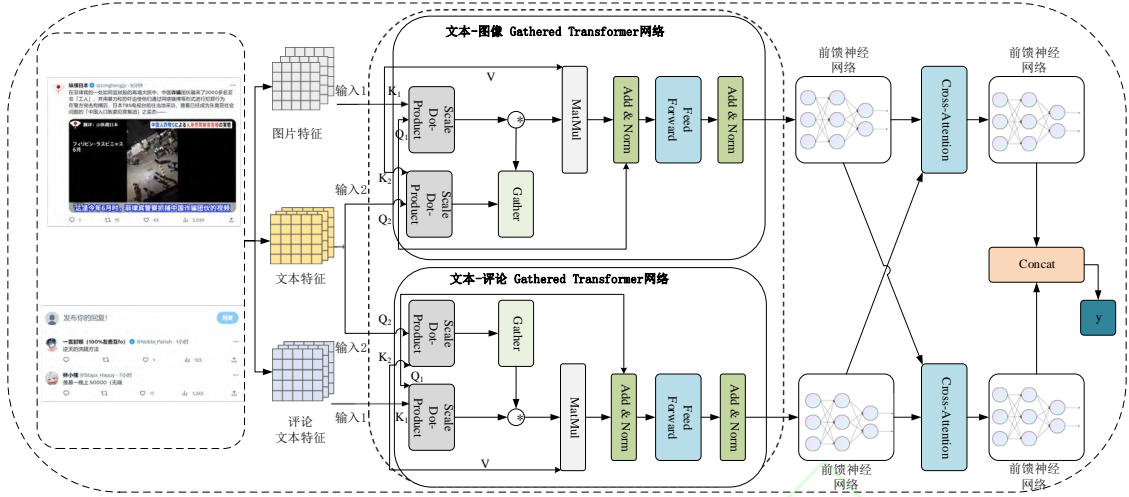


图 1 MDCGTN 模型

Fig.1 MDCGTN model

3.2 整体框架

图 1 展示了我们模型的总体框架。我们在模型中引入了多模态双协同 gather transformer 网络，通过充分融合和协同文本、图像与评论间的信息，并有效筛选模态信息中关键信息来提高虚假信息检测任务的性能。整个模型由以下组件构成：

(1) 文本-图像编码网络：对于给定帖子的文本和图像信息，我们分别使用BERT和ResNet提取文本内容和视觉内容的特征。

(2) 多模态 gather transformer 网络：由于不同模态信息间具有细腻内在联系，我们提出了一种双协同 gather transformer 网络，用于充分融合和协同多模态细粒度特征的有效方法。

(3) 双协同机制：由前馈神经网络和Cross-attention 网络构成，提高文本、视觉和用户评论间相互指导关系，提高多模态信息间互动增强能力。

(4) 虚假信息检测器：虚假信息检测器旨在将社交媒体帖子分为真实和虚假。它采用全连接层，并配有相应的激活函数，通过生成一个预测概率，最终确定帖子内容真实性。

3.3 文本-图像编码网络

正如任务描述中所陈述，我们的模型输入为多模态新闻 $P = \{T, V, C\}$ 。其中， T 表示文本内容， V 表示视觉内容， C 表示对应的用户评论。为了更好的获得社交媒体帖子细粒度的图像和文本特征表示。对于文本信息，我们使用预训练模型BERT来进行文本特征提取；对于图像信息，我们使用ResNet50来获取图像特征表示。

3.3.1 文本编码网络

为了精确地获取文本语义和上下文信息，我们采用BERT作为文本处理模型的核心模块。BERT已被证明在许多领域是有效的，如问题回答、翻译、阅读理解和文本分类等任务。

给定文本内容 T ，将 T 建模为单词序列 $T = \{w_1, w_2, \dots, w_m\}$ (m 表示文本中的单词数)，将处理后的文本特征表示为 $P_T = \{P_{t1}, P_{t2}, \dots, P_{tm}\}$ 。其中， P_{ti} 表示第 i 个单词 w_i 的特征。词的表示方法 P_{ti} 是由预先训练的BERT模型计算的：

