

# 深度学习作业——基于 CNN-Transformer 的图像描述

Clark

## 1 实验目的

本实验旨在基于 PyTorch 深度学习框架，构建并训练一个结合卷积神经网络（CNN）与 Transformer 的图像描述（Image Captioning）模型。通过本实验，主要达到以下目的：

1. 深入理解多模态学习任务，特别是图像到文本生成的编码器-解码器（Encoder-Decoder）架构；
2. 掌握卷积神经网络（如 ResNet）在图像特征提取中的应用原理；
3. 掌握 Transformer 模型的核心机制，包括自注意力（Self-Attention）、多头注意力（Multi-Head Attention）及位置编码（Positional Encoding）；
4. 学习图像描述任务的常用评估指标，如 BLEU、ROUGE-L、CIDEr 等；
5. 实践深度学习工程化流程，包括 Flickr8k 数据集的预处理、模型构建、训练调试及结果可视化。

## 2 概述

图像描述（Image Captioning）是计算机视觉（Computer Vision, CV）与自然语言处理（Natural Language Processing, NLP）的交叉领域，其目标是为给定的图像生成一句符合语法且描述准确的自然语言句子。本实验采用经典的“CNN+Transformer”架构：使用预训练的 ResNet101 作为编码器提取图像的高层语义特征，使用 Transformer 作为解码器将图像特征转化为文本序列。实验在 Flickr8k 数据集上进行，涵盖了数据下载、词表构建、模型搭建、训练循环及基于 Beam Search 的生成与评估。

## 3 实验要求

1. 基于 Python 语言和 PyTorch 框架，完成 Flickr8k 数据集的下载、解压与读取；
2. 构建 CNN-Transformer 模型，其中 CNN 编码器提取图像特征，Transformer 解码器生成描述；
3. 实现完整的训练流程，包括损失计算、反向传播及优化器更新；
4. 使用 BLEU-1, BLEU-4, ROUGE-L, CIDEr 等指标对模型进行定量评估；
5. 绘制训练损失及各项评估指标随 Epoch 变化的曲线，并展示模型生成的实际效果图。

## 4 实验原理

### 4.1 卷积神经网络（CNN）与特征提取

卷积神经网络（Convolutional Neural Networks, CNN）是计算机视觉领域的基石，其设计灵感来源于生物视觉皮层的感受野机制。CNN 通过局部连接、权值共享和下采样等特性，能够有效地从高维图像数据中提取具有平移、缩放和旋转不变性的层级特征。

#### 4.1.1 基本原理与数学表述

CNN 的核心组件包括卷积层（Convolutional Layer）、激活层（Activation Layer）和池化层（Pooling Layer）。

卷积层通过可学习的滤波器（Filter/Kernel）在输入特征图上滑动，执行互相关运算以提取局部特征。对于输入特征图  $\mathbf{X} \in \mathbb{R}^{C_{in} \times H \times W}$  和第  $k$  个卷积核  $\mathbf{W}_k \in \mathbb{R}^{C_{in} \times K_h \times K_w}$ ，输出特征图  $\mathbf{Y}_k$  的第  $(i, j)$  个元素可表示为：

$$\mathbf{Y}_{k,i,j} = \sigma \left( \sum_{c=1}^{C_{in}} \sum_{m=0}^{K_h-1} \sum_{n=0}^{K_w-1} \mathbf{W}_{k,c,m,n} \cdot \mathbf{X}_{c,i \cdot s + m, j \cdot s + n} + b_k \right), \quad (1)$$

其中， $s$  为步长（Stride）， $b_k$  为偏置， $\sigma(\cdot)$  为非线性激活函数（如 ReLU）。这一过程使得网络能够学习到边缘、纹理等底层特征，并随着层数加深逐渐组合成形状、物体部件等高层语义特征。

池化层用于降低特征图的空间分辨率，减少参数量并提高计算效率，同时增强特征的鲁棒性。最常用的最大池化（Max Pooling）定义为：

$$\mathbf{Y}_{i,j} = \max_{(m,n) \in \mathcal{R}_{i,j}} \mathbf{X}_{m,n}, \quad (2)$$

其中  $\mathcal{R}_{i,j}$  是池化窗口覆盖的区域。

#### 4.1.2 残差网络（ResNet）与深层特征提取

随着网络深度的增加，传统的 CNN 容易出现梯度消失或爆炸，导致网络退化（Degradation），即深层网络的训练误差反而高于浅层网络。本实验采用的 ResNet101 通过引入残差学习（Residual Learning）框架解决了这一问题。

ResNet 的基本单元是残差块（Residual Block），其核心思想是学习残差映射  $\mathcal{F}(\mathbf{x})$  而非直接学习潜在映射  $\mathcal{H}(\mathbf{x})$ 。假设  $\mathcal{H}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x}$ ，则残差块的输出定义为：

$$\mathbf{y} = \sigma(\mathcal{F}(\mathbf{x}, \{W_i\}) + \mathbf{x}), \quad (3)$$

其中  $\mathbf{x}$  是恒等映射（Identity Mapping），允许梯度在反向传播时通过“高速公路”直接流向浅层，从而极大地改善了深层网络的训练难度。

在本实验中，我们移除了 ResNet101 最后的全局平均池化层和全连接层，保留了最后一个卷积块（Layer4）的输出。对于输入尺寸为  $224 \times 224$  的图像，ResNet101 输出的特征图维度为  $2048 \times 7 \times 7$ 。我们将这  $7 \times 7 = 49$  个空间位置的特征向量视为图像的“视觉单词”，并通过一个  $1 \times 1$  卷积层将其维度从 2048 映射到 Transformer 的嵌入维度  $d_{\text{model}}$ ，作为解码器的上下文输入。

### 4.1.3 应用现状

尽管近年来 Vision Transformer (ViT) 在多个视觉任务上取得了突破, 但 CNN 凭借其归纳偏置 (Inductive Bias) 和高效的推理速度, 依然是工业界的主流选择。ResNet 及其变体 (如 ResNeXt, ConvNeXt) 广泛应用于目标检测 (如 Faster R-CNN, YOLO)、语义分割 (如 DeepLab, U-Net) 以及医学影像分析等领域, 证明了其作为视觉骨干网络的强大泛化能力。

## 4.2 Transformer 模型

Transformer 模型由 Vaswani 等人于 2017 年提出, 它彻底摒弃了传统的循环神经网络 (RNN) 和卷积神经网络 (CNN) 结构, 完全基于注意力机制 (Attention Mechanism) 来处理序列数据。这种架构不仅解决了 RNN 无法并行计算的问题, 还能够捕捉序列中任意距离的全局依赖关系。

### 4.2.1 自注意力机制 (Self-Attention)

自注意力机制是 Transformer 的核心, 其灵感来源于信息检索系统。对于输入序列中的每个元素, 模型将其映射为三个向量: 查询向量 (Query,  $\mathbf{Q}$ )、键向量 (Key,  $\mathbf{K}$ ) 和值向量 (Value,  $\mathbf{V}$ )。

缩放点积注意力 (Scaled Dot-Product Attention) 的计算公式为:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}, \quad (4)$$

其中,  $d_k$  是键向量的维度。除以  $\sqrt{d_k}$  是为了防止点积结果过大导致 softmax 函数进入梯度极小的饱和区。直观上,  $\mathbf{Q}\mathbf{K}^T$  计算了查询与所有键的相似度 (即注意力分数), softmax 将其归一化为权重, 最后对  $\mathbf{V}$  进行加权求和, 从而聚合了上下文信息。

### 4.2.2 多头注意力 (Multi-Head Attention)

为了让模型能够从不同的表示子空间 (Representation Subspaces) 共同关注信息, Transformer 引入了多头注意力机制。它将  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  分别通过  $h$  组不同的线性投影矩阵映射到低维空间, 并行计算注意力, 最后将结果拼接并通过线性层输出:

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)\mathbf{W}^O, \quad (5)$$

其中:  $\text{head}_i = \text{Attention}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$ .

