



融合PMA-Net与区域级CLIP提示的图像字幕生成模型设计

为了解决图像字幕生成中的视觉语义对齐难题，我们提出一种将原型记忆注意力网络（PMA-Net）与区域级CLIP提示机制融合的模型方案。该模型保留PMA-Net通过解码器自注意力层引入原型记忆的框架，同时结合FSGR方法的CLIP区域提示来增强图像区域的语义表示，使视觉特征与语言更加紧密对齐。下面将详细介绍模型结构、融合方法、创新点、应用场景、实现可行性，以及论文的大致结构。

模型整体结构组成

本方案的模型由以下几个主要模块构成：

- **图像编码器（Visual Encoder）**：采用预训练的CLIP视觉编码器（如ViT-L/14）来提取图像特征^①。整幅图像通过ViT获得一系列patch级别的视觉特征表示。这些特征为后续的区域语义提示提供基础。我们选择CLIP的ViT模型作为编码器，一方面因为其在跨模态表示上的高质量和强适应性，另一方面相较于传统目标检测特征，其提取的图像表示具有更低的计算开销^①。
- **区域级语义提示模块（Region-Level Prompt Module）**：这是融合FSGR机制的关键组件。该模块利用CLIP模型的文本编码能力，为图像中局部区域生成语义提示信息。具体而言，首先对图像编码器输出的patch特征进行局部对比学习优化，聚合同语义的图像patch，使其表示更加语义一致^②。接着，引入一个语义量化（Semantic Quantification）模块，通过预先构建的概念词池（例如大量类别标签组成的集合，采用模板如“a photo of {object}”经CLIP文本编码得到嵌入）将视觉特征映射到离散的语义空间^②。每个图像patch特征都与概念池中的文本嵌入计算相似度，选取最高相似的几个概念，使得语言桥接的视觉特征图得以生成^③。换言之，图像的区域特征被赋予了对应的文本语义标识。这些语义增强的区域特征将在后续作为视觉提示供解码器使用。
- **原型记忆模块（Prototype Memory Module）**：继承自PMA-Net，是模型的核心创新之一。不同于以往将记忆引入编码器的做法，我们将记忆库集成在解码器的每一层自注意力机制中^④。具体而言，在训练过程中，每个解码器自注意力层都会从当前样本的历史K/V激活中构建一对记忆银行（键记忆BK和值记忆BV），存储最近\$T\$个训练步骤中生成的所有键、值向量^⑤。然后通过聚类（如K-means）从记忆银行中提取原型向量集合\$(M_K, M_V)\$，作为该层自注意力的扩展键和值^{④ ⑥}。在实际的注意力计算中，解码器的查询\$Q\$不仅和当前输入序列的\$K,V\$交互，还会与这些记忆原型\$M_K, M_V\$进行注意力计算，从而让模型能够检索利用过去样本的激活信息^⑥。这一机制使模型在描述当前图像时，可以参考训练集中其他相似图像/句子的隐藏表示，仿佛有“经验”可循，从而丰富生成内容。
- **Transformer 解码器（Caption Decoder）**：采用基于Transformer的文本解码器，逐字生成图像描述句子。解码器在每一层包含自注意力、交叉注意力和前馈网络等子层。自注意力层通过上述记忆模块进行扩展，加入原型记忆的键值供注意力检索。交叉注意力层则将解码器当前隐藏状态作为查询，与图像编码器提供的区域语义特征图作为键和值进行交互，从视觉特征中提取与当前生成单词相关的信息。解码器的输出通过softmax形成描述文本。整个架构的数据流如下：
 1. 图像经编码器提取patch特征，经区域提示模块得到融合语义的视觉特征序列；
 2. 解码器逐步接收前文已生成的词作为输入，利用自注意力（含记忆）建模上下文并查询原型记忆库；
 3. 解码器利用交叉注意力从视觉特征序列中检索相关的区域语义信息，指导下一个词的生成；

4. 不断循环直到生成完整的描述句子。

上述模块通过精心的连接实现了视觉特征、语义提示和记忆信息的融合。例如，交叉注意力层的查询来自解码器上一层的输出隐状态，键/值由视觉编码器的区域特征和对应的语义提示embedding构成，实现视觉内容与文本生成的对齐；而在自注意力层，解码器查询同时关注自身序列历史以及记忆原型，从而结合**当前语境+过往经验**来形成新的隐藏表示。

