

# 基于跨模态引导和对齐的多模态预训练方法

才 华<sup>1,2</sup>, 易亚希<sup>1</sup>, 付 强<sup>3</sup>, 冉 越<sup>1</sup>, 孙俊喜<sup>4</sup>

(1. 长春理工大学电子信息工程学院, 吉林长春 130022; 2. 长春中国光学科学技术馆, 吉林长春 130117;  
3. 长春理工大学空间光电技术研究所, 吉林长春 130022; 4. 东北师范大学信息科学与技术学院, 吉林长春 130117)

**摘 要:** 现有的视觉语言多模态预训练方法仅在图像和文本的全局语义上进行特征对齐, 对模态间细粒度特征交互的探索不足. 针对这一问题, 本文提出了一种基于跨模态引导和对齐的多模态预训练方法. 该方法在模态特征提取阶段, 采用基于视觉序列压缩的双流特征提取网络, 在视觉编码器中联合图像和文本信息逐层引导视觉序列压缩, 缓解与文本无关的冗余视觉信息对模态间细粒度交互的干扰; 在模态特征对齐阶段, 对图像和文本特征进行细粒度关系推理, 实现视觉标记与文本标记的局部特征对齐, 增强对模态间细粒度对齐关系的理解. 实验结果表明, 本文方法能够更好地对齐视觉文本的细粒度特征, 在图文检索任务中, 微调后的图像检索和文本检索的平均召回率分别达到了 86.4% 和 94.88%, 且零样本图文检索的整体指标相较于经典图文检索算法 CLIP (Contrastive Language-Image Pre-training) 提升了 5.36%, 在视觉问答等分类任务中, 准确率也优于目前主流多模态预训练方法.

**关键词:** 多模态预训练; 跨模态引导; 视觉序列压缩; 双流特征提取; 细粒度关系推理; 局部特征对齐

**基金项目:** 国家自然科学基金 (No.61890963, No.U2341226); 吉林省人才专项 (No.20240602015RC); 西安市飞行器光学成像与测量技术重点实验室开放基金 (No.2023-13)

**中图分类号:** TP391

**文献标识码:** A

**文章编号:** 0372-2112(2024)10-3368-14

**电子学报 URL:** <http://www.ejournal.org.cn>

**DOI:** 10.12263/DZXB.20240271

## Multimodal Pretraining with Cross-Modal Guidance and Alignment

CAI Hua<sup>1,2</sup>, YI Ya-xi<sup>1</sup>, FU Qiang<sup>3</sup>, RAN Yue<sup>1</sup>, SUN Jun-xi<sup>4</sup>

(1. School of Electronic Information and Engineering, Changchun, Jilin 130022, China;

2. Changchun China Optics Science and Technology Museum, Changchun, Jilin 130117, China;

3. School of Opto-Electronic Engineer, Changchun University of Science and Technology, Changchun, Jilin 130022, China;

4. School of Information Science and Technology, Northeast Normal University, Changchun, Jilin 130117, China)

**Abstract:** Current multimodal pre-training techniques for visual languages predominantly focus on aligning global semantic features between images and text, yet they inadequately explore the granular feature interactions between modalities. Addressing this gap, this paper proposes a novel multimodal pre-training strategy informed by cross-modal guidance and alignment. Our method employs a dual-stream feature extraction network designed for visual sequence compression, to facilitate modality feature extraction. During this phase, a synergistic image-text guidance is integrated within the visual encoder, orchestrating the compression of visual sequences layer by layer. This approach mitigates the obfuscation of modality-specific fine-grained interactions by irrelevant visual information. Subsequently, in the modality feature alignment phase, we implement fine-grained relational reasoning on the image and textual features to achieve localized feature alignment among visual tokens and textual tokens. This advancement bolsters the model's comprehension of fine-grained alignment relationships. After fine-tuning, in the image-text retrieval tasks, our approach achieves an average recall rate of 86.4% for images and 94.88% for texts, which represents a significant 5.36% improvement in zero-shot image-text retrieval over the canonical CLIP (Contrastive Language-Image Pre-training) algorithm. Moreover, our method also surpasses existing mainstream multimodal pre-training methods in accuracy for classification tasks like visual question answering.

**Key words:** multimodal pre-training; cross-modal guidance; visual sequence compression; dual-stream feature extraction; fine-grained relational reasoning; localized feature alignment

**Foundation Item(s):** National Natural Science Foundation of China (No.61890963, No.U2341226); Jilin Province

Talent Development Special Fund (No.20240602015RC); Xi'an Key Laboratory of Aircraft Optical Imaging and Measurement Technology (No.2023-13)

## 1 引言

视觉语言多模态学习属于计算机视觉与自然语言处理的交叉领域,旨在通过同时处理视觉和语言两种模态信息,实现跨模态的联合学习与理解.在视频理解<sup>[1]</sup>、生物医学<sup>[2,3]</sup>、自动驾驶<sup>[4]</sup>等领域应用广泛.与单一模态信息处理不同,图像数据和文本数据在表达方式和信息密度上都存在较大差异,这种不可忽视的模态异构性和语义鸿沟给多模态模型学习视觉语言对齐关系带来了极大的挑战.

现有的多模态预训练模型在处理视觉语言模态的对齐问题上,主要采取两种方法.一种方法是基于图像区域和文本单词的对齐<sup>[5-7]</sup>,利用训练好的目标检测器识别图像中所有潜在对象,实现文本单词与具体视觉对象的细粒度特征对齐.然而,这种方法对下游任务的迁移受到目标检测器的类别限制,且识别到的潜在对象也可能与文本描述毫无关联.另一种方法基于图文对比学习 ITC (Image-Text Contrastive) 对齐<sup>[8-10]</sup>,专注于在图像和文本的全局语义上进行粗粒度对齐,克服了目标检测器的限制.但仅进行全局语义的跨模态特征对齐,容易导致多模态模型忽略视觉对象和文本单词之间的细粒度对齐关系,而这对于许多下游任务是至关重要的<sup>[11,12]</sup>.为此, FILIP<sup>[13]</sup> (Fine-grained Interactive Language-Image Pre-training) 提出了更精细的标记级对比对齐,但没有做进一步的模态交互融合,在复杂的视觉语言理解任务上表现较差. ALBEF<sup>[14]</sup> (ALign BEfore Fuse) 提出先对齐再融合的多模态学习策略,在多模态融合编码器中学习模态间细粒度交互关系,在各类多模态下游任务中取得了不错的效果,但其将模态间的深层次交互完全依托于融合过程中的注意力机制实现,缺乏更为显式的细粒度对齐手段.

此外,由于在现有的多模态数据集中,图像通常包含比简短文本描述更丰富的信息,这种模态间的信息不对称性导致多模态模型在使用 ViT<sup>[15]</sup> (Vision Transformer) 提取图像特征时,会不可避免地编码与文本内容无关的冗余背景信息,这不仅增加了网络模型的计算负担,同时也会阻碍文本和图像之间的细粒度交互<sup>[16]</sup>.尽管在计算机视觉领域,研究者对于 ViT 中的长视觉序列问题,提出了许多序列修剪<sup>[17,18]</sup>、序列重组<sup>[19]</sup>以及序列压缩<sup>[20]</sup>等算法,以减少视觉序列长度,提升计算效率.但这些方法仅使用单一的视觉信息指导视觉序列修剪,忽略了文本信息,因此并不适合直接应用于需要同时处理两种模态信息的视觉语言多模态模型.

针对上述问题,本文遵循先对齐再融合的预训练

策略,提出了一种基于跨模态引导和对齐的多模态预训练方法 (Multimodal pretraining with cross-modal Guidance and Alignment, MulGA),在视觉语言特征提取过程中,通过双模态信息的联合引导,逐步压缩视觉序列中与文本描述无关的冗余标记,缓解背景特征对模态交互的干扰;同时在模态对齐阶段进行细粒度关系推理,加强模型对视觉对象和文本单词之间深层关联的理解,实现更精准的特征对齐.

## 2 方法架构

本文提出的视觉语言多模态预训练模型 MulGA 整体网络框架如图 1 所示.该框架由基于视觉序列压缩的双流特征提取网络、细粒度关系推理模块以及多模态交互融合模块三个核心组件构成.

首先,输入的图像-文本对经过双流特征提取网络提取特征向量,并在细粒度关系推理模块中进行初步的特征细粒度对齐处理;随后,将初步对齐的图像特征和文本特征馈送到多模态交互融合模块,通过执行图像文本匹配 (Image-Text Matching, ITM) 和掩蔽语言建模 (Masked Language Modeling, MLM) 两个预训练任务,实现进一步的模态细粒度融合;最终输出深度融合后的多模态表征.具体流程如算法 1 所示.

### 2.1 基于视觉序列压缩的双流特征提取网络

本文采用双流特征提取架构,通过两个独立的单模态编码器分别提取输入的图像和文本特征,如图 2 所示.