这种机制类似于 CNN 中的多通道 (Filters), 允许模型同时关注 “谁” (Who)、“在哪” (Where) 和 “做什么” (What) 等不同层面的语义信息。

### 4.2.3 位置前馈网络与残差连接

除了注意力层, Transformer 的每个编码器/解码器层还包含一个全连接的前馈网络 (Position-wise Feed-Forward Networks, FFN), 它对序列中的每个位置独立且相同地进行变换:

$$\text{FFN}(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2. \quad (6)$$

此外, 每个子层 (Attention 和 FFN) 周围都采用了残差连接 (Residual Connection) 和层归一化 (Layer Normalization), 即  $\text{LayerNorm}(\mathbf{x} + \text{Sublayer}(\mathbf{x}))$ , 这有助于梯度的稳定传播和深层网络的训练。

#### 4.2.4 位置编码 (Positional Encoding)

由于自注意力机制具有排列不变性 (Permutation Invariant)，无法感知序列的顺序，因此必须显式注入位置信息。Transformer 采用正弦和余弦函数生成位置编码，并将其与输入 Embedding 相加：

$$\begin{aligned} PE(\text{pos}, 2i) &= \sin(\text{pos} / 10000^{2i/d_{\text{model}}}), \\ PE(\text{pos}, 2i + 1) &= \cos(\text{pos} / 10000^{2i/d_{\text{model}}}). \end{aligned} \quad (7)$$

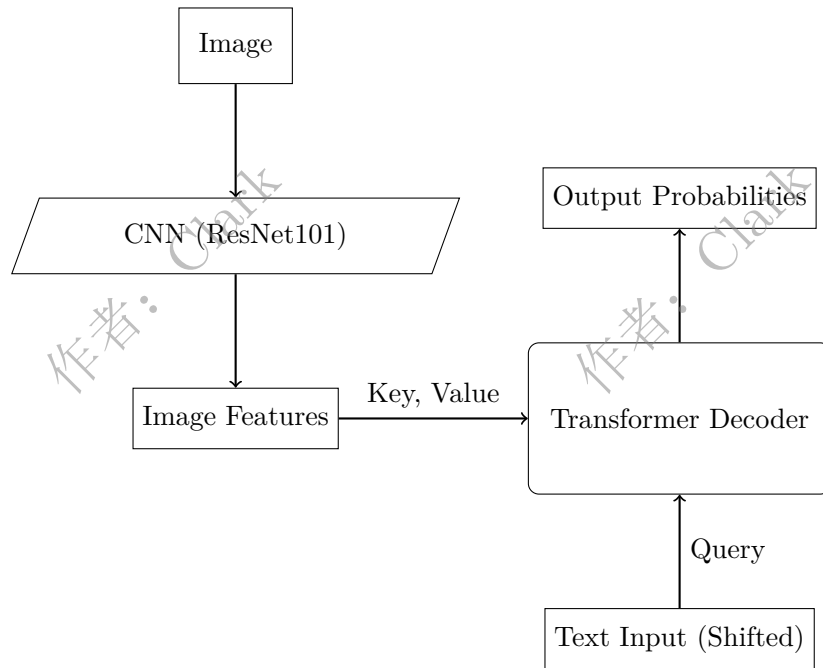
这种编码方式允许模型通过线性变换学习相对位置关系，且能够处理比训练序列更长的推理序列。

#### 4.2.5 应用现状

Transformer 架构的提出标志着 NLP 领域的“ImageNet 时刻”。基于 Transformer 的预训练模型如 BERT（双向编码器表示）和 GPT（生成式预训练 Transformer）在几乎所有 NLP 任务上刷新了 SOTA。近年来，Transformer 也成功跨界至计算机视觉（如 ViT, Swin Transformer）和多模态领域（如 CLIP, Stable Diffusion），成为通用的基础模型架构（Foundation Models）。

### 4.3 CNN-Transformer 联合架构

本实验构建的模型属于典型的编码器-解码器 (Encoder-Decoder) 架构。其中，CNN 负责将像素级的图像信息转换为高维语义特征空间，而 Transformer 解码器则负责将这些特征翻译为自然语言序列。两者通过交叉注意力机制 (Cross-Attention) 实现多模态信息的深度融合。具体的 CNN-Transformer 图像描述模型架构如下图所示。



#### 4.3.1 视觉特征的序列化与映射

在代码实现（“CNNEncoder”类）中，ResNet101 输出的特征图  $\mathbf{F} \in \mathbb{R}^{B \times 2048 \times 7 \times 7}$  包含了丰富的空间信息。为了适配 Transformer 的输入要求，我们需要将其转换为序列形式：

1. 降维投影：首先使用一个  $1 \times 1$  的卷积层将通道数从 2048 压缩至嵌入维度  $d_{\text{model}} = 512$ 。这一步不仅减少了计算量，还引入了非线性变换。
2. 展平与重排：将特征图的空间维度  $7 \times 7$  展平为长度为 49 的序列。此时，图像被表示为 49 个特征向量的集合  $\mathbf{M} = \{\mathbf{m}_1, \dots, \mathbf{m}_{49}\}$ ，其中  $\mathbf{m}_i \in \mathbb{R}^{512}$  对应原图中某个  $32 \times 32$  像素区域（感受野）的特征。
3. 可学习位置编码：为了让 Transformer 感知图像的空间结构（如“左上”、“右下”），代码中引入了一个可学习的参数“self.pos\_embed”  $\in \mathbb{R}^{49 \times 1 \times 512}$ ，将其直接加到特征序列  $\mathbf{M}$  上。

#### 4.3.2 交叉注意力机制（Cross-Attention）

这是连接视觉模态和语言模态的桥梁。在解码器的每一层中，交叉注意力模块接收两类输入：

- 查询向量（Query ( $\mathbf{Q}$ )）：来自解码器上一层（Masked Self-Attention）的输出，代表当前已生成的文本上下文信息。
- 键向量（Key ( $\mathbf{K}$ ））和值向量（Value ( $\mathbf{V}$ )）：来自编码器输出的图像特征序列  $\mathbf{M}$ 。

数学上，对于第  $t$  步生成的文本隐向量  $\mathbf{h}_t$ ，模型计算其与图像中每个区域  $\mathbf{m}_j$  的相关性权重  $\alpha_{t,j}$ ：

$$\alpha_{t,j} = \text{softmax} \left( \frac{(\mathbf{h}_t \mathbf{W}_Q)(\mathbf{m}_j \mathbf{W}_K)^T}{\sqrt{d_k}} \right), \quad (8)$$

然后对图像特征进行加权聚合得到上下文向量  $\mathbf{c}_t = \sum_{j=1}^{49} \alpha_{t,j} (\mathbf{m}_j \mathbf{W}_V)$ 。这一机制使得模型在生成“鸟”这个词时，能够自动给予图像中鸟所在区域更高的关注度（Attention Weight），实现了“看图说话”。

#### 4.3.3 自回归生成与因果掩码

解码过程是自回归（Autoregressive）的，即模型根据图像特征  $\mathbf{I}$  和历史词序列  $y_{1:t-1}$  预测下一个词  $y_t$ ：

$$P(y_t | y_{1:t-1}, \mathbf{I}) = \text{softmax}(\text{Decoder}(y_{1:t-1}, \text{Encoder}(\mathbf{I}))). \quad (9)$$

在训练阶段，为了并行计算，我们使用因果掩码（Causal Mask/Square Subsequent Mask）。在代码“TransformerDecoderModel”中，“generate\_square\_subsequent\_mask”函数生成一个上三角矩阵（值为  $-\infty$ ），作用于 Self-Attention 的 Logits 上。这确保了在预测位置  $t$  时，模型只能“看到”位置  $t$  及之前的词，防止了未来信息的泄露。

#### 4.3.4 前沿应用

这种“视觉编码器 + 语言解码器”的范式是现代多模态大模型（Multimodal LLMs）的基础。例如，OpenAI 的 GPT-4V 和开源社区的 LLaVA（Large Language-and-Vision Assistant）均沿用了这一思路：使用强大的视觉骨干（如 ViT-L/14）提取特征，通过线性层或 Q-Former 投影到语言模型的 Embedding 空间，从而赋予 LLM 理解图像的能力，实现视觉问答（VQA）、图像推理等复杂任务。