融合CLIP区域提示与记忆模块的方法

融合FSGR的区域提示机制，需要解决提示信息在模型各部分的交互方式。我们的设计在**交叉注意力和记忆检索**两个方面，实现了提示信息与原型记忆的融合：

- **交叉注意力中的提示对齐**：在解码器的交叉注意力子层，我们让视觉特征不仅包含传统视觉编码器输出，还包括**文本提示特征**。具体做法是：将区域语义提示模块产生的**文本嵌入**（例如“horse”、“grass”等概念的CLIP文本向量）与对应的视觉patch特征结合，形成复合的键和值矩阵供交叉注意力使用。一种实现方式是对每个图像patch，用其最相关的概念文本嵌入替换或融合原始视觉特征（例如拼接后过线性变换融合）。这相当于为视觉区域显式地附加一个“语义标签”。因此，当解码器通过交叉注意力查询图像信息时，注意力机制可以直接检索到**携带语义提示的视觉特征**。比如，当生成关于“马”的描述时，解码器的查询可注意到带有“horse”语义的区域特征，从而更准确地获取与“马”相关的视觉信息。这种在交叉注意力层引入提示的方式，确保了图像区域和文本单词之间的细粒度对齐，有效缩小视觉特征与语言描述间的语义鸿沟²。
- **记忆模块中的提示关联**：原型记忆模块主要作用于解码器的自注意力，用于跨样本借鉴信息。虽然记忆库本身不直接存储外部的CLIP提示，但**隐式地**，区域提示机制对记忆检索是有帮助的。一方面，解码器在交叉注意力获取了语义增强的视觉信息后，其隐藏状态会更加语义明确，这些隐藏状态通过自注意力的查询参与记忆检索时，**会倾向于匹配到含有相似语义模式的记忆原型**。举例来说，当当前图像有关“马和草地”时，解码器隐藏向量中将编码“horse”、“grass”的语义，那么自注意力可能从记忆库中检索到过去描述过类似场景（有马和草地）的原型向量，从而提供更相关的历史信息。另一方面，我们可以在设计上进一步**引导记忆检索融合提示**：例如，对记忆键\$M_K\$增加语义过滤，仅检索与当前图像语义标签相符的一 subset（这可由训练时统计每个原型常见语义实现）；或者将当前图像提取的主要语义标签作为条件，影响记忆聚类的选样。总之，通过上述方式，记忆模块虽然不直接使用CLIP提示向量作为输入，但在**查询与匹配环节都受到区域语义的上下文影响**，做到提示信息与记忆检索的协同。这样，当模型生成描述时，既有来自当前图像的细粒度语义提示，又有来自相似语境下过往经验的支持，形成“双轮驱动”的信息源。

需要强调的是，我们保持PMA-Net的原有机制，即**仅利用解码器自身产生的激活构建记忆库，不从编码器特征中生成记忆**⁴。记忆库中的原型键值依然来源于解码器历史隐藏状态集合（经过聚类得到），这与区域提示模块提供的图像语义特征属于不同来源的信息。两者在解码器不同注意力层分别发挥作用：交叉注意力侧重**视觉-语言对齐**，自注意力（结合记忆）侧重**语言上下文扩展**。我们通过合理设计，使**区域提示信息和原型记忆在解码器中互补融合**：提示提供准确的视觉语义对齐，记忆提供丰富的语言先验支持。

技术创新点及视觉-语义对齐贡献

本方案在结构设计上融合了两种机制，具有如下技术创新与贡献：

- **原型记忆与细粒度语义提示的首次结合**：以往的图像描述模型很少同时利用“跨样本记忆检索”和“区域级语义提示”这两类技术。我们的设计将PMA-Net提供的**原型记忆网络**与FSGR提供的**区域语义对齐**手段结合，构建新的框架。在架构层面，这是一次有益的尝试：**记忆模块**擅长从历史数据中挖掘全局模式，而**CLIP提示**擅长对齐局部视觉语义，两者结合能互相弥补不足。例如，记忆模块可以缓解描述生成