$$P_T = \{P_{t1}, P_{t2}, \dots, P_{tm}\} = \text{BERT}(W) \quad (1)$$

其中， $P_{ti} \in \mathbb{R}^{d_t}$ 是BERT中对应标记的输出层隐藏状态， d_t 为词嵌入的维度。

3.3.2 视觉编码网络

对给定视觉内容 V ，我们使用自下而上的注意力预训练模型ResNet50来提取区域特征。输

出是一组区域特征 $P_V = \{P_{v1}, P_{v2}, \dots, P_{vn}\}$ (n 表示图像中的区域数量)。其中, 每个 P_{vj} 被定义为第 j 个区域的平均池化卷积特征。在训练期间, 预训练模型参数保持固定。同时, 在处理给定的视觉内容 V 的过程中, 视觉特征提取器中倒数第二层池化层的操作可以表示为:

$$P_V = \{P_{v1}, P_{v2}, \dots, P_{vn}\} = \text{ResNet50}(V) \quad (2)$$

其中, $P_{vj} \in \mathbb{R}^{d_v}$, d_v 表示为图像嵌入层的维度。此外, 通过另外计入 2D 卷积层, 将嵌入维度 d_v 调整为 d_t , 以满足任务需求。

3.3.3 评论编码网络

对于给定评论内容 C , 我们同样使用 BERT 提取评论文本特征。给定内容评论 $C = \{c_1, \dots, c_o\}$ (o 表示评论数量), 转换后得到评论特征由 $P_{Ch} = \{P_{c1}, \dots, P_{co}\}$ 表示。其中, 每个 P_{Ch} 对应第 h 个评论 C_h 的特征, 评论特征 P_{Ch} 通过与训练的 BERT 计算得到:

$$P_{Ch} = \text{BERT}(c_h) \quad (3)$$

其中, $P_{Ch} \in \mathbb{R}^{d_t}$ 是 BERT 中隐藏层的池化输出, BERT 中隐藏层的池化输出, d_t 表示单词嵌入的维度。

3.4 多模态 gather transformer 网络

为了有效整合社交媒体帖子中文本、视觉和评论特征, 本文设计了多模态双协同 gather transformer 网络来构建多模态上下文信息, 并充分捕捉和融合高阶互补信息。如图 1 所示, 双协同 gather transformer 网络由于文本-图像 gather transformer、文本-评论 gather transformer 两个模块构成。

在 transformer 模型中, Attention 机制是不可或缺的, 它构成了 Sequence-to-Sequence 处理中编码器-解码器架构的核心, 有利于建立输入与输出之间的长程依赖关系。对于长度为 L 的输入序列 $X \in \mathbb{R}^{L \times D}$, 单头 self-attention 机制的表达公式如下:

$$\text{Att}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (4)$$

其中, $\text{softmax}(\cdot)$ 和 V 的乘积表示由查询值 Q 和键值 K 之间的交互得到的注意力分数对输入的加权组合。分母中缩放因子 \sqrt{d} 能有效抑制点积幅度增长, 进而稳定度量标准。

为了检索和捕捉模态信息中的关键部分, 充分理解和提取模态之间细粒度相互关系。我们通过添加一个多模态 gather 网络来增强传统 Transformer 网络的多头注意力架构。

3.4.1 文本-图像 gather transformer 网络

为了充分捕捉和检索文本与图像模态信息间细粒度内在联系, 我们设计了文本-图像 gather transformer 网络。

将经过预训练后得到的视觉特征 P_V 和文本特征 P_T 分别作为 input1 和 input2 (如图 1 所示) 输入到文本-图像 gather-transformer 网络当中。在自注意力机制中, 对于多模态输入 $P_V = \{P_{v1}, P_{v2}, \dots, P_{vn}\} \in \mathbb{R}^{d_v}$ 和 $P_T = \{P_{t1}, P_{t2}, \dots, P_{tm}\} \in \mathbb{R}^{d_t}$, 我们计算相似度矩阵 S :

$$S = \text{Softmax}\left(\frac{P_T P_V^T}{\sqrt{d_k}}\right) \quad (5)$$

其中, $S_{i,j}$ 表示 P_T 中第 i 个单词与 P_V 中第 j 个区域之间的相似度得分。将相似度矩阵 S 展开为元素个数为 N (N 为 n 与 m 的乘积) 的线性序列 s :

$$s = (s_0, s_1, \dots, s_{N-1}) \quad (6)$$

$$s_g = \text{gather}(s) \quad (7)$$

通过引入 Gather 机制, 筛选线性序列 s 中元素权重最高的 K 个元素, 得到新的线性序列 s_g 。最后, 将序列 s_g 还原为经 gather 机制处理后的相似度矩阵 S_g 。为方便理解, Gather 机制具体实现流程如图 2 所示。

