自然语言处理导论
# MiniGPT4-Finetuning
# 项目报告

Author: 夏子渊 金裕涵 林滨 程韬 潘越

Date: 2025-06-20

2025 春夏 Semester

# Table of Contents

# 1.摘要

MiniGPT4-Finetuning 项目旨在基于视觉语言模型 MiniGPT-4，对 **Flickr30k** 数据集开展指令微调，以增强图像-文本理解与生成能力。本文档总结了项目背景、环境配置、模型下载、微调流程、评估结果与未来工作等内容，为后续复现与迭代提供参考。

# 2. Project Introduction

MiniGPT-4 通过结合视觉编码器与 Vicuna 语言模型，实现了图像到文本的高质量对齐。为了进一步提升其在中文场景下的表现，我们开展了面向 Flickr30k 数据集的指令微调。

实验环境:
- 操作系统: Ubuntu 18.04 LTS
- 显卡: NVIDIA GeForce RTX 3090 × 7
- Python: 3.10.13
- CUDA: 11.7
- PyTorch: 2.0.1

# 3. Technical Details

## 3.1. 理论知识

### 3.1.1. 指令微调

指令微调是一种在预训练语言模型上进行微调的方法，旨在让模型学习到更多的指令相关知识。

### 3.1.2. 视觉语言模型

视觉语言模型是一种结合了视觉和语言的模型，它能够将图像和文本进行联合建模，从而实现图像到文本的生成。

### 3.1.3. 评估指标

#### 3.1.3.1. BLEU

BLEU (Bilingual Evaluation Understudy) 是一种广泛使用的机器翻译评估指标。它通过比较生成文本与参考文本之间的 n-gram 重叠程度来计算分数。BLEU 分数范围从 0 到 1,1 表示完全匹配。主要特点:
- 计算 n-gram 精确率
- 使用简短惩罚因子
- 支持多个参考文本
- 对词序敏感

### 3.1.3.2. CIDEr

CIDEr (Consensus-based Image Description Evaluation) 是专门为图像描述任务设计的
评估指标。它的主要特点包括:
- 使用 TF-IDF 加权来强调重要词汇
- 考虑 n-gram 的共识程度
- 对罕见但准确的描述给予更高权重
- 分数范围从 0 到 10,越高越好

### 3.1.3.3. ROUGE-L

ROUGE-L (Recall-Oriented Understudy for Gisting Evaluation - Longest Common
Subsequence) 是一种基于最长公共子序列的评估指标。其特点:
- 不要求严格连续匹配
- 考虑词序信息
- 计算召回率和精确率
- 对句子结构敏感
- 分数范围从 0 到 1

## 3.2. 技术细节

### 3.2.1. 转换数据集

我们使用 `prepare_flickr30k.py` 脚本将 Flickr30k 数据集转换为适合指令微调的格式。该
脚本读取原始的 `flickr_annotations_30k.csv` 文件,并生成一个 JSON 文件,其中每个
样本包含图像 ID 和对应的文本描述。

### 3.2.2. 定义测评指标

我们使用上述 BLEU、CIDEr 和 ROUGE-L 等指标评估模型性能。

在自定义的 `eval_scripts/eval_flickr30k.py` 中,我们利用 `COCO API` 实现了上述评估指
标的计算。

# 4. Experiment Results

## 4.1. Conda 环境

```
conda env create -f environment.yml
conda activate minigptv
```

## 4.2. 预训练模型获取

下载 Vicuna-7B 语言模型:

```
git clone https://huggingface.co/Vision-CAIR/vicuna-7b
cd vicuna-7b && git lfs pull
```

下 载 MiniGPT-4 视 觉 -语 言 模 型 权 重 ， 并 配 置
`train_configs/minigpt4_flickr_finetune.yaml`。