文本编码器由 6 层 Transformer<sup>[21]</sup> 块组成,每层包含一个多头自注意力 (Multi-head Self-Attention, MSA) 和一个前馈神经网络 (Feedforward Neural Network, FNN),层归一化和残差连接用于防止梯度消失.接收原始文本和掩蔽文本输入,将输入的文本单词转换为词嵌入向量  $\mathbf{t}_i$ ,并与位置嵌入  $\mathbf{T}_{\text{pos}}$  求和构造输入表示.需要注意的是,掩蔽文本输入仅作为后续 MLM 预训练任务使用,不参与文本编码器的梯度更新.

文本编码器计算过程,见式(1)~(3).

$$\mathbf{W}_0 = (\mathbf{t}_{\text{cls}}, \mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3, \dots, \mathbf{t}_N) + \mathbf{T}_{\text{pos}} \quad (1)$$

$$\mathbf{W}'_l = \text{LN}(\text{MSA}(\mathbf{W}_{l-1}) + \mathbf{W}_{l-1}) \quad (2)$$

$$\mathbf{W}_l = \text{LN}(\text{MLP}(\mathbf{W}'_l) + \mathbf{W}'_l) \quad (3)$$

其中,  $\mathbf{W}_l$  表示第  $l$  层编码器中序列化的文本特征矩阵.

视觉编码器由 12 层的 ViT 块组成,将输入图像  $I \in \mathbf{R}^{H \times W \times C}$  拆分成  $N$  个  $P \times P$  大小的图像块,其中  $N = H \times W / P^2$ .每个图像块被展平并通过一个可训练的线性投影映射成  $D$  维的图像块嵌入向量  $\mathbf{v}_i$ ,并与位置嵌入

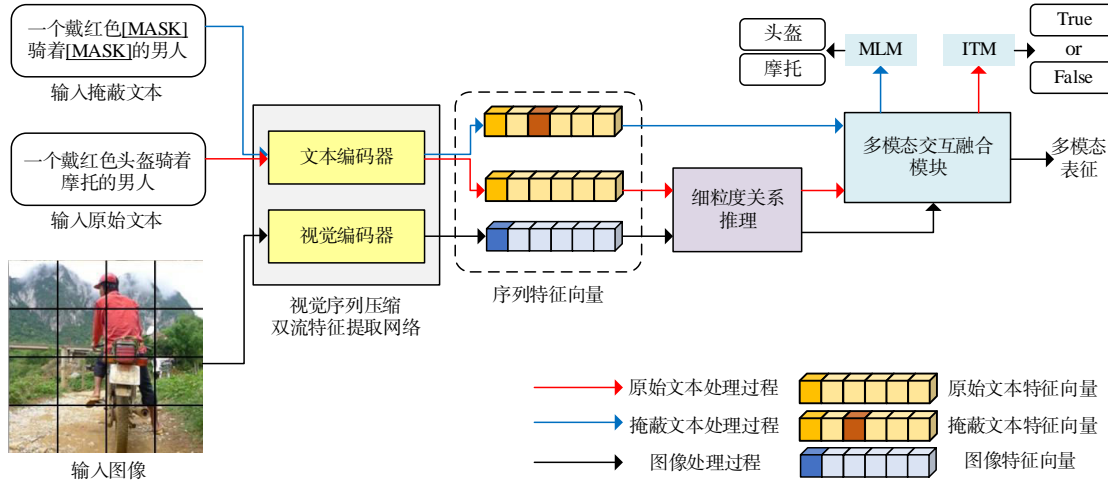


图1 MulGA模型整体框架

## 算法1 MulGA

输入: 图像  $I$ , 原始文本  $T$ , 掩蔽文本  $T_{\text{mask}}$ , 视觉标记保留比例  $\lambda$

输出: 多模态表征  $F_{\text{fusion}}$

```

1: Function MulGA( $I, T, T_{\text{mask}}, \lambda$ )
2: //基于视觉序列压缩的双流特征提取
3:  $W, W_{\text{mask}} \leftarrow \text{BERT}(T, T_{\text{mask}})$ ;
4: for  $i = 1$  to 12 do
5:   if  $i \leq 6$  then
6:      $Z_i \leftarrow \text{ViT}(I)$ ; //前六层保持不变
7:   else  $i \leq 9$  then
8:      $m_{\text{cls}} \leftarrow t_{\text{cls}} + v_{\text{cls}}$ ; //合并CLS标记
9:      $a_{\text{cls}} \leftarrow \text{MSA}(m_{\text{cls}}, Z_i)$ ;
10:    划分相关序列  $Z^{\text{R}}$  和非相关序列  $Z^{\text{NR}}$ ;
11:     $s_{x,y} \leftarrow \text{Similarity}(v_x \in Z^{\text{NR}}, v_y \in Z^{\text{R}})$ ;
12:    合并融合, 得到压缩视觉序列  $Z_{\text{out}}$ ;
13:   else
14:      $Z_i \leftarrow \text{ViT}(Z_{i-1})$ ; //后三层保持不变
15: //细粒度关系推理
16: for  $k = 1$  to  $M$  do
17:   //计算图像和文本特征的标记级相似度
18:    $S_I, S_T \leftarrow \text{Token\_Similarity}(Z^k, W^k)$ ;
19:   计算对比损失  $L_{\text{ite}} \leftarrow \text{ITC}(S_I, S_T)$ ;
20: //多模态交互融合
21:  $F_{\text{fusion}}, F_{\text{fusion}}^{\text{mask}} \leftarrow \text{MCA}(Z, W, W_{\text{mask}})$ ;
22: 执行预训练任务 ITM, MLM( $F_{\text{fusion}}, F_{\text{fusion}}^{\text{mask}}$ );
23: 计算总损失  $L_{\text{total}} \leftarrow L_{\text{ite}} + L_{\text{itm}} + L_{\text{mlm}}$ ;
24: return

```

$V_{\text{pos}}$  求和构造序列化的视觉特征矩阵  $Z_0$ , 如式(4)所示:

$$Z_0 = (v_{\text{cls}}, v_1, v_2, v_3, \dots, v_N) + V_{\text{pos}} \quad (4)$$

在输入数据中, 并不是所有的图像内容都有与之相匹配的文本描述, 有必要对视觉序列进行修剪操作, 减少与文本无关的冗余视觉信息. 但现有的主流视觉序列修

剪方法只利用了单一的视觉信息, 并不适合需要处理两种模态信息的多模态模型. 为此, 本文提出图像文本联合引导的视觉序列压缩 (image-text coordinated Visual Sequence Compression, VSC) 方法, 如图3所示.

保持视觉编码器中的前六层和后三层不变, 仅在中间三层逐层选择并压缩与文本描述不相关的视觉标记, 避免因压缩过多视觉标记导致图像结构信息丢失.

对于第  $l$  层视觉序列矩阵  $Z_l \in \mathbf{R}^{(n+1) \times D}$ , 首先将序列中图像 CLS 特征向量  $v_{\text{cls}}$  和文本编码器输出的文本 CLS 特征向量  $t_{\text{cls}}$  合并, 获得具有多模态语义的联合特征向量  $m_{\text{cls}}$ , 如式(5)所示:

$$m_{\text{cls}} = \frac{1}{2} (v_{\text{cls}} + t_{\text{cls}}) \quad (5)$$

然后将联合标记与视觉序列级联送入多头自注意力层, 计算联合标记对每个视觉标记的注意力得分  $a_{\text{cls}}$ , 作为视觉标记与文本描述相关程度的评估标准, 如式(6)所示:

$$a_{\text{cls}} = \text{softmax}\left(\frac{m_{\text{cls}} \cdot Z_l[1:]^T}{\sqrt{d}}\right) \quad (6)$$

依据  $a_{\text{cls}} = (a_1, a_2, \dots, a_n)$  和预先定义的视觉标记保留比例  $\lambda$  对特征序列进行分组. 将需要保留的高得分标记划分为相关序列  $Z^{\text{R}}$ , 可以忽视的低得分标记划分为非相关序列  $Z^{\text{NR}}$ .

对于非相关序列中的视觉标记, 直接删除会导致信息丢失, 简单融合又会生成额外的标记, 破坏视觉信息的整体性<sup>[20]</sup>. 为了避免这些问题, 本文采用一种序列压缩策略. 首先, 对非相关序列中的每一个视觉标记  $v_x = Z^{\text{NR}}[:, x]$ , 在相关标记  $v_y = Z^{\text{R}}[:, y]$  中找到其最近邻的视觉标记  $v_{\text{near}}$ .