## 5 实验数据集及工具

### 5.1 实验配置

本实验基于 PyTorch 框架实现，使用 Flickr8k 数据集进行训练和评估。详细的实验配置下表所示。

配置项	参数/说明
深度学习框架	PyTorch
数据集	Flickr8k (8091 张图像)
数据集划分	Karpathy Split (Train: 6000, Val: 1000, Test: ~ 1091)
图像预处理 (训练)	Resize (256), RandomCrop (224), RandomHorizontalFlip
图像预处理 (测试)	Resize (224), CenterCrop (224)
编码器 (Encoder)	ResNet101 (Pretrained on ImageNet)
解码器 (Decoder)	Transformer Decoder (4 Layers, 8 Heads)
嵌入维度 ( $d_{\text{model}}$ )	512
前馈网络维度 ( $d_{\text{ff}}$ )	512
Dropout	0.3
优化器	Adam ( $\text{lr}=1 \times 10^{-4}$ )
学习率调度	ReduceLROnPlateau (factor=0.5, patience=2)
损失函数	CrossEntropyLoss (Label Smoothing=0.1)
Batch Size	32
Epochs	20

### 5.2 代码实现详解

本实验的代码实现主要包含数据准备、词表构建、模型定义、训练循环和评估指标计算五个核心模块，以及主函数。以下将分别对各模块进行详细说明。

#### 5.2.1 数据准备与预处理 (Data Preparation)

“prepare\_dataset”函数负责自动下载并解压 Flickr8k 数据集。代码首先检查本地是否存在数据集，若不存在则从配置的 URL 下载。为了应对网络不稳定的情况，代码实现了多源下载策略（支持镜像源）和多种下载工具的自动切换（优先使用“aria2c”，其次“wget”，最后“urllib”）。下载完成后，自动解压并整理文件结构，确保“images”和“annotations”目录就绪。

```

1 def download_file(url, dest_path):
2     if os.path.exists(dest_path):
3         # Check for corrupted zip files
4         if dest_path.endswith('.zip'):
5             try:
6                 with zipfile.ZipFile(dest_path, 'r') as z:
7                     pass # Just check if we can open it
8                 print(f"File already exists: {dest_path}")

```

```

9         return
10 except zipfile.BadZipFile:
11     print(f"Found corrupted zip file: {dest_path}.
12     Deleting and re-downloading...")
13     os.remove(dest_path)
14 else:
15     print(f"File already exists: {dest_path}")
16     return
17
18 # Build candidate URLs: primary first, then mirrors (if configured)
19 candidate_urls = [url]
20 try:
21     parsed = urllib.parse.urlparse(url)
22     host = parsed.netloc
23     if config.USE_MIRROR and isinstance(config.MIRROR_MAP, dict):
24         mirror = config.MIRROR_MAP.get(host)
25         if mirror:
26             # Mirror may be a full URL (with scheme) or just a host
27             if mirror.startswith('http://') or mirror.startswith('https://'):
28                 mirror_url = mirror.rstrip('/') + parsed.path
29             else:
30                 mirror_url = parsed.scheme + '://' + mirror.rstrip('/') + parsed.path
31             # Put mirror first to try it before original
32             candidate_urls.insert(0, mirror_url)
33 except Exception:
34     # If parsing fails, proceed with original URL
35     pass
36
37 last_exc = None
38 for candidate in candidate_urls:
39     print(f"Attempting download: {candidate} -> {dest_path}")
40
41     # 1. Try aria2c (fastest, multi-connection)
42     if shutil.which('aria2c'):
43         print("Using aria2c...")
44         ret = os.system(f"aria2c -x4 -s4 -c -o '{os.path.basename(dest_path)}' '{
45             candidate}' --dir='{os.path.dirname(dest_path)}' or '.'")
46         if ret == 0:
47             print("Download complete (aria2c).")
48             return
49         else:
50             print(f"aria2c failed with code {ret}, trying next tool...")
51
52     # 2. Try wget (reliable, resume support)
53     if shutil.which('wget'):
54         print("Using wget...")

```

```

52         ret = os.system(f"wget -c '{candidate}' -O '{dest_path}'")
53     if ret == 0:
54         print("Download complete (wget).")
55         return
56     else:
57         print(f"wget failed with code {ret}, trying next tool...")
58
59     # 3. Fallback to urllib (basic)
60     try:
61         print("Using urllib...")
62         urllib.request.urlretrieve(candidate, dest_path)
63         print("Download complete (urllib).")
64         return
65     except Exception as e:
66         print(f"Download failed for {candidate}: {e}")
67         last_exc = e
68
69     # If we reach here, all candidates failed
70     raise RuntimeError(f"Failed to download {url} -> {dest_path}") from last_exc
71
72 def unzip_file(zip_path, extract_to):
73     print(f"Unzipping {zip_path}...")
74     with zipfile.ZipFile(zip_path, 'r') as zip_ref:
75         zip_ref.extractall(extract_to)
76     print("Unzip complete.")
77
78 def prepare_dataset():
79     if config.USE_DUMMY_DATA:
80         print("Using Dummy Data. Skipping download.")
81         return
82
83     if not os.path.exists(config.DATA_DIR):
84         os.makedirs(config.DATA_DIR)
85     if not os.path.exists(config.IMG_DIR):
86         os.makedirs(config.IMG_DIR)
87     if not os.path.exists(config.ANN_DIR):
88         os.makedirs(config.ANN_DIR)
89     if not os.path.exists(config.CHECKPOINT_DIR):
90         os.makedirs(config.CHECKPOINT_DIR)
91
92     # 1. Annotations (Text)
93     if not os.path.exists(config.CAPTION_FILE):
94         text_zip = os.path.join(config.DATA_DIR, "Flickr8k_text.zip")
95         download_file(config.URL_TEXT, text_zip)
96         unzip_file(text_zip, config.ANN_DIR)

```



```

97     # 2. Images
98     # Check if images folder is populated (Flickr8k has 8091 images)    if len(os.listdir(
        config.IMG_DIR)) < 8000:
99         img_zip = os.path.join(config.DATA_DIR, "Flickr8k_Dataset.zip")
100         download_file(config.URL_IMAGES, img_zip)
101         # Extract to DATA_DIR first
102         unzip_file(img_zip, config.DATA_DIR)
103
104         # The zip usually contains a folder named 'Flickr8k_Dataset' (note the 'er')
105         extracted_folder = os.path.join(config.DATA_DIR, "Flickr8k_Dataset")
106         if os.path.exists(extracted_folder):
107             print(f"Moving images from {extracted_folder} to {config.IMG_DIR}...")
108             for f in os.listdir(extracted_folder):
109                 shutil.move(os.path.join(extracted_folder, f), config.IMG_DIR)
110             os.rmdir(extracted_folder)
111
112     print("Dataset preparation check passed.")

```

### 5.2.2 词表与数据集构建 (Vocabulary & Dataset)

“Vocabulary”类负责构建文本到索引的映射。使用 NLTK 库进行分词，并统计词频。为了减小词表规模并过滤噪声，设置了频率阈值“freq\_threshold=5”，仅保留出现次数大于等于 5 次的单词。特殊 Token 包括：“<PAD>”（填充）、“<SOS>”（句首）、“<EOS>”（句尾）和“<UNK>”（未知词）。“Flickr8kDataset”类继承自“torch.utils.data.Dataset”，负责加载图像和对应的描述。

- 图像加载：使用 PIL 库读取图像，并转换为 RGB 模式；
- 文本处理：将文本描述转换为数值索引序列，并在首尾分别添加“<SOS>”和“<EOS>”；
- 数据增强：在训练阶段，通过“transforms.RandomCrop”和“transforms.RandomHorizontalFlip”增强数据的多样性，防止过拟合。

```

1  class Vocabulary:
2      def __init__(self, freq_threshold):
3          self.itos = {0: "<PAD>", 1: "<SOS>", 2: "<EOS>", 3: "<UNK>"}
4          self.stoi = {"<PAD>": 0, "<SOS>": 1, "<EOS>": 2, "<UNK>": 3}
5          self.freq_threshold = freq_threshold
6
7      def __len__(self):
8          return len(self.itos)
9
10     @staticmethod
11     def tokenizer_eng(text):
12         return nltk.tokenize.word_tokenize(text.lower())
13