中的模式偏差和遗忘问题（通过参考多样训练样本经验），而CLIP提示机制则解决了视觉特征语义抽象不足的问题，使模型更懂图像内容的细节含义。

- **增强的视觉-语言对齐：**该模型通过**显式的区域语义标注**，极大缩小了视觉特征和文本词汇之间的语义gap。传统注意力仅在视觉空间做软对齐，模型可能会将错误的标签赋予某些视觉区域。而引入CLIP提示后，每个区域特征都带有接近自然语言的含义，这使得解码器选择单词时有了明确的锚点，有效避免了语义偏差和内容幻觉。同时，**原型记忆**强化了视觉-语言对齐的一致性：当模型遇到稀有物体或复杂场景时，记忆库中类似场景的描述模式可供参考，帮助模型产生符合训练语义的描述，减少遗漏或误识别。例如，模型看到一个罕见乐器，通过CLIP提示获得其类别语义，又能从记忆中调取之前描述过该乐器的句子模式，从而准确地产生描述。这种多模态、多样本的对齐策略，提高了生成描述在细节和整体语义上的一致性。
- **结构设计上的创新：**我们在Transformer解码器中实现了**双流注意力融合**创新：一条流是**传统交叉注意力流**，提供图像视觉信息但经过语义富集；另一条流是**记忆自注意力流**，提供语言历史信息。相较原始PMA-Net，我们增加了**区域语义提示流**，但同时巧妙地保持了记忆流的位置和机制不变。这样的设计不仅保留了原模型验证有效的部分，还通过新增模块丰富了模型的表示能力。这种结构上的创新点在于：没有简单地把不同特征拼在一起，而是明确区分了语义提示与记忆的职责，让它们在Transformer不同子层各司其职又协同工作，体现了架构设计的合理性和新颖性。
- **提升跨模态任务性能：**通过上述创新，我们的模型在视觉-语义对齐问题上取得突破，预期将带来性能提升。基于报道，PMA-Net引入记忆可令COCO指标提升约3.7 CIDEr⁷；FSGR机制则实现了MSCOCO数据集新的SOTA表现，并在NoCaps开放域测试中表现出色⁸。融合两者后，模型有望在描述生成的准确性、丰富性上均超越现有方法。这证明了我们的创新设计在有效融合多来源信息、加强跨模态对齐方面的价值。

适用的应用场景与任务

该融合模型具备较强的通用性和开放域能力，适用于多种视觉语言任务：

- **通用图像字幕生成：**在传统的图片描述任务（如COCO数据集）上，本模型能够生成更加细致且准确的描述文本。区域级提示确保模型识别并写出图像中的关键对象和属性（例如物体种类、颜色、动作等），而记忆模块则帮助模型借鉴训练集中相似图像的表达，使描述更贴切人类语言风格。对于要求高描述质量的应用（如照片描述、新闻图片说明），该模型都能胜任。
- **开放域视觉描述：**由于引入了CLIP预训练知识，本模型能够处理开放域中的新奇概念和细粒度类别描述。这对于**开放域图像字幕**（Open-vocabulary Captioning）非常重要。在测试图像包含训练集中未见过的新对象时，CLIP提供的语义提示使模型仍能识别并命名该对象（因为CLIP的视觉语义空间覆盖广泛），避免出现“未知物体”的情况。例如，一张包含珍稀动物的照片，模型可直接利用CLIP提示获取该动物的物种名称并体现在描述中。原型记忆在这种情况下也会辅助模型调用相似语境的句型，确保语言连贯。
- **细粒度描述和密集字幕：**模型能够在需要细节的场景下工作，比如对复杂场景生成细粒度描述、或者对图像不同区域生成密集字幕（Dense Captioning）。区域提示机制本身强调了区域级别的表示，如果辅以适当的区域提取方法（如区域提案网络或分割），模型可扩展用于产生**区域性的描述**。每个候选区域都可通过CLIP提示获得语义标签，再由解码器配合记忆生成该区域的描述文本。这在应用上可以用于**图像内容解析、辅助盲人场景解说**等需要丰富细节的场景。
- **跨领域的视觉语言任务：**由于记忆模块存储了多样的数据模式，CLIP提供跨领域的知识，本模型对于**跨领域**（例如医疗、遥感等非自然图像领域）的描述任务也有潜力。只要有对应领域的数据用于训练或提供概念提示，模型能够利用CLIP的知识对专业领域物体命名，并利用记忆模块保持描述风格的一致性。