最近邻标记计算过程见式(7)和式(8).

$$s_{x,y} = \frac{v_x^T v_y}{\|v_x\| \|v_y\|} \quad (7)$$

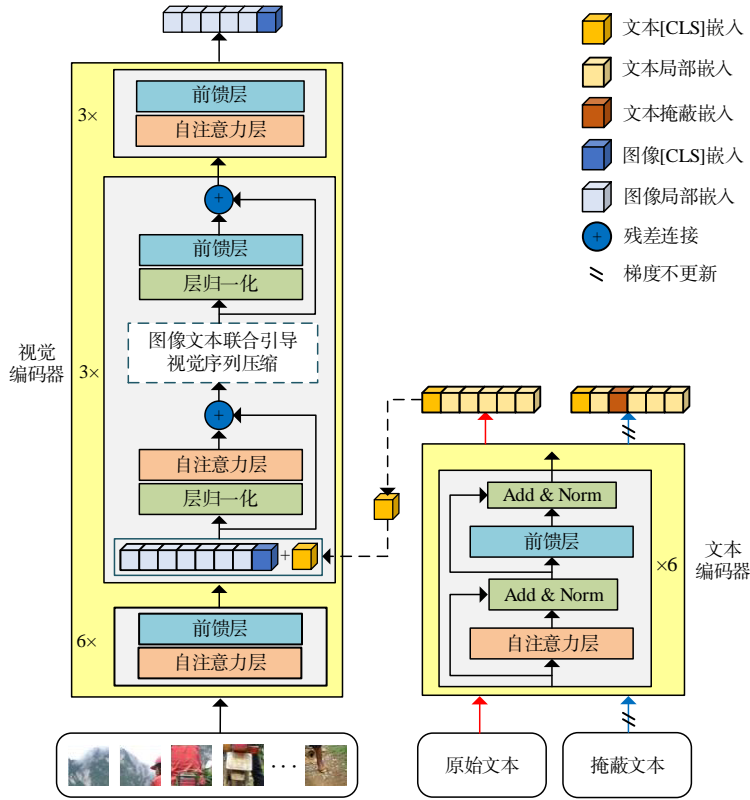


图2 双流特征提取架构示意图

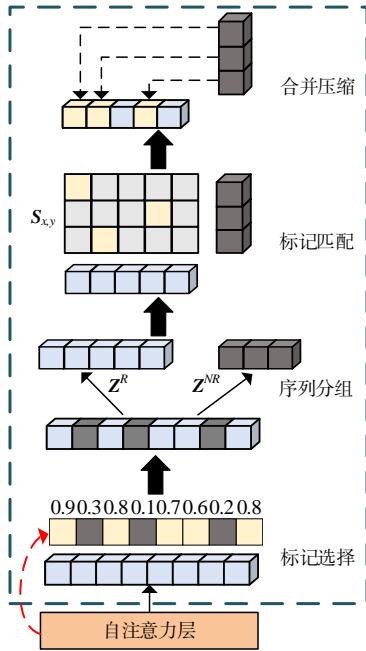


图3 图像文本联合引导视觉序列压缩示意图

$$\mathbf{v}_{\text{near}} = \arg \max_{\mathbf{v}_j \in \mathbf{Z}^R[\cdot, y]} s_{x,y} \quad (8)$$

随后,计算非相关标记与最近邻相关标记的融合权值,将非相关标记加权合并到其最近邻标记中,获得更新后的令牌  $\mathbf{v}_{\text{update}}$ . 在不损失信息也不增加额外标记

的情况下,实现视觉序列压缩,如式(9)~(11)所示:

$$\theta_x = \frac{\exp(s_{x,y})}{\sum_{\mathbf{v}_x \in \mathbf{Z}^{\text{NR}}[\cdot, x]} \exp(s_{x,y}) + e} \quad (9)$$

$$\theta_y = \frac{e}{\sum_{\mathbf{v}_x \in \mathbf{Z}^{\text{NR}}[\cdot, x]} \exp(s_{x,y}) + e} \quad (10)$$

$$\mathbf{v}_{\text{update}} = \theta_y \mathbf{v}_y + \sum_{\mathbf{v}_x \in \mathbf{Z}^{\text{NR}}[\cdot, x]} \theta_x \mathbf{v}_x \quad (11)$$

其中,  $\theta_x$  为每个非相关标记的融合权值,  $\theta_y$  为相关标记本身的融合权值.

最后,将当前层的视觉序列重构为  $\mathbf{Z}_l \in \mathbf{R}^{(k+1) \times D}$ , 并馈送到多层感知机模块进一步处理, 其中  $k = \lambda n$ .

### 2.2 基于对比学习的细粒度关系推理

细粒度关系推理模块(Fine-grained Relationship Inference Module, FRIM)利用视觉和语言表征的标记级相似度综合表示图像与文本以及文本与图像的相似度,并使用对比损失最大化正样本对的相似度、最小化负样本对的相似度,实现更细粒度和更具可解释性的跨模态对齐,具体流程如图4所示.

假设在一个训练批次中有  $M$  个图像-文本对  $\{I_k, T_k\}_{k=1}^M$ , 其中  $\{I_k, T_k\}$  为正样本对, 而  $I_k$  与其他文本组合  $\{I_k, T_m\}_{m \neq k}^M$  为负样本对. 对于该批次中的第  $k$  个图像-文本对  $\{I_k, T_k\}$ , 将其视觉特征  $\mathbf{Z}^k = (\mathbf{v}_{\text{cls}}, \mathbf{v}_1, \dots, \mathbf{v}_{n_1})$  和文

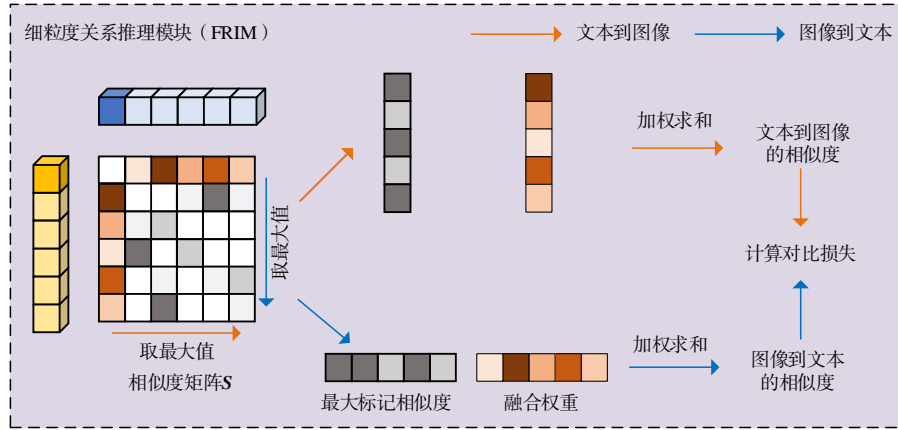


图4 细粒度关系推理模块

本特征  $\mathbf{W}^k = (t_{cls}, t_1, \dots, t_{n_2})$  经过单层感知机映射到公共嵌入空间, 用  $g_v(\mathbf{Z}^k)$  和  $g_t(\mathbf{W}^k)$  表示, 其中  $n_1$  和  $n_2$  为图像和文本的非填充标记数量. 为了实现图像与文本的细粒度对齐, 细粒度关系推理模块需要对每个特征标记都计算相似度, 最终获得一个  $(n_1 + 1) \times (n_2 + 1)$  的标记级相似性矩阵  $S$ , 如式(12)所示:

$$S = g_v(\mathbf{Z}^k)^T g_t(\mathbf{W}^k) \quad (12)$$

对于图像的第  $i$  个视觉标记  $v_i$ , 取其中除全局语义表示外  $n_2$  个相似度中最大的一个作为该视觉标记与整段文本的相似度, 如式(13)所示:

$$\max_{0 \leq j < n_2} g_v(v_i)^T g_t(t_j) \quad (13)$$

这样可以得到  $n_1$  个视觉标记与整段文本的相似度, 同理也能获得  $n_2$  个文本标记与整幅图像的相似度. 为了进一步量化  $I_k$  和  $T_k$  之间的匹配程度, 通过对全局语义表示与各标记之间的相似性进行 softmax 归一化处理, 评估各个标记在匹配过程中的贡献程度, 获得对应的融合权重  $\omega$ , 进而加权融合得到图像到文本的相似度  $S_I(\mathbf{Z}^k, \mathbf{W}^k)$  与文本到图像的相似度  $S_T(\mathbf{W}^k, \mathbf{Z}^k)$ .

$S_I(\mathbf{Z}^k, \mathbf{W}^k)$  和  $S_T(\mathbf{W}^k, \mathbf{Z}^k)$  的计算过程见式(14)和式(15).

$$S_I(\mathbf{Z}^k, \mathbf{W}^k) = \begin{cases} \sum_{i=1}^{n_1} \omega_i g_v(v_i)^T g_t(t_p), & \text{if } p = \arg \max_{0 \leq j < n_2} g_v(v_i)^T g_t(t_j) \\ 0, & \text{if } p \neq \arg \max_{0 \leq j < n_2} g_v(v_i)^T g_t(t_j) \end{cases} \quad (14)$$

$$S_T(\mathbf{W}^k, \mathbf{Z}^k) = \begin{cases} \sum_{j=1}^{n_2} \omega_j g_t(t_j)^T g_v(v_p), & \text{if } p = \arg \max_{0 \leq i < n_1} g_t(t_j)^T g_v(v_i) \\ 0, & \text{if } p \neq \arg \max_{0 \leq i < n_1} g_t(t_j)^T g_v(v_i) \end{cases} \quad (15)$$

需要注意的是, 与图文对比学习不同, 此时的  $S_I(\mathbf{Z}^k, \mathbf{W}^k)$  不一定等于  $S_T(\mathbf{W}^k, \mathbf{Z}^k)$ .