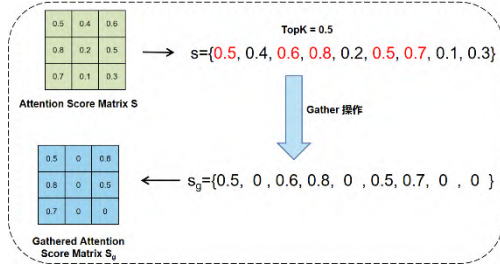


图2 gather 机制处理流程

Fig. 2 Process of the gather mechanism

基于此, 我们将文本-图像 gather transformer 网络中修改后的单头多模态自注意力机制定义为:

$$\begin{aligned} & \text{GatheredAtt}(Q, K, V, G) \\ &= \text{Softmax}\left(\frac{\mathbb{G}(QK^T)}{\sqrt{d}}\right)V \end{aligned} \quad (8)$$

其中, $Q = W_Q^1 P_T$, $K = W_K^1 P_V$, $V = W_V^1 P_V$ 。 W_Q^1 、 W_K^1 和 W_V^1 分别表示将输入投影到查询、键和值的不同线性变换。符号 $\mathbb{G}(\cdot)$ 表示 gather 机制的处理操作。

3.4.2 文本-评论 gather transformer 网络

为了充分捕捉和检索文本与评论模态信息间细粒度内在联系, 我们设计了文本-评论 gather transformer 网络。

将经过预训练后得到的评论特征 P_{Ch} 和文本特征 P_T 分别作为 input1 和 input2 (如图 1 所示) 输入到文本-评论 gather-transformer 网络当中。在自注意力机制中, 对于多模态输入 $P_{Ch} = \{P_{c1}, \dots, P_{co}\} \in \mathbb{R}^{d_t}$ 和 $P_T = \{P_{t1}, P_{t2}, \dots, P_{tm}\} \in \mathbb{R}^{d_t}$, 可计算相似度矩阵 S :

$$S = \text{Softmax}\left(\frac{P_T P_C^T}{\sqrt{d_t}}\right) \quad (9)$$

其中, $S_{q,p}$ 表示 P_T 中第 q 个单词与 P_C 中第 p 个评论之间的相似度得分。接下来, Gather 机制的筛选、检索流程与 3.4.1 部分相似。

基于此, 我们将文本-评论 gather

transformer 网络中改进后的单头多模态自注意力机制定义为:

$$\begin{aligned} & \text{GatheredAtt}(Q, K, V, G) \\ &= \text{Softmax}\left(\frac{\mathbb{G}(QK^T)}{\sqrt{d}}\right)V \end{aligned} \quad (10)$$

其中, $Q = W_Q^2 P_T$, $K = W_K^2 P_C$, $V = W_V^2 P_C$ 。 W_Q^2 、 W_K^2 和 W_V^2 分别表示将输入投影到查询、键和值的不同线性变换。同样, 符号 $\mathbb{G}(\cdot)$ 表示 gather 机制的处理操作。

3.5 双协同机制

设文本-图像 gather transformer 的输出为 P_{TV} , 文本-评论 gather transformer 的输出为 P_{TC} 。为了探索文本、图像和评论之间的内在关系, 我们通过两个 Cross-Attention 网络进一步对 P_{TV} 和 P_{TC} 进行编码, 使其充分协同和交互。

如图 1 所示, 对于上方的 Cross-Attention, 将 P_{TV} 作为输入 (即 Q 值), P_{TC} 作为输入 (即 K 值和 V 值), 对应的 Attention 操作如下式:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{P_{TV} P_{TC}^T}{\sqrt{d}}\right) P_{TC} \quad (11)$$

与此对应, 下方的 Cross-Attention 操作, 将 P_{TC} 作为输入 (即 Q 值), P_{TV} 作为输入 (即 K 值和 V 值)。

3.6 虚假信息检测器

在得到双协同机制中, Cross-Attention 模块输出的两个多模态表示 F_{ca1} 和 F_{ca2} 后, 我们将其连接起来得到最终特征表示:

$$F_f = \text{Concat}(F_{ca1} + F_{ca2}) \quad (12)$$

其中, 虚假信息检测器以最终多模态表示 F_f 作为输入, 它包含相应激活函数的全连接构成:

$$\hat{y} = \text{softmax}(W * F_f + b) \quad (13)$$

式中, $\text{softmax}(\cdot)$ 表示激活函数, \hat{y} 表示该帖子的预测概率, F_f 表示该帖子最终多模态表示。为了增强模型检测效果, 我们采用交叉熵损失函