这意味着从日常影像到专业影像的广泛场景下，本模型都能提供合理的文字描述，具有较好的泛化能力⁸。

总之，本设计既适用于常规的图像字幕任务（提升评价指标和描述质量），又具备开放域和细粒度扩展能力，能够满足新兴的多样化视觉-语言应用需求。

实现可行性分析

在实现方面，我们需要考量训练开销、对原始模型的改动难度以及复现可能性：

- **训练资源需求：**模型引入了CLIP大型视觉编码器和记忆模块，训练时计算和存储开销会有所增加。CLIP ViT-L/14本身较大，不过我们可以冻结CLIP编码器参数，通过在其中注入少量可训练参数来降低显存占用⁹。例如，采用视觉提示向量和适配器（adapter）微调CLIP，而不微调其主干权重⁹。FSGR论文已证明，只训练少量附加参数即可有效提取本地语义且保持CLIP知识，这使得训练一个融合模型成为可能。原型记忆模块方面，主要开销在于存储最近T个batch的键值以及定期执行K-Means聚类⁵。作者实现中利用了高效的Faiss库GPU版KNN和KMeans，加上存储长度T可控（如T=1500）¹。因此在现代GPU集群上（例如8卡以上），训练这样一个模型是可行的，但需要合理规划显存（尤其是patch特征、文本提示和记忆库的存储）以及训练时间（聚类操作在每个epoch会进行有限次数）。总的来说，虽然资源需求较标准Caption模型高一些，但在研究环境下是可承受的。
- **对原始PMA-Net的改动难度：**我们的设计在PMA-Net基础上增加了“区域提示”相关模块，但PMA-Net的核心流水并未被破坏。具体改动包括：集成CLIP模型作为图像编码器（PMA-Net原本也使用CLIP特征，因此变化不大，只是现在需要在线提取patch而非使用预提特征）¹；在交叉注意力前增加语义提示融合操作，这可以通过增加一个前馈网络或Attention层实现，将视觉patch和文本提示组合成新的键值矩阵。解码器结构基本保持，唯一期望调整的是交叉注意力的键值维度（由于融合了文本提示，可能需要在通道维上扩充或采用多头分别Attention视觉/文本）。这些改动在Transformer框架中是**局部的、可插拔的**，开发者可在PMA-Net开源代码基础上增添相应模块。原型记忆的代码逻辑几乎不变，只需确保在带提示的条件下仍能正常检索——这主要靠模型学习，无需修改算法。综上，改动难度中等，可基于现有代码扩展，无需推倒重来。
- **复现和实现风险：**PMA-Net和相关CLIP提取、MaskCLIP等代码都是开源可得的⁸。这为我们的复现提供了良好基础。我们可能需要复用FSGR的一些预处理步骤，例如获取概念词典/语义标签（FSGR利用MaskCLIP等方法自动生成patch语义标签¹⁰）。这些也都有现成工具，可以减轻开发负担。在实现过程中，需要注意的一个风险是：**多模块联合训练的稳定性**。记忆模块引入非参数化的动态检索过程，CLIP提示需要跨模型梯度传播，二者的交互可能收敛较慢。我们可以采取分阶段训练策略：先训练基础caption模型+CLIP提示模块（不启用记忆），待其收敛后再加载记忆模块继续训练；或者联合训练时降低记忆检索频率、平滑更新原型等技巧，以确保训练稳定。只要按照论文提供的细节调整，这些都是可以克服的。由于是对成熟方法的融合，复现实验应当较为顺利，可以预期得到与论文报告相近的性能提升。

总而言之，本方案在实现上具备**较高可行性**：所需组件大多已有成熟实现，我们的创新部分在于整合和改进，开发难度可控；训练资源需求增加但在可接受范围；有清晰的复现路径和公开数据支持。这些都为论文方案的落地实现提供了保障。