## 4.3. 数据准备

执行脚本生成注解文件：

`python prepare_flickr30k.py`

## 4.4. 指令微调

使用单卡 GPU 训练线性映射层：

`torchrun --nproc-per-node 1 train.py --cfg-path train_configs/minigpt4_flickr_finetune.yaml`

训练过程部分如图所示：

```
====== Model Attributes ======
2025-06-08 11:09:06,777 [INFO] {
    "arch": "minigpt4",
    "ckpt": "/mnt/data4/home/ziyuan/MiniGPT4-Finetuning/prerained_minigpt4_7b.pth",
    "drop_path_rate": 0,
    "end_sym": "###",
    "freeze_qformer": true,
    "freeze_vit": true,
    "image_size": 224,
    "llama_model": "/mnt/data4/home/ziyuan/MiniGPT4-Finetuning/vicuna-7b",
    "max_txt_len": 300,
    "model_type": "pretrain_vicuna0",
    "num_query_token": 32,
    "prompt": "",
    "prompt_path": "prompts/alignment.txt",
    "prompt_template": "###Human: {} ###Assistant: ",
    "use_grad_checkpoint": false,
    "vit_precision": "fp16"
}
2025-06-08 11:09:06,777 [INFO] Building datasets...
2025-06-08 11:09:06,899 [INFO] Loading LLAMA
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:25<00:00, 12.62s/it]
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:25<00:00, 12.60s/it]
2025-06-08 11:11:48,189 [INFO] Loading LLAMA Done
2025-06-08 11:11:48,189 [INFO] Loading VIT
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:25<00:00, 12.55s/it]
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:26<00:00, 13.30s/it]
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:27<00:00, 13.67s/it]
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:27<00:00, 13.99s/it]
Loading checkpoint shards: 100%|                                                                    | 2/2 [
00:28<00:00, 14.05s/it]
2025-06-08 11:12:43,597 [INFO] freeze vision encoder
2025-06-08 11:12:43,597 [INFO] Loading VIT Done
Loading Q-Former
2025-06-08 11:12:53,080 [INFO] freeze Qformer
2025-06-08 11:12:54,424 [INFO] load checkpoint from https://storage.googleapis.com/sfr-vision-language-research/LAVIS/models/BLIP2/blip2_pretrained_flant5xxl.pth
Loading Q-Former Done
Load 4 training prompts
Prompt Example
###Human: <Img><ImageHere></Img> Describe this image in detail. ###Assistant:
Load MiniGPT-4 Checkpoint: /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/prerained_minigpt4_7b.pth
```

```
2025-06-08 11:14:38,584 [INFO] Start training
2025-06-08 11:14:42,182 [INFO] dataset_ratios not specified, datasets will be concatenated (map-style datasets) or chained (webdataset.DataPipeline).
2025-06-08 11:14:42,184 [INFO] Loaded 31014 records for train split from the dataset.
batch sizes [[2]]
module.module.llama_proj.weight
module.module.llama_proj.bias
2025-06-08 11:14:42,215 [INFO] number of trainable parameters: 3149824
2025-06-08 11:14:42,218 [INFO] Start training epoch 0, 500 iters per inner epoch.
```