数Loss来训练本模型:

$$Loss = - \sum_{n=1}^N y_n \log(\hat{y}_n) \quad (14)$$

N 表示社交媒体帖子的数量, \hat{y}_n 表示第 n 个帖子的预测概率, y_n 表示第 n 个帖子的真实标签。

4 实验

4.1 实验准备

4.1.1 数据集

我们使用两个公开可用的数据集 ReCOVery^[18]和MMCoVaR^[17], 将MDCGTN模型与其他基准模型进行有效性验证。表1展示了我们实验的两个数据集包含虚假信息、真实信息、图像和评论数量情况。其中, ReCOVery包含从2020年1月到2020年5月中涵盖多模态信息的COVID-19新闻帖子。MMCoVaR包含特定COVID-19疫苗相关的多模态社交媒体帖子, 帖子编辑时间为2020年2月至2021年5月。根据^[15]的方法, 我们将两个数据集按8:2的比例划分为训练集和测试集。

表1 ReCOVery和MMCoVaR数据集

Table 1 ReCOVery and MMCoVaR dataset

News Articles	ReCOVery	MMCoVaR
# of Fake News	1364	1635
# of Real News	665	958
# of Images	1675	22357
# of Comments	140820	24183

4.1.2 评价指标

虚假信息检测任务属于二分类, 一般使用准确率 (Accuracy) 作为主要评估指标。但在样本数据不均衡情况下, 其可靠性可能会降低。因此, 在实验过程中, 我们增加了精度 (Precision)、召回率 (Recall) 和加权 F1 分数 (Weighted F1-score) 作为补充评估指标, 以解决数据不平衡导致的问题。

4.1.3 实验设置

对于社交媒体帖子多模态嵌入, 我们使用预

训练得到的 ResNet 来提取视觉特征, 使用预训练的 BERT 模型来提取帖子文本和评论中的文本特征。其中, 图像嵌入维度为 2048, 文本嵌入维度为 768。为了适应本模型, 通过引入一个 2D 卷积层, 将视觉区域特征维度从 2048 转换为 768。我们模型使用自适应矩估计 (Adam) 优化器进行训练, 一共训练 200 个 epoch, 学习率为 0.001, 小批量设置为 64。

4.2 比较方法

我们选取了 11 种最先进的模型进行比较:

(1) HSA^[37]: HSA (Hierarchical Social Analysis) 是一种利用媒体中的用户评论和层次化社交网络结构进行虚假信息检测的方法。

(2) ExFaux^[38]: ExFaux 是一种基于图的虚假图像解释框架, 可以为检测结果提供内容解释。

(3) dEFEND^[39]: dEFEND 是一种可解释的虚假信息检测方法。它利用新闻文本与用户评论之间的相关性, 既可以对虚假新闻进行分类, 又可以确定用户评论对分类结果提供解释的方式。

(4) BTIC^[40]: BTIC 采用基于 BERT 的多模态框架用于不可靠新闻检测, 通过对可疑文章中的文本和视觉信息使用对比学习策略来利用这些信息。

(5) SAFE^[23]: SAFE 是一种以相似性为中心的多模态方法, 用于虚假信息检测。它从新闻素材中提取文本和视觉元素, 并探索它们之间的相互关系以获得最终表示。

(6) EANN^[21]: EANN 包括多模态特征提取器、虚假新闻检测器和后处理鉴别器。

(7) SpotFake^[41]: SpotFake 使用预训练的 BERT 提取文本特征, 同时使用在 ImageNet 数据集上预训练的 VGG-19 提取图像特征, 进而辅助判断帖子的真实性。

(8) MVAE^[14]: MVAE 通过将多模态特征引入到双模态变分自编码器 (VAE) 中, 获得多模态表

示。再经过二元分类器, 对双模态 VAE 产生的多模态潜在表示进行分类。

(9) FMFN^[10]: FMFN 是一种细粒度多模态融合网络, 它利用缩放点积注意力机制来融合文本中单词的词嵌入, 以及表示图像不同特征的多个特征向量。

(10) MMTN^[42]: MMTN 是一种多模态掩码 Transformer 网络, 它利用掩码 Transformer 网络联合建模多模态信息的模态间和模态内关系, 并屏蔽无关上下文信息, 实现虚假信息检测任务。