最后计算图文对比损失, 实现图像和文本的细粒度交互和对齐, 如式(16)~(18)所示:

$$L_k^{I2I} = -\frac{1}{M} \log \frac{\exp(S_I(\mathbf{Z}^k, \mathbf{W}^k)/\tau)}{\sum_{m=1}^M \exp(S_I(\mathbf{Z}^k, \mathbf{W}^m)/\tau)} \quad (16)$$

$$L_k^{T2I} = -\frac{1}{M} \log \frac{\exp(S_T(\mathbf{W}^k, \mathbf{Z}^k)/\tau)}{\sum_{m=1}^M \exp(S_T(\mathbf{W}^k, \mathbf{Z}^m)/\tau)} \quad (17)$$

$$L_{itc} = \frac{1}{2} \sum_{k=1}^M (L_k^{I2I} + L_k^{T2I}) \quad (18)$$

其中,  $\tau$  是可学习的温度参数,  $L_k^{I2I}$  为图像  $I_k$  到文本的对比损失,  $L_k^{T2I}$  为文本  $T_k$  到图像的对比损失,  $L_{itc}$  为该训练批次的总对比损失.

### 2.3 多模态交互融合模块

多模态交互融合模块如图5所示, 输入对齐后图像特征、原始文本特征以及掩蔽文本特征, 其中掩蔽文本特征和原始文本特征经过多头自注意力层作为查询向量, 图像特征作为键向量和值向量, 在多头交叉注意力层 (Multi-head Cross-Attention, MCA) 进行融合. 通过图像文本匹配 (ITM) 和掩蔽语言建模 (MLM) 两个预训练任务训练融合编码器隐式挖掘图像和文本的细粒度交互关系, 如式(19)~(22)所示:

$$\mathbf{M}_0 = \mathbf{W} \quad (19)$$

$$\mathbf{M}_l'' = \text{LN}(\text{MSA}(\mathbf{M}_{l-1})) + \mathbf{M}_{l-1} \quad (20)$$

$$\mathbf{M}_l' = \text{LN}(\text{MCA}(\mathbf{M}_l'', \mathbf{Z})) + \mathbf{M}_l'' \quad (21)$$

$$\mathbf{M}_l = \text{LN}(\text{MLP}(\mathbf{M}_l')) + \mathbf{M}_l' \quad (22)$$

其中,  $\mathbf{W}$  是文本编码器的输出,  $\mathbf{Z}$  是视觉编码器的输出,  $\mathbf{M}_l$  是第  $l$  层编码器中序列化的多模态特征矩阵.

图像文本匹配 (ITM) 任务将融合编码器输出的文本 CLS 嵌入作为多模态全局语义表示, 经过一个全连

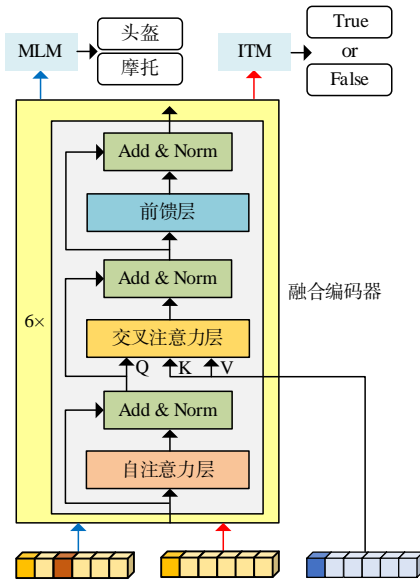


图5 多模态交互融合模块

接层和 Softmax 来预测图像和文本是否匹配. ITM 损失计算如式(23)所示:

$$L_{itm} = E_{(I,T) \sim M} (H(y^{itm}, p^{itm}(I, T))) \quad (23)$$

其中,  $p^{itm}$  为预测两类的概率,  $y^{itm}$  表示真值标签,  $y^{itm}$  当图像和文本相匹配时为 1, 不匹配时为 0.

掩蔽语言建模 (MLM) 任务利用图像信息和文本上下文信息预测被随机掩蔽的文本标记. 由于提高 MLM 任务中文本的掩蔽比例, 能够帮助模型更有效地利用视觉特征<sup>[22,23]</sup>. 因此本文以 30% 的概率随机掩蔽原始文本. MLM 损失计算如式(24)所示:

$$L_{mlm} = E_{(I, T_{mask}) \sim M} (H(y^{msk}, p^{msk}(I, T_{mask}))) \quad (24)$$

其中  $T_{mask}$  表示一个被掩蔽的文本标记,  $p^{msk}(I, T_{mask})$  表示模型对被掩蔽标记的预测概率,  $y^{msk}$  表示词汇概率分布.

MulGA 的预训练总损失可以表示为

$$L_{total} = \alpha L_{itc} + \beta L_{itm} + \gamma L_{mlm} \quad (25)$$

### 3 实验分析

本节详细介绍了 MulGA 在视觉语言检索和分类两种下游任务中的实验测试结果, 并与主流的预训练方法进行了全面比较, 展示了本文模型在各个任务上的良好表现.

#### 3.1 实验细节

MulGA 使用预训练的 ViT-B/16 作为视觉编码器, 将 BERT<sup>[24]</sup> 的前六层作为文本编码器, 后六层作为多模态融合编码器, 视觉标记保留比例  $\lambda$  设置为 70%. 在 4 张 NVIDIA Tesla V100 32 GB GPU 上预训练了 30 个周期. 为了公平比较, 除了由于计算资源受限, 降低了批次大小之外, 其余实验设置基本与 ALBEF 保持一致.

#### 3.2 预训练数据集

本文采用 MSCOCO<sup>[25]</sup> (Microsoft Common Objects in Context)、VG<sup>[26]</sup> (Visual Genome)、SBU<sup>[27]</sup> (SBU Captions) 和 CC3M<sup>[28]</sup> (Conceptual Caption) 4 个可公开访问的数据集作为基础数据集. 由于四个数据集中的部分图像数据在下载前已缺失, 最终 4 M 大小的基础数据集中, 图像总数约为 400 万张, 图像-文本对的数量约为 510 万套.

此外, 本文以随机采样的方式从 CC12 M 数据集中筛选数据, 构造了一个包含 6 M 图像文本对的数据集, 与采用大规模数据集预训练的模型进行对比实验.

#### 3.3 视觉语言检索任务评估

视觉语言检索任务包括图像到文本检索 (Image Retrieval, IR) 和文本到图像检索 (Text Retrieval, TR). 对于此项任务, 本文使用召回率 R@k (Recall@k)、平均精确率 mAP (Average Precision) 和  $F_1$  分数对模型进行全面评估.

##### 3.3.1 图文检索微调实验

MulGA 在 MSCOCO 和 Flickr30k<sup>[29]</sup> 上图文检索微调实验的 R@k 的评估结果如表 1 所示. 其中第二列代表预训练数据集中的图像数目, “—”表示该模型结果未报告.

在使用 4 M 公共预训练数据集的条件下, 与 ALBEF 模型相比, MulGA 在 MSCOCO 和 Flickr30 k 上的平均召回率分别提升了 2.9% 和 1.2%. 并在 Flickr30 k 上的文本检索召回率达到了最优水平. 当使用 6 M 预训练数据集时, MulGA 展现了与 BLIP (Bootstrapping Language-Image Pre-training) 相当的性能, 在 Flickr30 k 上的 TR@5 和 TR@10 达到 100, 相比于在更大数据集上预训练的 ALBEF, 在 MSCOCO 数据集的 TR@1 指标上实现了 1.9% 的提升.

微调检索实验的 mAP 与  $F_1$  分数评估结果如表 2 所示, 其中 mAP 为模型在 1、5、10 不同阈值下的精确率均值. 在 4 M 数据条件下, MulGA 在 MSCOCO 数据集上的 mAP 优于 TCL (Triple Contrastive Learning), 略低于使用 14 M 数据预训练的 ALBEF 和 BLIP, 在 Flickr30 k 上 mAP 和  $F_1$  分数仅次于 BLIP; 当预训练数据扩充到 6 M 时, mAP 增长了 1.14%, 模型性能超过 ALBEF, 与 BLIP 模型相近.

为了验证实验结果的可靠性, 本文基于 Wilcoxon 符号秩检验法, 以召回率作为各个模型的样本数据, 分析 MulGA 在图文检索微调任务中相较于其他模型的提升是否显著, 各个配对样本的 P 值以热力图的形式展现, 如图 6 所示.