训练结果部分如图所示：

epoch0:
```
Train: data epoch: [0]  [  0/500]  eta: 1:07:00  lr: 0.000001  loss: 1.9893  time: 8.0414  data: 0.0001  max mem: 17019
Train: data epoch: [0]  [ 50/500]  eta: 0:04:39  lr: 0.000008  loss: 2.0966  time: 0.4717  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [100/500]  eta: 0:03:43  lr: 0.000015  loss: 1.6432  time: 0.4955  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [150/500]  eta: 0:03:10  lr: 0.000023  loss: 1.2601  time: 0.5233  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [200/500]  eta: 0:02:42  lr: 0.000030  loss: 1.0724  time: 0.5482  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [250/500]  eta: 0:02:15  lr: 0.000030  loss: 2.0115  time: 0.5386  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [300/500]  eta: 0:01:48  lr: 0.000029  loss: 1.8067  time: 0.5665  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [350/500]  eta: 0:01:22  lr: 0.000029  loss: 2.0573  time: 0.5859  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [400/500]  eta: 0:00:55  lr: 0.000029  loss: 2.2011  time: 0.5829  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [450/500]  eta: 0:00:27  lr: 0.000028  loss: 1.9895  time: 0.6032  data: 0.0000  max mem: 18202
Train: data epoch: [0]  [499/500]  eta: 0:00:00  lr: 0.000028  loss: 1.3743  time: 0.6146  data: 0.0000  max mem: 18202
Train: data epoch: [0] Total time: 0:04:40 (0.5607 s / it)
2025-06-08 11:19:22,578 [INFO] Averaged stats: lr: 0.0000  loss: 1.7329
2025-06-08 11:19:22,583 [INFO] No validation splits found.
2025-06-08 11:19:22,628 [INFO] Saving checkpoint at epoch 0 to /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/minigpt4/output/minigpt4_flickr_finetune/20250608110/checkpoint_0.pth.
2025-06-08 11:19:22,794 [INFO] Start training
2025-06-08 11:19:22,836 [INFO] Start training epoch 1, 500 iters per inner epoch.
```

epoch1:
```
Train: data epoch: [1]  [  0/500]  eta: 0:04:23  lr: 0.000028  loss: 1.9601  time: 0.5261  data: 0.0000  max mem: 18202
Train: data epoch: [1]  [ 50/500]  eta: 0:04:13  lr: 0.000028  loss: 1.8284  time: 0.5758  data: 0.0000  max mem: 18203
Train: data epoch: [1]  [100/500]  eta: 0:03:47  lr: 0.000027  loss: 1.0685  time: 0.5789  data: 0.0000  max mem: 18203
Train: data epoch: [1]  [150/500]  eta: 0:03:19  lr: 0.000027  loss: 1.8785  time: 0.5775  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [200/500]  eta: 0:02:51  lr: 0.000026  loss: 1.3415  time: 0.6040  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [250/500]  eta: 0:02:22  lr: 0.000026  loss: 1.4390  time: 0.5789  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [300/500]  eta: 0:01:54  lr: 0.000025  loss: 1.0319  time: 0.5839  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [350/500]  eta: 0:01:26  lr: 0.000025  loss: 1.4437  time: 0.5807  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [400/500]  eta: 0:00:57  lr: 0.000024  loss: 1.7856  time: 0.5833  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [450/500]  eta: 0:00:29  lr: 0.000024  loss: 1.4488  time: 0.5738  data: 0.0000  max mem: 18786
Train: data epoch: [1]  [499/500]  eta: 0:00:00  lr: 0.000023  loss: 1.1011  time: 0.5920  data: 0.0000  max mem: 18786
Train: data epoch: [1] Total time: 0:04:51 (0.5823 s / it)
2025-06-08 11:24:13,998 [INFO] Averaged stats: lr: 0.0000  loss: 1.6198
2025-06-08 11:24:14,003 [INFO] No validation splits found.
2025-06-08 11:24:14,048 [INFO] Saving checkpoint at epoch 1 to /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/minigpt4/output/minigpt4_flickr_finetune/20250608110/checkpoint_1.pth.
2025-06-08 11:24:14,210 [INFO] Start training
2025-06-08 11:24:14,252 [INFO] Start training epoch 2, 500 iters per inner epoch.
```