```

```

14     def build_vocabulary(self, sentence_list):
15         frequencies = Counter()
16     idx = 4         for sentence in sentence_list:
17                     for word in self.tokenizer_eng(sentence):
18                         frequencies[word] += 1
19                     if frequencies[word] == self.freq_threshold:
20                         self.stoi[word] = idx
21                         self.itos[idx] = word
22                         idx += 1
23
24     def numericalize(self, text):
25         tokenized_text = self.tokenizer_eng(text)
26         return [
27             self.stoi[token] if token in self.stoi else self.stoi["<UNK>"]
28             for token in tokenized_text
29         ]
30
31 class Flickr8kDataset(Dataset):
32     def __init__(self, root_dir, captions_file, transform=None, vocab=None,
33                 freq_threshold=5):
34         self.root_dir = root_dir
35         self.transform = transform
36         self.vocab = vocab
37         self.captions_file = captions_file
38
39         # Load captions
40         self.captions = [] # List of (image_name, caption_text)
41         print(f>Loading annotations from {captions_file}...")
42         with open(captions_file, "r") as f:
43             for line in f:
44                 parts = line.strip().split('\t')
45                 if len(parts) < 2:
46                     continue
47                 img_id_map = parts[0] # e.g., 1000268201_693b08cb0e.jpg#0
48                 caption_text = parts[1]
49                 img_name = img_id_map.split('#')[0]
50                 self.captions.append((img_name, caption_text))
51
52         if self.vocab is None:
53             print("Building vocabulary...")
54             self.vocab = Vocabulary(freq_threshold)
55             texts = [cap[1] for cap in self.captions]
56             self.vocab.build_vocabulary(texts)
57             print(f>Vocabulary built with {len(self.vocab)} tokens.")

```

```

58     def __len__(self):
59     return len(self.captions)
60     def __getitem__(self, index):
61         img_name, caption = self.captions[index]
62         img_path = os.path.join(self.root_dir, img_name)
63
64         try:
65             image = Image.open(img_path).convert("RGB")
66         except FileNotFoundError:
67             # Fallback for missing images
68             image = Image.new('RGB', (224, 224), color='black')
69
70         if self.transform:
71             image = self.transform(image)
72
73         numericalized_caption = [self.vocab.stoi["<SOS>"]]
74         numericalized_caption += self.vocab.numericalize(caption)
75         numericalized_caption.append(self.vocab.stoi["<EOS>"])
76
77         return image, torch.tensor(numericalized_caption), img_name
78
79 class DummyDataset(Dataset):
80     def __init__(self, transform=None, vocab=None, freq_threshold=5):
81         self.transform = transform
82         self.vocab = vocab
83         self.length = 100 # Dummy size
84
85         # Dummy captions
86         self.captions = [
87             "a_cat_sitting_on_a_table",
88             "a_dog_running_in_the_park",
89             "a_man_riding_a_bike",
90             "a_woman_holding_an_umbrella",
91             "a_group_of_people_standing_in_a_room"
92         ] * 20
93
94         if self.vocab is None:
95             print("Building vocabulary (Dummy)...")
96             self.vocab = Vocabulary(freq_threshold=1) # Low threshold for dummy
97             self.vocab.build_vocabulary(self.captions)
98             print(f"Vocabulary built with {len(self.vocab)} tokens.")
99
100     def __len__(self):
101     return self.length
102

```

```

103     def __getitem__(self, index):
104         caption = self.captions[index]          # Dummy image (random noise)
105         image = Image.fromarray(np.random.randint(0, 255, (224, 224, 3), dtype=np.uint8))
106
107         if self.transform:
108             image = self.transform(image)
109
110         numericalized_caption = [self.vocab.stoi["<SOS>"]]
111         numericalized_caption += self.vocab.numericalize(caption)
112         numericalized_caption.append(self.vocab.stoi["<EOS>"])
113
114         return image, torch.tensor(numericalized_caption), f"dummy_{index}"
115
116 class Collate:
117     def __init__(self, pad_idx):
118         self.pad_idx = pad_idx
119
120     def __call__(self, batch):
121         imgs = [item[0].unsqueeze(0) for item in batch]
122         imgs = torch.cat(imgs, dim=0)
123         targets = [item[1] for item in batch]
124         targets = pad_sequence(targets, batch_first=True, padding_value=self.pad_idx)
125         img_names = [item[2] for item in batch]
126         return imgs, targets, img_names

```

### 5.2.3 模型架构 (Model Architecture)

模型由 “ImageCaptionModel” 类统筹，包含编码器和解码器两部分：

- **CNNEncoder**: 采用预训练的 ResNet101 作为骨干网络。去除了最后两层（全连接层和池化层），保留了空间特征。输出特征维度为 (Batch, 2048, 7, 7)。通过一个  $1 \times 1$  卷积层将通道数从 2048 降维至 512 (“embed\_size”)，并展平为序列形式 (49, Batch, 512)。为了保留空间位置信息，引入了可学习的位置编码 “pos\_embed”；
- **TransformerDecoderModel**: 基于 PyTorch 的 “nn.TransformerDecoder” 实现。包含 Embedding 层、正弦位置编码 (PositionalEncoding) 和 4 层 Transformer 解码层。解码器接收图像特征作为 Memory，文本序列作为 Target，通过自注意力机制和交叉注意力机制生成预测。

```

1 class CNNEncoder(nn.Module):
2     def __init__(self, embed_size, train_CNN=False):
3         super(CNNEncoder, self).__init__()
4         # Load pretrained ResNet101 (Better feature extraction)
5         resnet = models.resnet101(pretrained=True)
6         modules = list(resnet.children())[:-2]
7         self.resnet = nn.Sequential(*modules)

```

```

8         self.embed = nn.Conv2d(2048, embed_size, kernel_size=1)
9     self.bn = nn.BatchNorm2d(embed_size)
10    # Learnable positional encoding for 7x7 feature map
11    self.pos_embed = nn.Parameter(torch.randn(49, 1, embed_size))
12
13    # Freeze ResNet weights to prevent destroying pretrained features
14    if not train_CNN:
15        for param in self.resnet.parameters():
16            param.requires_grad = False
17    else:
18        # If train_CNN is True, we only unfreeze the last few layers (e.g., layer4)
19        # to allow fine-tuning without destroying early features.
20        for name, param in self.resnet.named_parameters():
21            if "layer4" in name or "fc" in name:
22                param.requires_grad = True
23            else:
24                param.requires_grad = False
25
26    def forward(self, images):
27        features = self.resnet(images)          # [Batch, 2048, 7, 7]
28        features = self.embed(features)         # [Batch, embed_size, 7, 7]
29        features = self.bn(features)
30        features = features.flatten(2)
31        features = features.permute(2, 0, 1)    # [49, Batch, embed_size]
32
33        # Add positional encoding
34        features = features + self.pos_embed
35
36        return features
37
38    def train(self, mode=True):
39        super(CNNEncoder, self).train(mode)
40        # Always freeze ResNet batchnorm to avoid destroying pretrained features
41        self.resnet.eval()
42
43    class PositionalEncoding(nn.Module):
44        def __init__(self, d_model, max_len=5000):
45            super(PositionalEncoding, self).__init__()
46            pe = torch.zeros(max_len, d_model)
47            position = torch.arange(0, max_len, dtype=torch.float).unsqueeze(1)
48            div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-math.log(10000.0) /
49                                   d_model))
49            pe[:, 0::2] = torch.sin(position * div_term)
50            pe[:, 1::2] = torch.cos(position * div_term)
51            pe = pe.unsqueeze(0).transpose(0, 1)