分析结果表明, 本文方法与 UNITER (UNiversal Image-Text Representation)、ALBEF、TCL 等模型相比,

表1 图文检索微调实验-召回率

单位:%

模型	数据集	MSCOCO						Flickr30 k					
		Text Retrieval			Image Retrieval			Text Retrieval			Image Retrieval		
	图像	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
UNITER <sup>[6]</sup>	4 M	65.7	88.6	93.8	59.2	79.9	88.0	87.3	98	99.2	75.6	94.1	96.8
ImageBERT <sup>[30]</sup>	6 M	66.4	89.8	94.4	50.5	78.7	87.1	87.0	97.6	99.2	73.1	92.6	96.0
OSCAR <sup>[7]</sup>	4 M	70.0	91.1	95.5	54.0	80.8	88.5	—	—	—	—	—	—
ALBEF <sup>[14]</sup>	4 M	73.1	91.4	96.0	56.8	81.5	89.2	94.3	99.4	99.8	82.8	96.7	98.4
METER <sup>[31]</sup>	4 M	76.2	93.2	96.8	57.1	82.7	90.1	94.3	99.6	99.9	82.2	96.3	98.4
TCL <sup>[32]</sup>	4 M	75.6	92.8	96.7	59.0	83.2	89.9	94.9	99.5	99.8	84.0	96.7	98.5
ViLTA <sup>[33]</sup>	4 M	73.3	91.8	95.9	59.5	83.1	89.7	94.5	99.8	99.8	85.2	97.2	98.8
MaskVLM <sup>[34]</sup>	4 M	76.3	93.8	96.8	60.1	83.6	90.4	95.3	99.8	100	84.9	97.4	98.6
VL-Match <sup>[35]</sup>	4 M	76.6	93.8	97.1	60.2	83.6	90.1	96.4	99.8	100	86.0	97.5	99.0
ALIGN <sup>[9]</sup>	1.8 B	77.0	93.5	96.9	59.9	83.3	89.8	95.3	99.8	100	84.9	97.4	98.6
ALBEF <sup>[14]</sup>	14 M	77.6	94.3	97.2	60.7	84.3	90.5	95.9	99.8	100	85.6	97.5	98.9
BLIP <sup>[36]</sup>	14 M	80.6	95.2	97.6	63.1	85.3	91.1	96.6	99.8	100	87.2	97.5	98.8
本文方法	4 M	77.1	93.6	97.1	60.0	83.9	90.5	96.4	99.9	100	85.6	97.5	98.9
本文方法	6 M	79.1	95.2	97.8	60.9	84.3	91.3	97.2	100	100	86.5	97.0	98.4

表2 图文检索微调实验-精确率、F<sub>1</sub>分数

单位:%

模型	数据集	MSCOCO				Flickr30 k			
		T-mAP	T-F <sub>1</sub>	I-mAP	I-F <sub>1</sub>	T-mAP	T-F <sub>1</sub>	I-mAP	I-F <sub>1</sub>
TCL <sup>[32]</sup>	4 M	34.60	41.36	28.20	34.35	41.59	48.73	37.73	44.71
ALBEF <sup>[14]</sup>	14 M	35.39	42.23	28.86	35.08	41.59	49.11	38.31	45.34
BLIP <sup>[36]</sup>	14 M	36.47	43.34	29.76	36.03	42.39	49.56	38.90	45.93
本文方法	4 M	35.18	41.99	28.61	34.82	42.13	49.29	38.33	45.36
本文方法	6 M	35.97	42.87	28.96	35.22	42.43	49.57	38.58	45.58

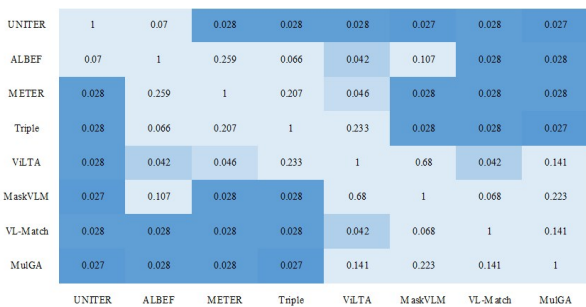


图6 微调检索实验统计显著性验证

图文检索具有显著性提升,在 ViLTA (Vision-Language pretraining through Textual Augmentation) 和 MaskVLM (Masked Vision and Language Modeling) 等模型上,提升并不明显.分析结果基本与 R@k、mAP 和 F<sub>1</sub> 分数三个评估指标的评估结果一致,证明本文实验结果具有一定可靠性.

### 3.3.2 图文检索零样本实验

为了评估 MulGA 的零样本迁移能力,本文将在 MSCOCO 数据集上微调的模型直接迁移到 Flickr30 k 数

据集上进行了图文检索零样本实验.表3展示了 MulGA 在 Flickr30 k 上零样本检索的评估结果.

表3 图文检索零样本实验-召回率

单位:%

模型	数据集	Flickr30 k					
		图像	TR@1	TR@5	TR@10	IR@1	IR@5
ImageBERT <sup>[30]</sup>	6 M	70.7	90.2	94.0	54.3	79.6	87.5
UNITER <sup>[6]</sup>	4 M	83.6	95.7	97.7	68.7	89.2	93.9
MaskVLM <sup>[34]</sup>	4 M	87.0	97.9	99.3	75.0	92.5	95.8
CLIP <sup>[8]</sup>	400 M	88.0	98.7	99.4	68.7	90.6	95.2
ALIGN <sup>[9]</sup>	1.2 B	88.6	98.7	99.7	75.7	93.8	96.8
ALBEF <sup>[14]</sup>	4 M	90.5	98.8	99.7	76.8	93.7	96.7
METER <sup>[31]</sup>	4 M	90.9	98.3	99.5	79.6	95.0	97.3
TCL <sup>[32]</sup>	4 M	93.0	99.1	99.6	79.6	95.1	97.4
VL-Match <sup>[35]</sup>	4 M	93.3	99.3	99.8	82.0	95.1	97.4
ALBEF <sup>[14]</sup>	14 M	94.1	99.5	99.7	82.8	96.3	98.1
BLIP <sup>[36]</sup>	14 M	94.8	99.7	100	84.9	96.8	98.6
本文方法	4 M	91.1	99.1	99.8	80.2	95.3	97.3
本文方法	6 M	93.2	99.6	99.9	82.8	96.0	98.1

相较于 ImageBERT、CLIP 和 ALIGN, MulGA 用更少的训练数据获得了更好的性能. 与同等训练数据的模型相比, 零样本检索性能也优于 ALBEF、METER (Multi-modal End-to-end Transformer) 等模型. 但对比不同数据体量下的两种最优算法 VL-Match 和 BLIP, MulGA 在 TR@1 和 IR@1 上的表现仍有较大差距, 推测可能是固定的视觉压缩比例导致在进行零样本迁移时无法准确压缩图像冗余信息, 一定程度的主体结构信息被破坏, 最终影响了模型性能.

零样本检索实验的 mAP 与  $F_1$  分数评估结果在表 4 中显示. MulGA 零样本迁移表现与 TCL 模型相当. 将预训练数据集扩充到 6 M 大小时, 模型的零样本检索性能会有显著提升. 这一结果表明, 通过扩大预训练数据集规模, 可以有效增强模型在处理零样本迁移任务时的泛化能力.

表 4 图文检索零样本实验-精确率、 $F_1$  分数 单位: %

模型	数据集	Flickr30 k			
		图像	T-mAP	T- $F_1$	I-mAP
TCL <sup>[32]</sup>	4 M	40.93	48.04	36.13	42.83
ALBEF <sup>[14]</sup>	14 M	41.32	48.47	37.30	44.26
BLIP <sup>[36]</sup>	14 M	41.58	48.74	38.04	45.03
本文方法	4 M	40.30	47.43	36.33	43.17
本文方法	6 M	41.04	48.19	37.27	44.21

为了验证零样本实验结果的可靠性, 本文对实验数据进行了显著性分析, 统计显著性验证结果如图 7 所示.

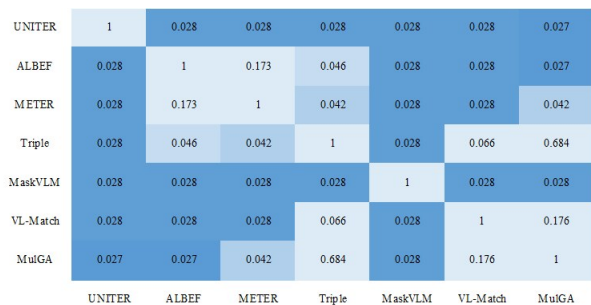


图 7 零样本检索实验统计显著性验证

从图中结果可以得出结论, MulGA 在零样本检索实验中与 UNITER、ALBEF、MaskVLM 相比, 具有显著的性能提升, 但与 METER、VL-Match (enhancing Vision-Language with token-level Matching) 模型相比提升不具备显著性, 总体表现与 TCL 模型相近.

### 3.4 视觉语言分类任务评估

为了进一步验证 MulGA 的性能, 本文在视觉问答、视觉推理和视觉蕴含三个视觉语言分类任务上进行了实验.

#### 3.4.1 视觉问答实验

视觉问答 (Visual Question Answering, VQA) 任务要求模型根据给定的图像和问题, 生成或从答案集中选择出正确的答案. 本文在 VQAv2<sup>[37]</sup> 数据集上微调 and 评估 MulGA 的视觉问答能力. 评估结果在表 5 中展示.