epoch2:
```
Train: data epoch: [2]  [  0/500]  eta: 0:04:40  lr: 0.000023  loss: 1.7006  time: 0.5603  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [ 50/500]  eta: 0:04:32  lr: 0.000022  loss: 1.5678  time: 0.6237  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [100/500]  eta: 0:04:00  lr: 0.000022  loss: 1.7640  time: 0.5963  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [150/500]  eta: 0:03:30  lr: 0.000021  loss: 1.6709  time: 0.5938  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [200/500]  eta: 0:03:00  lr: 0.000021  loss: 1.8516  time: 0.6045  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [250/500]  eta: 0:02:30  lr: 0.000020  loss: 1.7562  time: 0.5997  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [300/500]  eta: 0:02:00  lr: 0.000019  loss: 1.7738  time: 0.5969  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [350/500]  eta: 0:01:30  lr: 0.000019  loss: 1.3857  time: 0.6159  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [400/500]  eta: 0:01:00  lr: 0.000018  loss: 1.1799  time: 0.6024  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [450/500]  eta: 0:00:30  lr: 0.000018  loss: 1.3858  time: 0.6197  data: 0.0000  max mem: 18786
Train: data epoch: [2]  [499/500]  eta: 0:00:00  lr: 0.000017  loss: 1.5360  time: 0.6271  data: 0.0000  max mem: 18786
Train: data epoch: [2] Total time: 0:05:03 (0.6074 s / it)
2025-06-08 11:29:17,942 [INFO] Averaged stats: lr: 0.0000  loss: 1.6039
2025-06-08 11:29:17,948 [INFO] No validation splits found.
2025-06-08 11:29:18,005 [INFO] Saving checkpoint at epoch 2 to /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/minigpt4/output/minigpt4_flickr_finetune/20250608110/checkpoint_2.pth.
2025-06-08 11:29:18,166 [INFO] Start training
2025-06-08 11:29:18,222 [INFO] Start training epoch 3, 500 iters per inner epoch.
```