(11) DGExplain^[15]: DGExplain 是一种生成性方法, 通过分析视觉和文本信息的模态间联系, 进而识别与 COVID-19 相关多模态帖子的真实性。

4.3 实验结果分析

4.3.1 定量分析

实验结果如表 2 和表 3 所示, 通过观察表中的数据我们可以得到以下结论:

(1) 基于特征融合的方法在两个数据集上表现效果一般, 如 MVAE、SpotFake、EANN 和 ExFaux。这些方法的表现不佳可归因于其忽视了社交媒体帖子中视觉和文本元素之间的复杂关系。它们通过直接从混合的多模态内容中推断待检测内容的可信度, 而不是考虑这种相互关联。

(2) DGExplain 超过其他基线方法, 这表明它精确生成跨模态特征, 结合生成的和原始的多模态特征之间的一致性评估, 进而辅助检测虚假信息。同时, DGExplain 对内容-评论图的结合大大增强了生成的、原始的多模态特征和用户评论的整合, 进而有助于虚假信息的检测。

(3) 所提出的 MDCGTN 在两个数据集上均超过了所有基准方法。结果表明, 所提出的基于多模态双协同 gather transformer 网络能够充分捕捉模态间和模态内的细粒度特征关系, 检索和收集模态信息中关键部分, 最终达到不错检测效果。

表 2 不同方法在 ReCOVery 数据集上的实验结果

Table 2 Detection results of different methods on ReCOVery dataset

Methods	Acc (%)	P (%)	R (%)	F1(%)
HSA	77.9	73.7	73.6	73.6
ExFaux	76.3	71.9	69.5	70.4
dEFEND	85.6	82.6	81.3	82.3
BTIC	76.3	71.9	69.5	70.4
SAFE	83.1	80.3	78.9	79.5
EANN	84.7	81.6	83.4	82.4
SpotFake	68.1	63.7	65.0	64.1
MVAE	82.5	81.3	75.5	77.4
FMFN	87.4	85.3	84.5	84.9
MMTN	88.2	88.3	82.9	85.0
DGExplain	89.7	89.0	86.1	87.3
MDCGTN(Ours)	90.15	89.24	88.01	88.53

表 3 不同方法在 MMCoVaR 数据集上的实验结果

Table 3 Detection results of different methods on MMCoVaR dataset

Methods	Acc (%)	P (%)	R (%)	F1(%)
HSA	80.3	78.2	78.5	78.4
ExFaux	76.9	78.4	69.4	70.7
dEFEND	85.6	84.7	83.1	83.8
BTIC	82.9	82.3	79.1	80.3
SAFE	78.8	77.3	74.9	75.7
EANN	83.3	81.9	81.0	81.4
SpotFake	69.9	67.0	62.0	62.3
MVAE	81.5	80.5	83.4	80.8
FMFN	87.3	87.1	85.5	86.2
MMTN	88.4	87.7	87.6	87.7
DGExplain	89.5	89.6	87.1	88.1
MDCGTN(Ours)	90.40	90.20	89.13	89.59

4.3.2 消融实验

为了验证我们模型的有效性, 我们比较了 MDCGTN 的几个变式:

(1) $MDCGTN \rightarrow V$: 除掉视觉信息, 仅保留文本信息;

(2) $MDCGTN \rightarrow C$: 除掉评论信息, 仅使用帖子文本信息和视觉信息;

(3) $MDCGTN \rightarrow G$: 除掉 gather 机制, 仅使用传统的 Transformer 模块;

(4) $MDCGTN \rightarrow D$: 除掉双协同机制, 并利用自注意力网络替换交叉注意力网络。

表 4 和表 5 展示了模型的不同变式在两个数据集上的实验结果。

表 4 模型不同变式在 ReCOVery 数据集上的实验结果

Table 4 Performance comparison of different variants of the model on ReCOVery dataset

Methods	Acc (%)	P (%)	R (%)	F1(%)
$MDCGTN \neg V$	85.47	85.17	81.50	82.50
$MDCGTN \neg C$	87.83	87.14	84.42	85.50
$MDCGTN \neg G$	88.22	87.59	85.09	86.13
$MDCGTN \neg D$	87.98	87.04	85.32	86.28
$MDCGTN$	90.15	89.24	88.01	88.53