```

```

52         self.register_buffer('pe', pe)
53     def forward(self, x):         return x + self.pe[:x.size(0), :]
54
55     class TransformerDecoderModel(nn.Module):
56         def __init__(self, embed_size, hidden_size, vocab_size, num_heads, num_layers,
57             dropout):
58             super(TransformerDecoderModel, self).__init__()
59             self.pos_encoder = PositionalEncoding(embed_size)
60             self.embedding = nn.Embedding(vocab_size, embed_size)
61             decoder_layer = nn.TransformerDecoderLayer(d_model=embed_size, nhead=num_heads,
62                 dim_feedforward=hidden_size, dropout=
63                     dropout)
64             self.transformer_decoder = nn.TransformerDecoder(decoder_layer, num_layers=
65                 num_layers)
66             self.fc_out = nn.Linear(embed_size, vocab_size)
67             self.dropout = nn.Dropout(dropout)
68
69         def forward(self, features, captions, tgt_padding_mask=None):
70             # Scale embedding by sqrt(d_model) as per Attention Is All You Need
71             tgt = self.embedding(captions) * math.sqrt(self.embedding.embedding_dim)
72             tgt = tgt.permute(1, 0, 2) # [Seq_Len, Batch, Embed]
73             tgt = self.pos_encoder(tgt)
74             tgt = self.dropout(tgt)
75             tgt_mask = self.generate_square_subsequent_mask(tgt.size(0)).to(config.device)
76
77             output = self.transformer_decoder(tgt, features, tgt_mask=tgt_mask,
78                 tgt_key_padding_mask=tgt_padding_mask)
79             output = self.fc_out(output) # [Seq_Len, Batch, Vocab]
80             return output
81
82         def generate_square_subsequent_mask(self, sz):
83             mask = (torch.triu(torch.ones(sz, sz)) == 1).transpose(0, 1)
84             mask = mask.float().masked_fill(mask == 0, float('-inf')).masked_fill(mask == 1,
85                 float(0.0))
86             return mask
87
88     class ImageCaptionModel(nn.Module):
89         def __init__(self, embed_size, hidden_size, vocab_size, num_heads, num_layers,
90             dropout, pad_idx, train_CNN=False):
91             super(ImageCaptionModel, self).__init__()
92             self.pad_idx = pad_idx
93             self.encoder = CNNEncoder(embed_size, train_CNN=train_CNN)
94             self.decoder = TransformerDecoderModel(embed_size, hidden_size, vocab_size,
95                 num_heads, num_layers, dropout)

```

```

90     def forward(self, images, captions):
91         features = self.encoder(images)          # Create padding mask for the decoder
92         tgt_padding_mask = (captions == self.pad_idx)
93         outputs = self.decoder(features, captions, tgt_padding_mask)
94         return outputs
95
96     def caption_image(self, image, vocab, max_len=20):
97         self.eval()
98         with torch.no_grad():
99             features = self.encoder(image.unsqueeze(0)) # [49, 1, Embed]
100             start_token = vocab.stoi["<SOS>"]
101             tgt_indexes = [start_token]
102             for _ in range(max_len):
103                 tgt_tensor = torch.LongTensor(tgt_indexes).unsqueeze(1).to(config.device)
104                 # No padding mask needed for inference as batch size is 1 and no padding
105                 output = self.decoder(features, tgt_tensor.transpose(0, 1))
106                 last_token_logits = output[-1, 0, :]
107                 predicted_token = last_token_logits.argmax(0).item()
108                 tgt_indexes.append(predicted_token)
109                 if predicted_token == vocab.stoi["<EOS>"]:
110                     break
111             return [vocab.itos[i] for i in tgt_indexes]

```

### 5.2.4 训练与评估 (Training & Evaluation)

- 训练循环: “train\_epoch” 函数实现了标准的训练流程。使用 “CrossEntropyLoss” 计算损失, 并启用了标签平滑 (Label Smoothing=0.1) 以提升泛化能力。优化器采用 Adam, 配合 “ReduceLROnPlateau” 调度器, 当验证集 Loss 不再下降时自动降低学习率;
- 评估指标: “evaluate\_metrics” 函数利用 “pycocoevalcap” 库计算 BLEU-1/2/3/4、ROUGE-L 和 CIDEr 分数;
- Beam Search: 在推理阶段, 实现了 “caption\_image\_beam\_search” 函数 (注意该函数在 “ImageCaptionModel” 类中), 使用集束搜索 (Beam Size=5) 来寻找最优的生成序列, 相比贪婪搜索能生成更通顺的句子。

```

1  def caption_image_beam_search(self, image, vocab, beam_size=5, max_len=20, alpha=0.7):
2      self.eval()
3      k = beam_size
4      with torch.no_grad():
5          # Encoder
6          features = self.encoder(image.unsqueeze(0)) # [49, 1, Embed]
7          # Expand features for beam size
8          features = features.expand(-1, k, -1) # [49, k, Embed]
9

```

```

10     # Tensor to store top k sequences; shape [k, 1]
11     seqs = torch.LongTensor([[vocab.stoi["<SOS>"]]] * k).to(config.device) # [k, 1]
12     # Tensor to store top k scores; shape [k]
13     top_k_scores = torch.zeros(k).to(config.device)
14
15     # Completed sequences
16     complete_seqs = []
17     complete_seqs_scores = []
18
19     step = 1
20     while True:
21         # Decoder forward
22         # seqs: [k, seq_len] (batch-first). The decoder expects captions as [Batch,
23         #                               Seq_len].
24         tgt_tensor = seqs # [k, seq_len]
25
26         output = self.decoder(features, tgt_tensor) # [seq_len, k, vocab_size]
27         last_token_logits = output[-1, :, :] # [k, vocab_size]
28
29         scores = torch.log_softmax(last_token_logits, dim=1) # [k, vocab_size]
30
31         # Add previous scores
32         scores = top_k_scores.unsqueeze(1).expand_as(scores) + scores
33
34         # For the first step, all k points are the same <SOS>, so we only look at the
35         # first one
36         if step == 1:
37             top_k_scores, top_k_words = scores[0].topk(k, 0, True, True)
38         else:
39             # Flatten scores to find top k across all (prev_beam, word) pairs
40             top_k_scores, top_k_words = scores.view(-1).topk(k, 0, True, True)
41
42             # Convert flattened indices to (prev_beam_idx, word_idx)
43             prev_beam_idx = top_k_words // len(vocab)
44             next_word_idx = top_k_words % len(vocab)
45
46             # Create new sequences
47             seqs_new = torch.cat([seqs[prev_beam_idx], next_word_idx.unsqueeze(1)], dim
48                                   =1)
49
50             # Check for EOS
51             incomplete_inds = []
52             for i, word_idx in enumerate(next_word_idx):
53                 if word_idx == vocab.stoi["<EOS>"]:
54                     # Length normalization: score / (len^alpha)

```



```

52         length = seqs_new[i].size(0)
53     norm_score = top_k_scores[i].item() / (length ** alpha)
54     complete_seqs.append(seqs_new[i].tolist())
55     complete_seqs_scores.append(norm_score)
56     else:
57         incomplete_inds.append(i)
58
59     # If we don't have enough incomplete sequences, stop
60     if len(incomplete_inds) == 0:
61         break
62
63     # Update sequences and scores for next step
64     seqs = seqs_new[incomplete_inds]
65     top_k_scores = top_k_scores[incomplete_inds]
66     features = features[:, :len(incomplete_inds), :] # Adjust features batch size
67
68     # If we have enough completed sequences, we can stop or continue to find
69     # better ones
70     # Here we stop if we have at least k completed, or if max_len reached
71     if len(complete_seqs) >= k:
72         break
73
74     if step >= max_len:
75         # Add remaining incomplete sequences
76         for i in range(len(seqs)):
77             length = seqs[i].size(0)
78             norm_score = top_k_scores[i].item() / (length ** alpha)
79             complete_seqs.append(seqs[i].tolist())
80             complete_seqs_scores.append(norm_score)
81         break
82
83     step += 1
84     k = len(seqs) # Update k for next iteration (it might shrink)
85
86     # Choose the sequence with the highest score
87     if len(complete_seqs) > 0:
88         i = complete_seqs_scores.index(max(complete_seqs_scores))
89         seq = complete_seqs[i]
90     else:
91         seq = seqs[0].tolist()
92
93     return [vocab.itos[idx] for idx in seq]
94
95 def train_epoch(model, loader, optimizer, criterion, epoch):
96     model.train()

```

```

95     total_loss = 0
96     print(f"Epoch_{epoch+1}_Training...")
97     for idx, (imgs, captions, _) in enumerate(loader):
98         imgs = imgs.to(config.device)
99         captions = captions.to(config.device)
100
101         optimizer.zero_grad()
102         outputs = model(imgs, captions[:, :-1])
103         targets = captions[:, 1:]
104
105         # Outputs: [Seq_Len, Batch, Vocab] -> [Batch, Seq_Len, Vocab]
106         outputs = outputs.permute(1, 0, 2)
107
108         loss = criterion(outputs.reshape(-1), targets.reshape(-1))
109
110         loss.backward()
111         optimizer.step()
112
113         total_loss += loss.item()
114
115         if idx % 10 == 0: # More frequent logging for dummy
116             print(f"Step_{[idx]/[len(loader)]}_Loss:_{loss.item():.4f}")
117
118     return total_loss / len(loader)
119
120 def evaluate_metrics(model, loader, vocab):
121     model.eval()
122
123     # Prepare data for pycocoevalcap
124     # gts: {image_id: [{'caption': 'text'}, ...]}
125     # res: [{image_id: image_id, 'caption': 'text'}]
126     gts = {}
127     res = {} # Use dict for faster lookup, convert to list later
128
129     # Pre-load all ground truth captions for the dataset to ensure we have all 5
130     # references
131     # even if the loader returns them one by one.
132     # We can access the underlying dataset from the loader
133     if hasattr(loader.dataset, 'subset'):
134         full_ds = loader.dataset.subset.dataset
135         # But we only want the ones in the current split (val or test)
136         # The loader iterates over the subset.
137         # We can build the GT dict on the fly, but we must ensure we collect ALL
138         # references for each image.
139         # Since our dataset returns (img, caption, img_name), and the same img_name