论文结构大纲（建议）

最后，我们提供一份合理的论文结构大纲，以清晰组织上述内容，突出本研究的贡献：

1. **引言（Introduction）**：介绍图像字幕生成任务背景和挑战，概述当前视觉-语言对齐的难点；引出原型记忆和CLIP提示两种思路，并指出各自局限；提出本文的融合模型方案，强调其创新之处和潜在贡献；最后给出主要贡献点陈述。
2. **相关工作（Related Work）**：梳理与本研究相关的文献。包括：(a) 图像字幕生成方面：自注意力模型、记忆增强模型（如PMA-Net⁶、AMA等）、开放词汇字幕；(b) 视觉-语义对齐方面：CLIP在图像字幕/生成中的应用、细粒度语义对齐方法（如RegionCLIP、FSGR²等）；(c) 讨论我们方法与现有工作的区别，例如我们同时引入记忆和区域提示，实现不同模块优势互补，这是此前方法未曾探索的组合。
3. **方法（Methodology）**：详细介绍所提出模型的结构与方法。
 4. **3.1 总体架构概览**：给出模型框架图，描述模型包含的编码器、区域提示模块、记忆模块、解码器组件，以及它们的连接关系（对应我们的模型整体结构部分）。
 5. **3.2 区域语义提示模块**：阐述该模块的设计和原理。包括图像patch提取、局部对比学习细节（如何使相似patch聚类提高一致性）、语义量化过程（概念池的构建，匹配策略），以及与交叉注意力的融合方式。可附公式描述交叉注意力键值的组成： $K=[K_{\text{text}\{\text{vision}\}}; K_{\text{text}\{\text{sem}\}}]$, $V=[V_{\text{text}\{\text{vision}\}}; V_{\text{text}\{\text{sem}\}}]$$ 等等。
 6. **3.3 原型记忆自注意力机制**：介绍PMA-Net的记忆机制如何嵌入解码器层⁴。给出记忆键值生成公式、原型形成算法（K-means聚类获取\$M_K, M_V\$）⁵。重点突出我们在解码器中而非编码器中引入记忆的决策，以及记忆如何在推理时工作（例如使用训练得到的固定原型集合）。讨论记忆模块与区域提示的协同作用机制。
 7. **3.4 损失函数与训练策略**：列出模型训练用到的目标函数。如交叉熵损失用于字幕词生成，可能还包括对比损失（用于局部对比学习阶段，拉近正确patch-语义对，推远错误匹配）²。如果有的话，提及采用强化学习优化CIDEr的策略（沿用PMA-Net后期Self-critical训练）。还需说明训练分阶段或联合训练的方案、超参数设置等。
8. **实验（Experiments）**：
 9. **4.1 实验设置**：说明实验所用数据集（如MS COCO、NoCaps等）、评估指标（BLEU, CIDEr, SPICE等）、实现细节（模型超参数，同步Batch大小，记忆库T值，使用的CLIP模型版本等）。
 10. **4.2 与现有方法的比较**：给出主要结果表格，将本模型性能与主流模型（如Transformer基线、M2Transformer、M^2+CLIP方法、记忆模型PMA-Net、FSGR单独模型等）进行比较。突出本方法在CIDEr等核心指标上的提升，并通过**加粗**标注新的SOTA成绩。
 11. **4.3 消融实验（Ablation Study）**：为了验证各模块贡献，设计消融实验。例如：(a) 去掉记忆模块，仅保留CLIP提示，性能如何变化；(b) 去掉CLIP区域提示，仅用PMA-Net记忆，性能变化；(c) 记忆库大小、原型数对结果的影响；(d) 是否使用对比学习目标对区域提示效果的影响。通过定量分析证明我们的设计各部分都是有益的。
 12. **4.4 质性分析（Qualitative Analysis）**：展示一些例子，直观说明模型效果。例如图像对应的生成描述与其他模型对比；可视化交叉注意力的热力图，看到模型确实关注到了对应语义的图像区域；展示记忆模块检索到的训练样本原型的实例，以说明记忆如何帮助生成。特别地，可以举例说明当没有区域提示时模型出了某种错误描述，而加入提示后纠正了，对应attention图也更合理，证明语义对齐的改进。
 13. **讨论（Discussion）**：根据实验结果，对模型的适用范围、局限性进行讨论。例如，本模型在开放域上表现好，说明了CLIP提示的价值；但对于非常细粒度的领域专业名词，如果概念池覆盖不到，可能仍有

困难——这提示我们概念池需要扩充或结合知识库。记忆模块可能在训练早期引入噪声，需要平衡，这也是未来改进方向之一。讨论也可涉及模型推理效率（记忆检索会增加推理时间，如何权衡等）。

14. 结论 (Conclusion) : 总结全文，重申本文提出了融合原型记忆和区域CLIP提示的图像字幕模型，在语义对齐上取得了技术创新和实绩提升。指出本研究对未来多模态学习的启示，例如结合大模型知识与显式记忆是提升模型智能的有效途径。最后给出未来工作展望，比如扩展到视频字幕、对话式描述等。

通过上述结构，大纲覆盖了从理论动机、方法设计、实验验证到讨论总结的完整内容，清晰呈现了本研究的脉络与贡献。希望这份结构化报告有助于读者理解我们设计的新颖之处，以及其在视觉-语言对齐任务中的潜在价值。

6 2

1 4 5 With a Little Help from Your Own Past: Prototypical Memory Networks for Image Captioning

<https://openaccess.thecvf.com/content/ICCV2023/papers/>

Barraco_With_a_Little_Help_from_Your_Own_Past_Prootypical_Memory_ICCV_2023_paper.pdf

2 3 8 9 Lei Zhang's research works | Hefei University of Technology and other places

<https://www.researchgate.net/scientific-contributions/Lei-Zhang-2225267979>

6 7 Paper page - With a Little Help from your own Past: Prototypical Memory Networks for Image Captioning

<https://huggingface.co/papers/2308.12383>

10 GitHub - gjc0824/FSGR

<https://github.com/gjc0824/FSGR>