表 5 视觉问答和视觉推理实验结果 单位: %

模型	数据集 图像	VQAv2		NLVR2	
		test-dev	test-std	dev	test
UNITER <sup>[6]</sup>	4 M	72.7	72.9	77.2	77.9
GLIPv2 <sup>[38]</sup>	20 M	73.1	73.3	—	—
OSCAR <sup>[7]</sup>	4 M	73.2	73.4	78.1	78.4
ALBEF <sup>[14]</sup>	4 M	74.5	74.7	80.2	80.5
GRIT-VLP <sup>[39]</sup>	4 M	75.1	75.3	80.7	81.6
VL-Match <sup>[35]</sup>	4 M	75.1	75.2	82.0	82.2
MaskVLM <sup>[34]</sup>	4 M	75.5	75.4	81.6	82.0
ALBEF <sup>[14]</sup>	14 M	75.8	76.0	82.6	83.1
METER <sup>[31]</sup>	4 M	76.4	76.4	82.4	82.2
BLIP <sup>[36]</sup>	14 M	77.5	77.6	82.7	82.3
SimVLM <sup>[40]</sup>	1.8 B	77.9	78.1	81.7	81.8
OFA <sup>[41]</sup>	54 M	78.0	78.1	—	—
本文方法	4 M	76.9	76.6	82.3	82.7
本文方法	6 M	77.6	77.9	82.7	82.9

基于 4 M 数据集, MulGA 在 Test-Dev 和 Test-Std 上的准确率优于所有同等训练体量的对比算法, 甚至超越了使用更大数据集进行预训练的 GLIPv2 (Grounded Language-Image Pre-training v2), 其中与基础的 ALBEF 模型相比, 准确率分别提升了 3.22% 和 2.54%. 在 6 M 预训练数据的条件下, MulGA 在视觉问答下游任务中的平均准确率达到了 77.75%.

#### 3.4.2 视觉推理实验

视觉推理 (Natural Language for Visual Reasoning, NLVR) 任务侧重于预测文本标题是否与一对图像的内容相符. 本文在 NLVR2<sup>[42]</sup> 数据集上进行微调实验, 并使用准确率进行评估.

如表 5 所示, MulGA 在视觉推理任务上的表现超越了所有对比算法. 使用 4 M 数据预训练的 MulGA, 在 NLVR2 的 test 测试集上准确率达到了最佳水平. 与大模型相比, 使用 6 M 数据预训练的 MulGA 表现出了与 ALBEF 和 BLIP 相当的性能, 在 dev 和 test 测试集上分别比 SimVLM (Simple Visual Language Model) 提升了 1.22% 和 1.36%.

#### 3.4.3 视觉蕴含实验

视觉蕴含 (Visual Entailment, VE) 任务是一个判断一幅图像和一段文本之间的关系是蕴含的、中性的还是矛盾的三分类问题. 本文在 SNLI-VE<sup>[43]</sup> 数据集上微调模型, 并评估分类的准确率, 结果如表 6 所示.

表6 视觉蕴含实验结果 单位:%

模型	数据集	VE	
		Test	Val
UNITER <sup>[6]</sup>	4 M	78.3	78.6
UniVL <sup>[44]</sup>	2.84 M	79.7	80.0
ALBEF <sup>[14]</sup>	4 M	80.3	80.1
MaskVLM <sup>[34]</sup>	4 M	80.4	80.7
ALBEF <sup>[14]</sup>	14 M	80.9	80.8
METER <sup>[31]</sup>	4 M	81.2	80.9
VL-Match <sup>[35]</sup>	4 M	81.3	80.4
ViLTA <sup>[33]</sup>	4 M	81.7	81.5
本文方法	4 M	81.2	80.8
本文方法	6 M	81.5	81.6

在4 M数据的条件下, MulGA在test和val上的平均准确率达到81.0%, 相较于具有相似网络架构和训练范式的ALBEF模型, 实现了0.99%的性能提升. 使用6 M数据训练的MulGA在val上取得81.6%的最佳表现.

### 3.4.4 分类任务显著性分析

为了分析MulGA在视觉语言分类任务上是否具有显著性提升, 本文将视觉问答、视觉推理和视觉蕴含三个下游任务的测试结果整合作为样本数据, 对视觉语

言分类任务的数据结果进行了统计显著性验证测试, 检验结果如图8所示.

UNITER	1	0.027	0.027	0.028	0.028	0.028	0.027	0.027
ALBEF-4M	0.027	1	0.116	0.027	0.058	0.043	0.066	0.028
MaskVLM	0.027	0.116	1	0.027	0.752	0.027	0.028	0.028
METER	0.028	0.027	0.027	1	0.08	0.221	0.6	0.027
VL-Match	0.028	0.058	0.752	0.08	1	0.046	0.058	0.027
MulGA-4M	0.028	0.043	0.027	0.221	0.046	1	0.416	0.028
ALBEF-14M	0.027	0.066	0.028	0.6	0.058	0.416	1	0.075
MulGA-6M	0.027	0.028	0.028	0.027	0.027	0.028	0.075	1

图8 视觉语言分类统计显著性验证

MulGA在4 M训练数据条件下, 与其他同体量模型都满足显著性差异条件, 当数据扩充到6 M大小时, MulGA相较于4 M的基础模型P值为0.028. 证明适当扩充预训练数据集规模, 能够显著提升模型在下游任务上的性能表现.

### 3.5 细粒度匹配分析

本文通过对融合编码器中交叉注意力层的注意力进行可视化, 检验MulGA模型的细粒度匹配能力. 图9展示了MulGA与ALBEF交叉注意力可视化的对比结果, 其中左侧为ALBEF, 右侧为MulGA.



图9 细粒度匹配对比结果

结果表明, ALBEF模型在进行微调后虽然能够正确匹配图像对象和对应的文本字词, 但会受到图像背

景信息的干扰, 注意力权重较为分散. 而MulGA经过视觉序列压缩和细粒度关系推理后, 更好地学习了图

像和文本细粒度对齐关系,对实义单词所指代的具体图像实例聚焦更为准确,在准确识别图像对象、属性、动作和数量等细粒度特征方面具有优越性.

### 3.6 消融实验

#### 3.6.1 标记保留比例分析

在视觉序列压缩过程中,每层的视觉标记保留比例  $\lambda$  是在训练前预设的超参数.为了研究不同标记保留比例对于视觉序列压缩的影响,本文使用不同的视觉标记保留  $\lambda$  比例预训练 5 个周期,并在 Flickr30 k 上

比较零样本图文检索的性能变化,如图 10 所示.

随着视觉标记保留比例  $\lambda$  的增加,模型在图文检索任务上的总体效果表现为先上升后下降的趋势,当  $\lambda$  减小到 60% 时,召回率大幅下降,推测是由于视觉标记保留比例过低,图像的大部分信息被压缩,破坏了图像主体信息的完整性,从而影响检索性能.将  $\lambda$  增加到 90% 时,视觉序列压缩对模型性能提升不大.而当  $\lambda$  为 70% 时,图文检索各指标均优于其他测试值,因此,本文设定每层视觉标记保留比例为 70%.

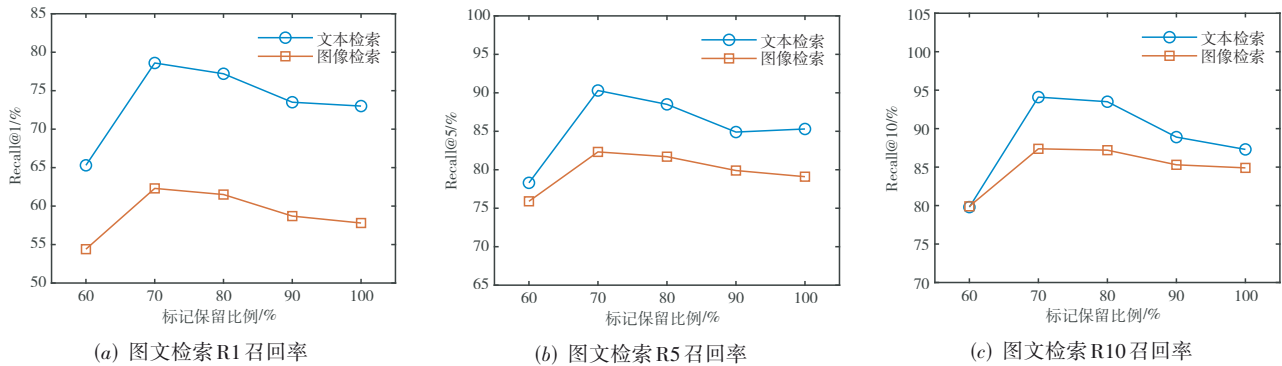


图 10 不同保留比例的图文检索性能

#### 3.6.2 图像文本联合引导分析

为了分析视觉序列压缩方法中,图像文本联合引导对模型有效性的影响,本文在视觉问答、视觉推理、视觉蕴含三个下游任务上,分别测试了仅使用视觉信息引导和仅文本信息引导的视觉序列压缩效果.实验结果在表 7 中展示.

表 7 不同模态信息指导下的实验结果 单位:%

模态信息	VQAv2	NLVR2	VE
文本信息	74.2	80.4	80.6
视觉信息	73.3	78.1	80.0
视觉信息+文本信息	75.9	81.7	80.4

与使用融合了两种模态的多模态信息引导视觉序列压缩相比,无论是视觉信息还是文本信息,仅依赖单一模态信息会产生较差的结果.因此,使用图像信息和文本信息来联合指导视觉编码器压缩视觉序列中的冗余信息,在视觉语言多模态学习中更具有有效性.

#### 3.6.3 损失函数权重分析

为了研究 ITC、ITM 和 MLM 三个损失对模型的贡献程度,本文以 1:1:0.8、1:0.8:1、0.8:1:1 和 1:1:1 四种不同的损失函数权重,分别进行 5 个周期的预训练,并在 Flickr30k 数据集上做零样本图文检索实验,分析三种损失对模型的影响,实验结果如表 8 所示.