表 5 模型不同变式在 MMCoVaR 数据集上的实验结果

Table 5 Performance comparison of different variants of the model on MMCoVaR dataset

Methods	Acc (%)	P (%)	R (%)	F1(%)
$MDCGTN \neg V$	82.35	84.79	78.02	79.54
$MDCGTN \neg C$	87.67	87.89	85.95	86.47
$MDCGTN \neg G$	87.57	87.48	85.94	86.47
$MDCGTN \neg D$	89.15	89.36	87.39	88.12
$MDCGTN$	90.40	90.20	89.13	89.59

从表 4 和表 5 中，我们可以从实验结果中得到以下结论：

(1) 视觉信息的影响：为了评估视觉信息的效果，我们对比了 $MDCGTN$ 和 $MDCGTN \neg V$ 在两个数据集上的性能。从结果中，我们可观察到所提出的 $MDCGTN$ 优于 $MDCGTN \neg V$ ，这表明视觉信息能为本模型提供有价值的补充信息。

(2) 用户评论的影响：通过对比 $MDCGTN$ 和 $MDCGTN \neg C$ 在两个数据集上的表现，我们可以得到 $MDCGTN$ 的性能优于 $MDCGTN \neg C$ 的结论，进而表明用户评论能为本模型提供有效的补充信息。

(3) gather transformer 网络的影响：为了评估 gather transformer 网络的效果，我们对比了 $MDCGTN$ 和 $MDCGTN \neg G$ 在两个数据集上的性能。从结果中，我们可观察到所提出的 $MDCGTN$ 优于 $MDCGTN \neg G$ ，这表明 gather transformer 网络能为本模型提供有价值的补充信息。

(4) 双协同机制的影响：通过对比 $MDCGTN$ 和 $MDCGTN \neg D$ 在两个数据集上的性能。从结果中，我们可观察到所提出的 $MDCGTN$ 优于 $MDCGTN \neg D$ ，这表明双协同机制的有效性。

4.3.3 参数分析

图 3 展示了参数 topK 值的变化对模型性能的影响。在 ReCOVery 数据集上，我们的模型准确率在 topK=0.1 时到达最高，在此之后的 topK 值的准确性呈下降趋势，不能达到同样水平。在 MMCoVaR 数据集上，我们注意到模型的准确率同样在 topK=0.1 时到达最高。因此，可以得出当 topK=0.1 时，本模型能够达到最高性能。



图 3 TopK 值对两个数据集中虚假信息检测的准确率和 F1 分数的影响

Fig. 3 The influence of TopK on accuracy and F1 scores pertaining to fake news detection across two datasets.

结束语 在本文中，我们设计了一个多模态双协同 gather transformer 模型对社交媒体虚假信息进行检测。现有的方法大多缺乏对多模态信息中关键部分提取，未能有效实现多模态信息间相互指导，多模态信息间缺乏彼此交互性。为了解决上述挑战，我们提出了多模态双协同

gather transformer 网络来建模多模态信息的模态间和模态内关系, 并对关键信息进行筛选。我们的方法包含三个主要组件: (1) 我们利用 ResNet 学习视觉表示, 并利用 BERT 学习文本表示。(2) 我们使用 gather transformer 网络更好的探索和捕捉模态间和模态内的细粒度联系, 有效检索和筛选模态信息中关键部分。(3) 我们引入一种双协同机制, 提高多模态信息间指导性, 提高文本与视觉特征间交互增强能力。在两个基准数据集上的实验表明, 我们提出的方法更加有效。下一步, 将讨论如何更好利用用户评论信息, 以及如何更好实现多模态信息融合。