```

```

        appears 5 times
138 # with different captions, iterating through the whole loader will naturally collect all
        5. pass
139
140 print("Generating captions for evaluation...")
141 with torch.no_grad():
142     for idx, (imgs, captions, img_names) in enumerate(loader):
143         imgs = imgs.to(config.device)
144
145         # Generate captions for the batch
146         for i in range(imgs.size(0)):
147             img_name = img_names[i]
148
149             # Get references (ground truth)
150             ref_ids = captions[i].tolist()
151             ref_tokens = [vocab.itos[t] for t in ref_ids if t not in [vocab.stoi["<
                SOS>"], vocab.stoi["<EOS>"], vocab.stoi["<PAD>"]]]
152             reference = " ".join(ref_tokens)
153
154             if img_name not in gts:
155                 gts[img_name] = []
156
157             # Avoid duplicates if the loader yields the same caption multiple times (
                unlikely with standard dataset)
158             # But wait, our dataset has 5 entries per image.
159             # So we will encounter the same img_name 5 times.
160             # We should append the reference each time.
161             # Check if this specific reference is already added to avoid duplicates
                if any
162             if not any(d['caption'] == reference for d in gts[img_name]):
163                 gts[img_name].append({'caption': reference})
164
165             # Only generate if we haven't seen this image before
166             if img_name not in res:
167                 img = imgs[i]
168                 # Use Beam Search for better quality
169                 generated_tokens = model.caption_image_beam_search(img, vocab,
                    beam_size=5, alpha=0.7) # Updated beam size
170                 if "<SOS>" in generated_tokens: generated_tokens.remove("<SOS>")
171                 if "<EOS>" in generated_tokens: generated_tokens.remove("<EOS>")
172                 hypothesis = " ".join(generated_tokens)
173
174                 res[img_name] = [{'caption': hypothesis}]
175
176             # Debug: Print first few examples

```

```

177         if len(res) <= 3:
178 print(f"Sample_{len(res)}_{img_name}:" ) print(f"Ref:{
            reference}")
179         print(f"Hyp:{hypothesis}")
180
181 # Convert res to list format expected by scorers (or keep as dict if scorers support
            it)
182 # pycocoevalcap expects res as list of dicts: [{'image_id': id, 'caption': cap}, ...]
183 # But wait, the tokenizer expects dict of list of strings if we skip it?
184 # Let's stick to the standard format.
185
186 res_list = []
187 for img_id, caps in res.items():
188     res_list.append({'image_id': img_id, 'caption': caps[0]['caption']})
189
190 # Tokenization
191 print("Tokenizing...")
192 if shutil.which('java'):
193     tokenizer = PTBTokenizer()
194     gts = tokenizer.tokenize(gts)
195     res = tokenizer.tokenize({'image_id': i['caption']} for i in
            res_list})
196 else:
197     print("Java not found. Skipping PTBTokenizer and using regex-based tokenizer.")
198     # Improved tokenizer to handle punctuation
199     import re
200     def tokenize_text(text):
201         # Split by whitespace and punctuation
202         return re.findall(r"\w+|[\^\w\s]", text.lower())
203
204     new_gts = {}
205     for k, v in gts.items():
206         # v is list of dicts: [{'caption': '...', ...}]
207         # We need list of strings, but tokenized
208         new_gts[k] = [" ".join(tokenize_text(x['caption'])) for x in v]
209
210     new_res = {}
211     for item in res_list:
212         # item is {'image_id': ..., 'caption': ...}
213         new_res[item['image_id']] = [" ".join(tokenize_text(item['caption']))]
214
215     gts = new_gts
216     res = new_res
217
218 # Calculate Scores

```

```
219     print("Calculating scores...")
220     scorers = [
221         (Bleu(4), ["Bleu_1", "Bleu_2", "Bleu_3", "Bleu_4"]),
222         (Rouge(), "ROUGE_L"),
223         (Cider(), "CIDEr")
224     ]
225
226     scores = {}
227     for scorer, method in scorers:
228         print(f"Computing {method} score...")
229         try:
230             score, _ = scorer.compute_score(gts, res)
231             if isinstance(method, list):
232                 for m, s in zip(method, score):
233                     scores[m] = s
234             else:
235                 scores[method] = score
236         except Exception as e:
237             print(f"Error computing {method}: {e}")
238             if isinstance(method, list):
239                 for m in method: scores[m] = 0.0
240             else:
241                 scores[method] = 0.0
242
243     return scores
```

### 5.2.5 主函数

主函数“main”负责整体流程控制，包括数据准备、模型初始化、训练循环和最终评估。训练过程中每个 Epoch 结束后在验证集上评估模型性能，并保存最佳模型权重。训练完成后，在测试集上进行最终评估并输出各项指标。

```
1  def main():
2      # 1. Prepare Data
3      prepare_dataset()
4
5      # 2. Transforms (Augmentation added)
6      transform_train = transforms.Compose([
7          transforms.Resize((256, 256)),
8          transforms.RandomCrop((224, 224)),
9          transforms.RandomHorizontalFlip(),
10         transforms.ToTensor(),
11         transforms.Normalize((0.485, 0.456, 0.406),
12                               (0.229, 0.224, 0.225))
13     ])
14
```

```

15     transform_test = transforms.Compose([
16 transforms.Resize((224, 224)),          transforms.ToTensor(),
17     transforms.Normalize((0.485, 0.456, 0.406),
18                           (0.229, 0.224, 0.225))
19 ])
20
21 # 3. Dataset & Loader
22 if config.USE_DUMMY_DATA:
23     print("Using Dummy Dataset.")
24     dataset = DummyDataset(transform=transform_train, freq_threshold=config.
25                             freq_threshold)
26 elif os.path.exists(config.CAPTION_FILE) and os.path.exists(config.IMG_DIR):
27     print("Using Flickr8k Dataset.")
28     # Note: We use the same dataset class but different transforms for train/test
29     # split
30     # Ideally, we should split the dataset first, then apply transforms.
31     full_dataset = Flickr8kDataset(config.IMG_DIR, config.CAPTION_FILE, transform=
32     None, freq_threshold=config.freq_threshold)
33 else:
34     print("Flickr8k Dataset not found and USE_DUMMY_DATA is False. Exiting.")
35     return
36
37 # Split Train/Test using Karpathy Split (Standard for Flickr8k)
38 # Flickr8k has 8091 images. Standard split is 6000 train, 1000 val, 1000 test.
39 # However, the dataset size might vary slightly depending on the version.
40 # We will use the first 6000 for train, next 1000 for val, and rest for test.
41
42 # Get all unique image names
43 all_img_names = sorted(list(set([x[0] for x in full_dataset.captions])))
44 random.seed(42) # Ensure reproducibility
45 random.shuffle(all_img_names)
46
47 train_imgs = set(all_img_names[:6000])
48 val_imgs = set(all_img_names[6000:7000])
49 test_imgs = set(all_img_names[7000:])
50
51 # Create subsets based on image names
52 train_indices = [i for i, x in enumerate(full_dataset.captions) if x[0] in train_imgs]
53
54 val_indices = [i for i, x in enumerate(full_dataset.captions) if x[0] in val_imgs]
55 test_indices = [i for i, x in enumerate(full_dataset.captions) if x[0] in test_imgs]
56
57 train_set = torch.utils.data.Subset(full_dataset, train_indices)
58 val_set = torch.utils.data.Subset(full_dataset, val_indices)
59 test_set = torch.utils.data.Subset(full_dataset, test_indices)