对比第一行、第三行和第四行数据,当降低 ITC 损失和 MLM 损失比重时,对模型零样本检索的性能的影响

较大,图像检索和文本检索的召回率都有明显下降;对比第二行与第四行数据,当降低 ITM 损失比重时,TR@1 和 TR@5 分别降低了 1.78% 和 2.1%,IR@5 降低了 0.7%.分析结果证明,三个损失都会对模型作出积极贡献.

#### 3.6.4 视觉序列压缩与细粒度关系推理消融

本文比较了不添加任何组件、单 FRIM 组件、单 VSC 组件和 VSC+FRIM 组合 4 种不同情况下模型的微调检索性能.同时,在单张 GPU 上计算一个训练批次数据前向传播的浮点运算次数 (Floating Point Operations, FLOPs),实验结果在表 9 中展示.需要注意的是,为了保证先对齐再融合的统一训练范式,单 VSC 组件测试并没有直接删除 FRIM,而是使用原始的 ITC 进行替换.

表 8 损失函数不同权重占比下的实验结果 单位:%

$\alpha:\beta:\gamma$	TR@1	TR@5	TR@10	IR@1	IR@5	IR@10
1:1:0.8	73.0	84.7	90.4	60.5	80.2	85.3
1:0.8:1	77.2	88.4	94.3	62.3	81.7	87.7
0.8:1:1	73.3	84.9	88.9	58.7	75.9	80.9
1:1:1	78.6	90.3	94.1	62.3	82.3	87.4

表 9 在 Flickr30 k 上使用不同模块组合进行预训练的图文检索评估

ITC	VSC	FRIM	FLOPs/G	TR@1/%	TR@5/%	IR@1/%	IR@5/%
√			988.6	88.1	95.5	74.8	92.9
		√	1 031.3	91.2	98.5	76.9	94.9
√	√		917.7	89.2	96.4	76.3	94.1
	√	√	931.9	91.6	99.1	79.1	95.7

第一行与第二行相比,将 ITC 更换为 FRIM,虽然增加了少量计算量,但图文检索的召回率都获得了提升;第三行在第一行的基础上添加了 VSC,显著降低了模型的计算复杂度,在提升模型效率的同时改善了检索性能;第四行组合了 VSC 与 FRIM,与第三行相比,计算成本略微增加,在文本检索和图像检索的性能分别提升了 2.81% 和 2.6%,实验结果表明,本文提出的 VSC 和

FRIM 两个关键组件能够有效提升多模态模型的性能和效率.

### 3.7 可视化

为了验证本文提出的图像文本联合引导的视觉序列压缩方法的有效性,从 MSCOCO 数据集中随机采样了部分数据作为样本输入,可视化视觉编码器中 3 个视觉序列压缩层所选择的相关视觉标记,如图 11 所示.

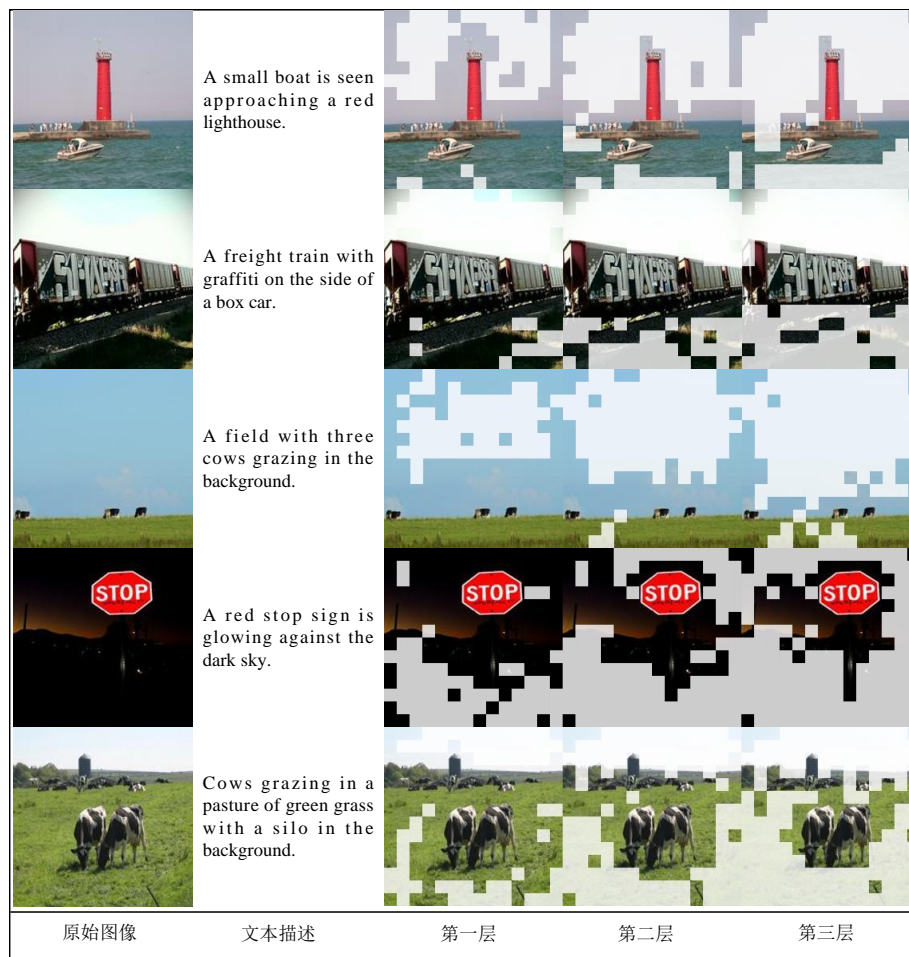


图 11 相关视觉标记选择可视化

随着层数增加,视觉相关标记逐渐减少,图像中与文本描述不相关的背景信息被逐渐过滤,模型更加关注文本所描述的主题内容.证明本文的视觉序列压缩方法能够在图像信息和文本信息的联合指导下,有效地减少视觉序列中的冗余信息,保留与文本描述相关的主题内容.

## 4 结论

本文提出了一种基于跨模态引导和对齐的视觉语言多模态预训练方法,成功解决了现有模型在模态间缺乏显式细粒度对齐机制的问题.通过引入视觉序列压缩模块和细粒度关系推理模块,显著提升了模型在

多个视觉语言任务中的性能,尤其在图文检索和视觉问答任务中表现出色.该方法在增强模态间细粒度特征对齐方面展现出强大能力,为多模态预训练模型的发展提供了新的思路.

### 参考文献

- [1] ABDU S A, YOUSEF A H, SALEM A. Multimodal video sentiment analysis using deep learning approaches, a survey[J]. Information Fusion, 2021, 76: 204-226.
- [2] ACOSTA J N, FALCONE G J, RAJPURKAR P, et al. Multimodal biomedical AI[J]. Nature Medicine, 2022, 28 (9): 1773-1784.