参考文献

- [1] CASTILLO C, MENDOZA M, POBLETE B. Information credibility on twitter[C]//Proceedings of the 20th international conference on World wide web,2011:675-684.
- [2] KWON S, CHA M, JUNG K, et al. Prominent features of rumor propagation in online social media[C]//2013 IEEE 13th international conference on data mining,IEEE,2013:1103-1108.
- [3] LIU X, NOURBAKHSI A, LI Q, et al. Real-time rumor debunking on twitter[C]//Proceedings of the 24th ACM international on conference on information and knowledge management,2015:1867-1870.
- [4] MA J, GAO W, WEI Z, et al. Detect rumors using time series of social context information on microblogging websites[C]//Proceedings of the 24th ACM international on conference on information and knowledge management,2015: 1751-1754.
- [5] MA J, GAO W, MITRA P, et al. Detecting rumors from microblogs with recurrent neural networks[J],2016.
- [6] YU F, LIU Q, WU S, et al. A Convolutional Approach for Misinformation Identification[C]//IJCAI,2017:3901-3907.
- [7] WANG J, WANG Y C, HUANG M J. False information in social networks: Definition, detection and control[J].Comput. Sci, 2021, 48: 263-277. (in Chinese)
- 王剑, 王玉翠, 黄梦杰. 社交网络中的虚假信息: 定义, 检测及控制[J]. 计算机科学, 2021.
- [8] HAO X, MING L. Deepfake Video Detection Based on 3D Convolutional Neural Networks[J].Computer Science,2021, 48(7):86-92. (in Chinese)
- 邢豪, 李明. 基于 3D CNNs 的深度伪造视频篡改检测[J]. 计算机科学, 2021, 48(7): 86-92.
- [9] PROCTER R, CRUMP J, KARSTEDT S, et al. Reading the riots: What were the police doing on Twitter?[J].Policing and society,2013,23(4): 413-436.
- [10] WANG J, MAO H, LI H. FMFN: Fine-grained multimodal fusion networks for fake news detection[J].Applied Sciences,2022,12(3): 1093.
- [11] QIAN S S, ZHANG T Z, XU C S. Survey of Multimedia Social Events Analysis[J]. Comput. Sci,2021,48(3): 97-112. (in Chinese)
- 钱胜胜, 张天柱, 徐常胜. 多媒体社会事件分析综述[J].计算机科学, 2021, 48(3): 97-112.
- [12] WU X K, ZHAO T F. Application of natural language processing in social communication: A review and future perspectives[J].Computer Science,2020,47(6): 184-193. (in Chinese)
- 吴小坤, 赵甜芳. 自然语言处理技术在社会传播学中的应用研究和前景展望[J].计算机科学, 2020, 47(6): 184-193.
- [13] HAN Z M, ZHENG C Y, DUAN D G, et al. Associated Users Mining Algorithm Based on Multi-information Fusion Representation Learning[J].Computer Science,2019,46(4): 77r82. (in Chinese)
- 韩忠明, 郑晨烨, 段大高, 等. 基于多信息融合表示学习的关联用户挖掘算法[J].计算机科学, 2019, 46(4): 77r82.
- [14] KHATTAR D, GOUD J S, GUPTA M, et al. Mvae: Multimodal variational autoencoder for fake news detection[C]//The world wide web conference,2019:2915-2921.
- [15] SHANG L, KOU Z, ZHANG Y, et al. A duo-generative approach to explainable multimodal covid-19 misinformation detection[C]//Proceedings of the ACM Web Conference 2022,2022: 3623-3631.
- [16] LIN H, MA J, CHENG M, et al. Rumor detection on twitter with claim-guided hierarchical graph attention networks[J].arXiv preprint arXiv:2110.04522, 2021.
- [17] CHEN M, CHU X, SUBBALAKSHMI K P. MMCoVaR: multimodal COVID-19 vaccine focused data repository for fake news detection