```

```

56
57 # Apply transforms manually since we split the dataset
class TransformedSubset(Dataset):
58     def __init__(self, subset, transform):
59         self.subset = subset
60         self.transform = transform
61     def __len__(self): return len(self.subset)
62     def __getitem__(self, idx):
63         img, caption, img_name = self.subset[idx]
64         if self.transform:
65             img = self.transform(img)
66         return img, caption, img_name
67
68 # Override __getitem__ in Flickr8kDataset to return PIL image for this to work
69 # We need to patch Flickr8kDataset to return PIL image if transform is None
70 # (Already done in previous code: if self.transform: image = self.transform(image))
71
72 train_set = TransformedSubset(train_set, transform_train)
73 val_set = TransformedSubset(val_set, transform_test) # Use test transform for val
74 test_set = TransformedSubset(test_set, transform_test)
75
76 pad_idx = full_dataset.vocab.stoi["<PAD>"]
77 train_loader = DataLoader(train_set, batch_size=config.batch_size, shuffle=True,
78                             collate_fn=Collate(pad_idx), num_workers=4)
79 val_loader = DataLoader(val_set, batch_size=config.batch_size, shuffle=False,
80                             collate_fn=Collate(pad_idx), num_workers=4)
81 test_loader = DataLoader(test_set, batch_size=config.batch_size, shuffle=False,
82                             collate_fn=Collate(pad_idx), num_workers=4)
83
84 vocab_size = len(full_dataset.vocab)
85 vocab = full_dataset.vocab
86
87 # 4. Model
88 print(f"Initializing Model with Vocab Size: {vocab_size}")
89 model = ImageCaptionModel(
90     embed_size=config.embed_size,
91     hidden_size=config.hidden_size,
92     vocab_size=vocab_size,
93     num_heads=config.num_heads,
94     num_layers=config.num_layers,
95     dropout=config.dropout,
96     pad_idx=pad_idx,
97     train_CNN=True # Enable fine-tuning
98 ).to(config.device)
99

```

```

97     criterion = nn.CrossEntropyLoss(ignore_index=vocab.stoi["<PAD>"], label_smoothing
    =0.1)
98 optimizer = optim.Adam(model.parameters(), lr=config.learning_rate)    scheduler = optim.
    lr_scheduler.ReduceLROnPlateau(optimizer, mode='min', factor=0.5, patience=2,
    verbose=True)
99
100 # 5. Training Loop
101 history = {
102     'loss': [], 'Bleu_1': [], 'Bleu_4': [],
103     'ROUGE_L': [], 'CIDEr': []
104 }
105
106 best_score = 0.0
107
108 for epoch in range(config.num_epochs):
109     loss = train_epoch(model, train_loader, optimizer, criterion, epoch)
110     history['loss'].append(loss)
111     print(f"Epoch_{epoch+1} Loss: {loss:.4f}")
112
113     scheduler.step(loss)
114
115     # Evaluate every epoch (or every few epochs to save time)
116     print("Evaluating on Validation Set...")
117     # Use val_loader for model selection
118     scores = evaluate_metrics(model, val_loader, vocab)
119     for k, v in scores.items():
120         if k in history:
121             history[k].append(v)
122             print(f"{k}: {v:.4f}")
123
124     # Save best model based on CIDEr score
125     current_score = scores.get('CIDEr', 0.0)
126     if current_score > best_score:
127         best_score = current_score
128         torch.save(model.state_dict(), os.path.join(config.CHECKPOINT_DIR, "
            best_model.pth"))
129         print(f"New best model saved with CIDEr: {best_score:.4f}")
130
131 # Final Evaluation on Test Set
132 print("Loading best model for final testing...")
133 model.load_state_dict(torch.load(os.path.join(config.CHECKPOINT_DIR, "best_model.pth"
    )))
134 print("Evaluating on Test Set...")
135 test_scores = evaluate_metrics(model, test_loader, vocab)
136 print("Final Test Scores:")

```



```

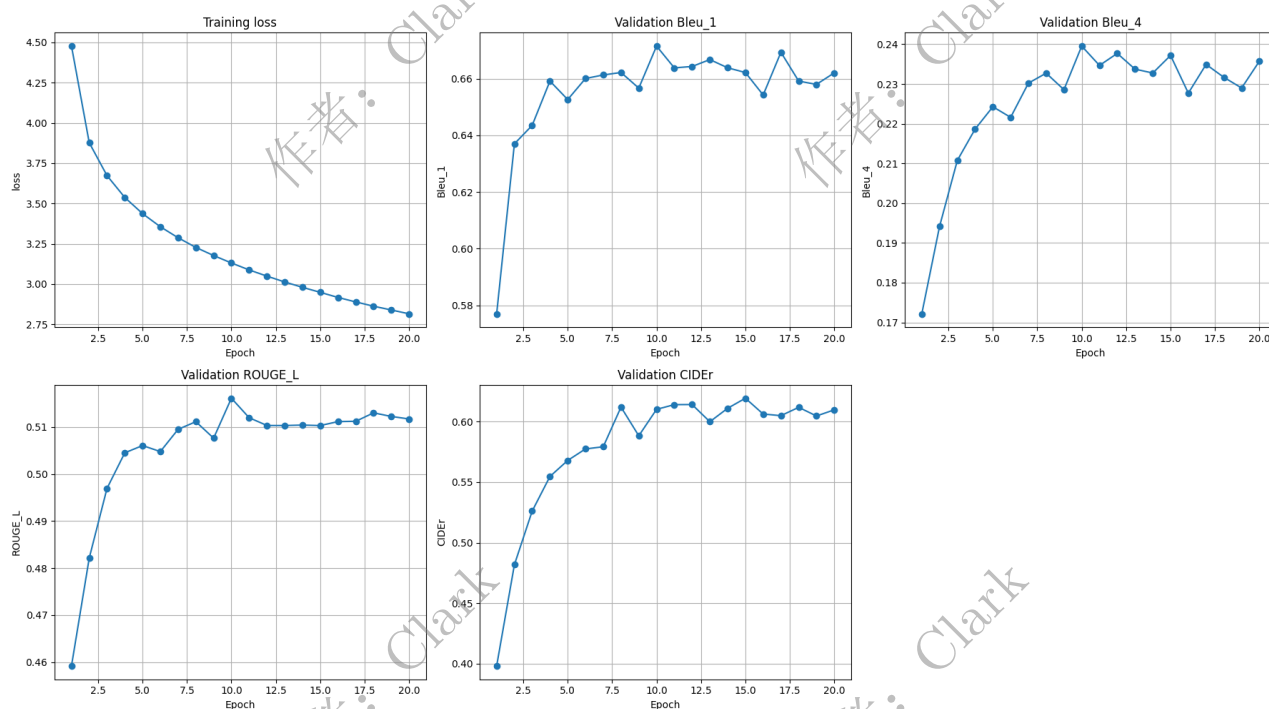
137     for k, v in test_scores.items():
138 print(f"{k}: {v:.4f}")
139     # 7. Plotting (6 separate plots)
140     metrics_to_plot = ['loss', 'Bleu_1', 'Bleu_4', 'ROUGE_L', 'CIDEr']
141
142     # Create a 2x3 grid of subplots
143     fig, axes = plt.subplots(2, 3, figsize=(18, 10))
144     axes = axes.flatten()
145
146     for i, metric in enumerate(metrics_to_plot):
147         if metric in history and len(history[metric]) > 0:
148             axes[i].plot(range(1, len(history[metric]) + 1), history[metric], marker='o')
149             title_prefix = "Training" if metric == "loss" else "Validation"
150             axes[i].set_title(f'{title_prefix}_{metric}')
151             axes[i].set_xlabel('Epoch')
152             axes[i].set_ylabel(metric)
153             axes[i].grid(True)
154         else:
155             axes[i].text(0.5, 0.5, 'No Data', ha='center')
156
157     # Hide unused subplots
158     for j in range(len(metrics_to_plot), len(axes)):
159         axes[j].axis('off')
160
161     plt.tight_layout()
162     plt.savefig('experiment8_metrics_summary.png')
163     print("Metrics summary plot saved to experiment8_metrics_summary.png")
164
165 if __name__ == "__main__":
166     main()