- [3] 樊琳, 龚勋, 郑岑洋. 基于文本引导下的多模态医学图像分析算法[J]. 电子学报, 2024, 52(7): 2498-2512.  
FAN L, GONG X, ZHENG C Y. A Multi-modal medical image analysis algorithm based on text guidance[J]. Acta Electronica Sinica, 2024, 52(7): 2498-2512. (in Chinese)
- [4] HUANG K L, SHI B T, LI X, et al. Multi-modal sensor fusion for auto driving perception: a survey[EB/OL]. (2022-02-06) [2024-03-25]. <https://doi.org/10.48550/arXiv.2202.02703>.
- [5] TAN H, BANSAL M. LXMERT: Learning cross-modality encoder representations from Transformers[C]//Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). Stroudsburg: Association for Computational Linguistics, 2019: 5100-5111.
- [6] CHEN Y C, LI L J, YU L C, et al. UNITER: Universal image-text representation learning[M]//Computer Vision — ECCV 2020. Cham: Springer International Publishing, 2020: 104-120.
- [7] LI X J, YIN X, LI C Y, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks[M]//Computer Vision — ECCV 2020. Cham: Springer International Publishing, 2020: 121-137.
- [8] RADFORD A, KIM J W, HALLACY C, et al. Learning transferable visual models from natural language supervision[EB/OL]. (2021-02-26) [2024-03-25]. <http://arxiv.org/abs/2103.00020>.
- [9] JIA C, YANG Y F, XIA Y, et al. Scaling up visual and vision-language representation learning with noisy text supervision[C]//Proceedings of the 38th International Conference on Machine Learning. New York: PMLR, 2021, 139: 4904-4916.
- [10] YU J H, WANG Z R, VASUDEVAN V, et al. Coca: Contrastive captioners are image-text foundation models[EB/OL]. (2022-06-14) [2024-03-25]. <https://doi.org/10.48550/arXiv.2205.01917>.
- [11] BAO H B, WANG W H, DONG L, et al. VLMO: Unified vision-language pre-training with mixture-of-modality-experts[C]//Proceedings of the 36th International Conference on Neural Information Processing Systems (NeurIPS). New York: Curran Associates, 2022, 35: 32897-32912.
- [12] 李志欣, 凌锋, 张灿龙, 等. 融合两级相似度的跨媒体图像文本检索[J]. 电子学报, 2021, 49(2): 268-274.  
LI Z X, LING F, ZHANG C L, et al. Cross-media image-text retrieval with two level similarity[J]. Acta Electronica Sinica, 2021, 49(2): 268-274. (in Chinese)
- [13] YAO L W, HUANG R H, HOU L, et al. FILIP: fine-grained interactive language-image pre-training[C/OL]//The Tenth International Conference on Learning Representations. (2022-01-29) [2024-03-25]. <https://openreview.net/forum?id=cpDhcsEDC2>.
- [14] LI J N, SELVARAJU R, GOTMARE A, et al. Align before fuse: Vision and language representation learning with momentum distillation[C]//Proceedings of the 35th International Conference on Neural Information Processing Systems (NeurIPS). New York: Curran Associates, 2021, 34: 9694-9705.
- [15] DOSOVITSKIY A, BEYER L, KOLESNIKOV A, et al. An image is worth 16x16 words: Transformers for image recognition at scale[EB/OL]. (2020-10-22) [2024-03-25]. <https://doi.org/10.48550/arXiv.2010.11929>.
- [16] LI C L, XU H Y, TIAN J F, et al. mPLUG: Effective and efficient vision-language learning by cross-modal skip-connections[C]//Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Stroudsburg: Association for Computational Linguistics, 2022: 7241-7259.
- [17] TANG Y H, HAN K, WANG Y H, et al. Patch slimming for efficient vision Transformers[C]//2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway: IEEE, 2022: 12165-12174.
- [18] RAO Y M, ZHAO W L, LIU B L, et al. DynamicViT: Efficient vision Transformers with dynamic token sparsification[C]//Proceedings of the 35th International Conference on Neural Information Processing Systems (NeurIPS). New York: Curran Associates, 2021, 34: 13937-13949.
- [19] LIANG Y W, GE C J, TONG Z, et al. EViT: Expediting vision Transformers via token reorganizations[C/OL]//The Tenth International Conference on Learning Representations. (2022-01-29) [2024-03-25]. [https://openreview.net/forum?id=BjyvwnXXVn\\_](https://openreview.net/forum?id=BjyvwnXXVn_).
- [20] WEI S Y, YE T Z, ZHANG S, et al. Joint token pruning and squeezing towards more aggressive compression of vision Transformers[C]//2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway: IEEE, 2023: 2092-2101.
- [21] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[C]//Proceedings of the 31st International Conference on Neural Information Processing Sys-

- tems (NeurIPS). New York: Curran Associates, 2017: 6000-6010.
- [22] BAO H B, DONG L, PIAO S H, et al. BEIT: BERT pre-training of image Transformers[C//The Tenth International Conference on Learning Representations. (2022-01-29) [2024-03-25]. <https://openreview.net/forum?id=p-BhZSz59o4>.
- [23] 汤嘉, 郭燕, 叶名玮, 等. 面向多视角对比学习和语义增强的多模态预训练方法[J]. 计算机科学, 2024, 51(1): 168-174.  
TANG J, GUO Y, YE M W, et al. Multimodal pre-training method for multi-view contrastive learning and semantic enhancement[J]. Computer Science, 2024, 51(1): 168-174. (in Chinese)
- [24] DEVLIN J, CHANG M W, LEE K, et al. BERT: Pre-training of deep bidirectional Transformers for language understanding[C//Proceedings of NAACL-HLT. Stroudsburg: ACL, 2019: 4171-4186.
- [25] LIN T Y, MAIRE M, BELONGIE S, et al. Microsoft COCO: Common objects in context[C//Computer Vision — ECCV 2014. Cham: Springer International Publishing, 2014: 740-755.
- [26] KRISHNA R, ZHU Y K, GROTH O, et al. Visual genome: Connecting language and vision using crowd-sourced dense image annotations[J]. International Journal of Computer Vision, 2017, 123(1): 32-73.
- [27] ORDONEZ V, KULKARNI G, BERG T L. Im2Text: Describing images using 1 million captioned photographs [C//Proceedings of the 25th International Conference on Neural Information Processing Systems. New York: Curran Associates, 2011, 24: 1143-1151.
- [28] SHARMA P, DING N, GOODMAN S, et al. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning[C//Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Stroudsburg: Association for Computational Linguistics, 2018: 2556-2565.
- [29] YOUNG P, LAI A, HODOSH M, et al. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions[J]. Transactions of the Association for Computational Linguistics, 2014, 2: 67-78.
- [30] QI D, SU L, SONG J, et al. ImageBERT: Cross-modal pre-training with large-scale weak-supervised image-text data[EB/OL]. (2020-01-23) [2024-03-25]. <https://doi.org/10.48550/arXiv.2001.07966>.
- [31] DOU Z Y, XU Y C, GAN Z, et al. An empirical study of training end-to-end vision-and-language Transformers [C//2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway: IEEE, 2022: 18166-18176.
- [32] YANG J Y, DUAN J L, TRAN S, et al. Vision-language pre-training with triple contrastive learning[C//2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway: IEEE, 2022: 15671-15680.
- [33] WANG W H, YANG Z, XU B, et al. ViLTA: Enhancing vision-language pre-training through textual augmentation [C//2023 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway: IEEE, 2023, 3158-3169.
- [34] KWON G, CAI Z W, RAVICHANDRAN A, et al. Masked vision and language modeling for multi-modal representation learning[EB/OL]. (2022-08-03) [2024-3-25]. <https://doi.org/10.48550/arXiv.2208.02131>.
- [35] BI J Y, CHENG D X, YAO P, et al. VL-match: Enhancing vision-language pretraining with token-level and instance-level matching[C//2023 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway: IEEE, 2023, 2584-2593.
- [36] LI J N, LI D X, XIONG C M, et al. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation[C//Proceedings of the 39th International Conference on Machine Learning. New York: PMLR, 2022, 162: 12888-12900.
- [37] AGRAWAL A, LU J S, ANTOL S, et al. VQA: Visual question answering[J]. International Journal of Computer Vision, 2017, 123(1): 4-31.
- [38] BYUN J, HWANG T, FU J L, et al. GRIT-VLP: Grouped mini-batch sampling for efficient vision and language pre-training[C//Computer Vision — ECCV 2022. Cham: Springer Nature Switzerland, 2022: 395-412.
- [39] ZHANG H T, ZHANG P C, HU X W, et al. GLIPv-2: Unifying localization and VL understanding[EB/OL]. (2022-06-12) [2024-3-25]. <https://doi.org/10.48550/arXiv.2206.05836>.
- [40] WANG Z R, YU J H, YU A W, et al. SimVLM: Simple visual language model pretraining with weak supervision [EB/OL]. (2021-08-24) [2024-3-25]. <https://doi.org/10.48550/arXiv.2108.10904>.
- [41] WANG P, YANG A, MEN R, et al. OFA: Unifying archi-

tures, tasks, and modalities through a simple sequence-to-sequence learning framework[C]//Proceedings of the 39th International Conference on Machine Learning. New York: PMLR, 2022, 162: 23318-23340.

- [42] SUHR A, ZHOU S, ZHANG A, et al. A corpus for reasoning about natural language grounded in photographs [C]//Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. Stroudsburg: Association for Computational Linguistics, 2019: 6418-6428.
- [43] XIE N, LAI F, DORAN D, et al. Visual entailment: A novel task for fine-grained image understanding[EB/OL]. (2019-01-20) [2024-03-25]. <https://doi.org/10.48550/arXiv.1901.06706>.
- [44] 刘天义, 吴祖焯, 陈静静, 等. 面向视觉语言理解与生成的多模态预训练方法[J]. 软件学报, 2023, 34(5): 2024-2034.
- LIU T Y, WU Z X, CHEN J J, et al. Multimodal pre-training method for vision-language understanding and generation[J]. Journal of Software, 2023, 34(5): 2024-2034. (in Chinese)



**冉越** 男, 2000年4月出生于河南省南阳市. 现为长春理工大学电子信息工程学院硕士研究生. 主要研究方向为计算机视觉和视觉语言多模态.

E-mail: ry13523068581@163.com



**孙俊喜** 男, 1971年6月出生于河北省唐山市. 现为东北师范大学信息科学与技术学院教授、博士生导师. 主要研究方向为AI视觉与智能感知技术.

E-mail: sunjx100@nenu.edu.cn

#### 作者简介



**才华** 男, 1977年2月出生于吉林省辉南县. 现为长春理工大学副教授、博士生导师. 获吉林省科技进步奖1项. 在国内外发表学术论文100余篇. 研究方向为计算机视觉与自然语言处理, 目前主持国家级省部级项目多项.

E-mail: caihua@cust.edu.cn



**易亚希** 男, 2001年2月出生于湖南省湘潭市. 现为长春理工大学电子信息工程学院硕士研究生. 主要研究方向为计算机视觉和视觉语言多模态.

E-mail: 2022100885@mails.cust.edu.cn



**付强** 男, 1984年8月出生于吉林省长春市. 现为长春理工大学空间光电技术研究所副所长. 主要研究方向为光学传输特性测试与多维度成像探测.

E-mail: fuqiang@cust.edu.cn