- and a baseline architecture for classification[C]//Proceedings of the 2021 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining,2021:31-38.
- [18] ZHOU X, MULAY A, FERRARA E, et al. Recovery: A multimodal repository for covid-19 news credibility research[C]//Proceedings of the 29th ACM international conference on information & knowledge management,2020:3205-3212.
- [19] SINGHAL S, KABRA A, SHARMA M, et al. Spotfake+: A multimodal framework for fake news detection via transfer learning (student abstract)[C]//Proceedings of the AAAI conference on artificial intelligence,2020, 34(10): 13915-13916.
- [20] LIU J, FENG K, PAN J Z, et al. MSRD: Multimodal Web Rumor Detection Method[J]. 2020.
- [21] WANG Y, MA F, JIN Z, et al. Eann: Event adversarial neural networks for multi-modal fake news detection[C]//Proceedings of the 24th acm sigkdd international conference on knowledge discovery & data mining,2018: 849-857.
- [22] YAO L, MAO C, LUO Y. Graph convolutional networks for text classification[C]//Proceedings of the AAAI conference on artificial intelligence,2019, 33(01): 7370-7377.
- [23] ZHOU X, WU J, ZAFARANI R. Safe: similarity-aware multi-modal fake news detection (2020)[J].Preprint. arXiv,2020, 200304981: 2.
- [24] WEI Z X, LIANG J M. Design of Image Retrieval System Based on Speech Recognition[J].Applied Mechanics and Materials,2012,220:2371-2374.
- [25] XUE J, WANG Y, TIAN Y, et al. Detecting fake news by exploring the consistency of multimodal data[J].Information Processing & Management,2021,58(5):102610.
- [26] XU K, BA J, KIROS R, et al. Show, attend and tell: Neural image caption generation with visual attention[C]//International conference on machine learning,PMLR,2015:2048-2057.
- [27] BAHDANAU D, CHO K, BENGIO Y. Neural machine translation by jointly learning to align and translate[J].arXiv preprint arXiv:1409.0473, 2014.
- [28] CHEN J, ZHANG H, HE X, et al. Attentive collaborative filtering: Multimedia recommendation with item-and component-level attention[C]//Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval,2017:335-344.
- [29] WANG S, HU L, CAO L, et al. Attention-based transactional context embedding for next-item recommendation[C]//Proceedings of the AAAI conference on artificial intelligence,2018,32(1).
- [30] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[J].Advances in neural information processing systems,2017,30.
- [31] BAEVSKI A, AULI M. Adaptive input representations for neural language modeling[J]. arXiv preprint arXiv:1809.10853, 2018.
- [32] DAI Z, YANG Z, YANG Y, et al. Transformer-xl: Attentive language models beyond a fixed-length context[J].arXiv preprint arXiv:1901.02860,2019.
- [33] RAE J W, POTAPENKO A, JAYAKUMAR S M, et al. Compressive transformers for long-range sequence modelling[J].arXiv preprint arXiv:1911.05507,2019.
- [34] RADFORD A, NARASIMHAN K, SALIMANS T, et al. Improving language understanding by generative pre-training[J]. 2018.
- [35] DEVLIN J, CHANG M W, LEE K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J].arXiv preprint arXiv:1810.04805,2018.
- [36] CHEN T, LI X, YIN H, et al. Call attention to rumors: Deep attention based recurrent neural networks for early rumor detection[C]//Trends and Applications in Knowledge Discovery and Data Mining: PAKDD 2018 Workshops, BDASC, BDM, ML4Cyber, PAISI, DaMEMO, Melbourne, VIC, Australia, June 3,2018, Revised Selected Papers 22.Springer International Publishing,2018: 40-52.
- [37] GUO H, CAO J, ZHANG Y, et al. Rumor detection with hierarchical social attention network[C]//Proceedings of the 27th ACM international conference on information and knowledge management,2018:943-951.
- [38] KOU Z, ZHANG D Y, SHANG L, et al. Exfaux: A weakly supervised approach to explainable fauxtography detection[C]//2020 IEEE International Conference on Big Data (Big Data),IEEE,2020:631-636.
- [39] SHU K, CUI L, WANG S, et al. defend: Explainable fake news detection[C]//Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining,2019:395-405.

- [40] ZHANG W, GUI L, HE Y. Supervised contrastive learning for multimodal unreliable news detection in covid-19 pandemic[C]//Proceedings of the 30th ACM International Conference on Information & Knowledge Management,2021:3637-3641.
- [41] SINGHAL S, SHAH R R, CHAKRABORTY T, et al. Spotfake: A multi-modal framework for fake news detection[C]//2019 IEEE fifth international conference on multimedia big data (BigMM),IEEE, 2019: 39-47.
- [42] WANG J, QIAN S, HU J, et al. Positive Unlabeled Fake News Detection Via Multi-Modal Masked Transformer Network[J].IEEE Transactions on Multimedia,2023.

向旺, 出生于 1996 年, 硕士研究生, 主要研究方向为自然语言处理和多媒体计算分析。

王金光, 出生于 1989 年, 博士研究生, 主要研究方向为数据挖掘和多媒体信息分析。

王一飞, 出生于 1997 年, 硕士研究生, 主要研究方向为计算机视觉和自然语言处理。

钱胜胜, 出生于 1991 年, 博士, 副研究员, 主要研究方向为数据挖掘和多媒体内容分析。



XIANG Wang, born in 1996, postgraduate. His main research interests include natural language processing and multimedia computing analysis .



QIAN Sheng-sheng, born in 1991, Ph.D., Associate Professor. His main research interests include data mining and multimedia content analysis.