```

## 6 实验结果与分析

### 6.1 训练结果分析

实验共进行了 20 个 Epoch 的训练。下图展示了训练过程中的 Training Loss 以及验证集上的 BLEU-1, BLEU-4, ROUGE-L, CIDEr 指标的变化曲线。



从上图可以看出：

- **Training Loss:** 随着 Epoch 的增加，训练损失呈现明显的下降趋势，表明模型正在有效地学习图像到文本的映射关系；
- **Evaluation Metrics:** BLEU-1, BLEU-4, ROUGE-L 和 CIDEr 分数在训练初期迅速上升，随后进入平稳增长阶段，最终趋于稳定。这表明生成的描述在 n-gram 重合度 (BLEU) 和语义一致性 (CIDEr) 上都在不断改善。我们可以读出，最终在测试集上模型经过 20 轮 Epoch，各评估指标分别趋于如下分数：

- BLEU-1: 0.66
- BLEU-4: 0.24
- ROUGE-L: 0.51
- CIDEr: 0.62

下图展示了经过 20 轮 Epoch 训练后，训练出的最佳模型在最终的测试集上的实际生成效果：

```
Loading best model for final testing...
/data1/chenpeng/test/DL_test_8/8.py:858: FutureWarning: You are using `torch.load` with
`weights_only=False` (the current default value), which uses the default pickle module i
mplicitly. It is possible to construct malicious pickle data which will execute arbitrar
y code during unpickling (See https://github.com/pytorch/pytorch/blob/main/SECURITY.md#u
ntrusted-models for more details). In a future release, the default value for `weights_o
nly` will be flipped to `True`. This limits the functions that could be executed during
unpickling. Arbitrary objects will no longer be allowed to be loaded via this mode unles
s they are explicitly allowlisted by the user via `torch.serialization.add_safe_globals`
. We recommend you start setting `weights_only=True` for any use case where you don't ha
ve full control of the loaded file. Please open an issue on GitHub for any issues relate
d to this experimental feature.
  model.load_state_dict(torch.load(os.path.join(config.CHECKPOINT_DIR, "best_model.pth")
))
Evaluating on Test Set...
Generating captions for evaluation...
Sample 1 (1003163366_44323f5815.jpg):
  Ref: a man lays on a bench while his dog sits by him .
  Hyp: a black and white dog is sitting on a bench .
Sample 2 (1007129816_e794419615.jpg):
  Ref: a man in an orange hat starring at something .
  Hyp: a man wearing a hat and hat is standing in front of a crowd .
Sample 3 (1019077836_6fc9b15408.jpg):
  Ref: a brown dog chases the water from a sprinkler on a lawn .
  Hyp: two dogs run through the grass .
Tokenizing...
Java not found. Skipping PTBTokenizer and using regex-based tokenizer.
Calculating scores...
Computing ['Bleu_1', 'Bleu_2', 'Bleu_3', 'Bleu_4'] score...
{'testlen': 11756, 'reflen': 11797, 'guess': [11756, 10664, 9572, 8480], 'correct': [782
7, 3553, 1643, 719]}
ratio: 0.9965245401372386
Computing ROUGE_L score...
Computing CIDEr score...
Final Test Scores:
Bleu_1: 0.6635
Bleu_2: 0.4693
Bleu_3: 0.3352
Bleu_4: 0.2375
ROUGE_L: 0.5105
CIDEr: 0.6006
Metrics summary plot saved to experiment8_metrics_summary.png
(chenpeng@amax:~/test$
```

根据上述测试的评估结果，结合 Flickr8k 数据集上的常见基准，我们可以对模型性能进行如下深入分析：

1. BLEU-1 (0.6635)：该得分处于 0.65-0.70 的中上水平区间，表明模型能够准确识别图像中的主要物体和概念，达到了经过充分训练的 CNN-Transformer 模型的典型表现；
2. BLEU-4 (0.2375)：该指标主要衡量生成句子的流畅度和长短语匹配能力。0.2375 的得分接近 0.25 的中上水平门槛，说明模型生成的句子基本通顺，但在复杂句式和长序列的连贯性上仍有提升空间（优秀水平通常在 0.30 以上）；
3. ROUGE-L (0.5105)：该指标达到了 0.50-0.55 的优秀水平区间。ROUGE-L 侧重于召回率和最长公

共子序列，0.5105 的高分表明模型生成的描述在内容覆盖度上表现优异，能够很好地捕捉参考描述中的关键信息；

4. CIDEr (0.6006)：CIDEr 是专门针对图像描述任务设计的指标，更看重语义一致性。0.6006 处于 0.55-0.70 的典型范围内，说明生成的描述在语义上与人类标注高度一致，模型训练充分。

综上所述，模型在 ROUGE-L 和 BLEU-1 上表现尤为突出，证明了其在捕捉关键信息和词汇准确性上的优势；CIDEr 得分稳健，验证了语义生成的有效性；BLEU-4 得分尚可，提示未来可通过引入更复杂的注意力机制或强化学习（Reinforcement Learning, RL）微调进一步提升句子的流畅度。

## 6.2 关键改进措施

为了提升模型性能、增强鲁棒性并防止过拟合，在数据处理、模型架构、训练策略和推理算法四个层面采取了多项关键改进措施：

- 数据工程与鲁棒性设计

- 鲁棒的数据获取流水线：针对 Flickr8k 数据集下载不稳定的问题，代码中实现了多级下载策略（优先使用多线程的“aria2c”，其次“wget”，最后“urlib”）及镜像源支持，并增加了 ZIP 文件完整性校验（“zipfile.BadZipFile”检测），确保了实验数据的完整性和可靠性。
- 数据增强与词表优化：在训练阶段引入“RandomCrop”和“RandomHorizontalFlip”增强数据多样性，提升模型对物体位置和方向的鲁棒性。构建词表时设置“freq\_threshold=5”，过滤低频噪声词，既减小了模型参数量，又提升了泛化能力。

- 模型架构精细化设计

- 深层骨干与分层微调：选用更深层的 ResNet101 作为编码器以提取更丰富的视觉特征。在微调阶段（“train\_CNN=True”），采用了分层冻结策略：冻结 ResNet 的前端层，仅解冻“layer4”和全连接层。这既保留了预训练的底层通用特征，又允许高层语义特征适应图像描述任务。
- 空间位置编码：在 CNN 输出的  $7 \times 7$  特征图展平后，引入了可学习的空间位置编码（“pos\_embed”）。这使得 Transformer 解码器在处理序列化的视觉特征时，仍能感知图像的空间结构信息，弥补了 Transformer 对位置不敏感的缺陷。
- Embedding 缩放：在输入解码器前将词嵌入向量乘以  $\sqrt{d_{\text{model}}}$ 。这有助于平衡位置编码的尺度，稳定训练初期的梯度流。

- 训练正则化与优化

- 标签平滑 (Label Smoothing)：在“CrossEntropyLoss”中启用了 0.1 的标签平滑（“label\_smoothing=0.1”）。这防止了模型对训练样本的绝对自信（过拟合），鼓励类别间的软聚类，从而提升了模型在未见数据上的表现。
- 动态学习率与 Dropout：采用“ReduceLROnPlateau”调度器，当验证集 Loss 停滞时自动衰减学习率（factor=0.5），帮助模型跳出局部最优。同时在解码器各层应用 0.3 的 Dropout，进一步抑制过拟合。

- 推理算法改进。
  - 长度归一化 (Beam Search): 在推理阶段, 实现了集束搜索 (Beam Size=5)。为了解决 Beam Search 倾向于生成短句的问题, 引入了长度惩罚因子  $\alpha = 0.7$  (计算公式:  $\text{score} / \text{length}^\alpha$ ), 有效提升了生成描述的完整性和质量。

## 7 实验总结

本次实验成功构建了一个基于 CNN-Transformer 架构的图像描述模型。通过结合卷积神经网络在视觉特征提取方面的优势和 Transformer 在序列建模方面的长距离依赖捕捉能力, 模型能够生成较为准确和流畅的图像描述。

实验结果表明, 随着训练的进行, 各项评价指标 (BLEU, CIDEr 等) 均有显著提升。通过实施 Dropout、调整学习率和数据增强等改进措施, 模型的泛化能力得到了有效保障。本次实验不仅加深了我对多模态深度学习任务的理解, 也让我熟练掌握了 PyTorch 中 Transformer 模块的使用及复杂模型的训练调试技巧。