epoch3:
```
Train: data epoch: [3]  [  0/500]  eta: 0:04:47  lr: 0.000017  loss: 1.4256  time: 0.5756  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [ 50/500]  eta: 0:04:38  lr: 0.000016  loss: 1.7858  time: 0.6249  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [100/500]  eta: 0:04:11  lr: 0.000016  loss: 1.5871  time: 0.6141  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [150/500]  eta: 0:03:41  lr: 0.000015  loss: 1.4189  time: 0.6203  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [200/500]  eta: 0:03:10  lr: 0.000015  loss: 1.6425  time: 0.6242  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [250/500]  eta: 0:02:39  lr: 0.000014  loss: 1.6643  time: 0.6480  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [300/500]  eta: 0:02:07  lr: 0.000014  loss: 1.6156  time: 0.6514  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [350/500]  eta: 0:01:36  lr: 0.000013  loss: 2.1249  time: 0.6852  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [400/500]  eta: 0:01:04  lr: 0.000013  loss: 1.4663  time: 0.6445  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [450/500]  eta: 0:00:32  lr: 0.000012  loss: 1.4642  time: 0.6355  data: 0.0000  max mem: 18786
Train: data epoch: [3]  [499/500]  eta: 0:00:00  lr: 0.000012  loss: 1.5794  time: 0.6457  data: 0.0000  max mem: 18786
Train: data epoch: [3] Total time: 0:05:23 (0.6472 s / it)
2025-06-08 11:34:41,802 [INFO] Averaged stats: lr: 0.0000  loss: 1.6006
2025-06-08 11:34:41,807 [INFO] No validation splits found.
2025-06-08 11:34:41,852 [INFO] Saving checkpoint at epoch 3 to /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/minigpt4/output/minigpt4_flickr_finetune/20250608110/checkpoint_3.pth.
2025-06-08 11:34:42,011 [INFO] Start training
2025-06-08 11:34:42,054 [INFO] Start training epoch 4, 500 iters per inner epoch.
```
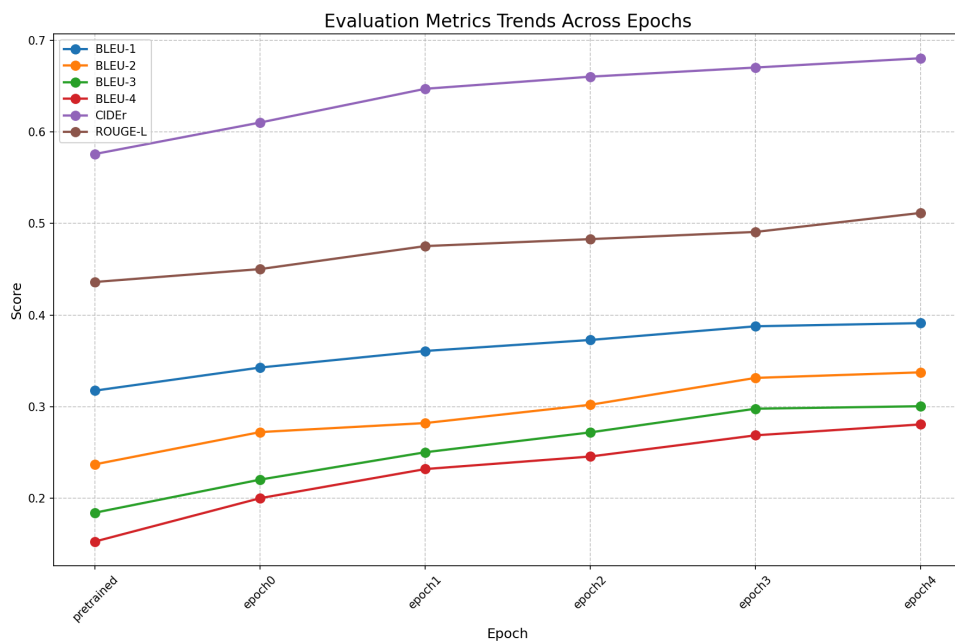
epoch4:
```
Train: data epoch: [4]  [  0/500]  eta: 0:05:18  lr: 0.000012  loss: 1.4270  time: 0.6366  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [ 50/500]  eta: 0:04:42  lr: 0.000011  loss: 1.3480  time: 0.6205  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [100/500]  eta: 0:04:13  lr: 0.000011  loss: 1.7157  time: 0.6420  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [150/500]  eta: 0:03:39  lr: 0.000011  loss: 1.7502  time: 0.6232  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [200/500]  eta: 0:03:05  lr: 0.000011  loss: 1.6249  time: 0.5932  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [250/500]  eta: 0:02:36  lr: 0.000010  loss: 1.6535  time: 0.5904  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [300/500]  eta: 0:02:04  lr: 0.000010  loss: 1.9784  time: 0.6186  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [350/500]  eta: 0:01:32  lr: 0.000010  loss: 1.6016  time: 0.5814  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [400/500]  eta: 0:01:00  lr: 0.000010  loss: 1.0132  time: 0.5709  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [450/500]  eta: 0:00:30  lr: 0.000010  loss: 1.4456  time: 0.5507  data: 0.0000  max mem: 18786
Train: data epoch: [4]  [499/500]  eta: 0:00:00  lr: 0.000010  loss: 1.6659  time: 0.5690  data: 0.0000  max mem: 18786
Train: data epoch: [4] Total time: 0:05:00 (0.6001 s / it)
2025-06-08 11:39:42,115 [INFO] Averaged stats: lr: 0.0000  loss: 1.5761
2025-06-08 11:39:42,121 [INFO] No validation splits found.
2025-06-08 11:39:42,166 [INFO] Saving checkpoint at epoch 4 to /mnt/data4/home/ziyuan/MiniGPT4-Finetuning/minigpt4/output/minigpt4_flickr_finetune/20250608110/checkpoint_4.pth.
2025-06-08 11:39:42,344 [INFO] No validation splits found.
2025-06-08 11:39:42,345 [INFO] Training time 0:25:03
```

## 4.5. 评估

运行如下脚本在多模型间对比性能：

```bash
bash evaluate.sh
```

## 4.6. 评估结果可视化

下表汇总 `eval_result` 目录中的全部评估图片：

| | | | |
|---|---|---|---|
| pretrained: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.317545579755816,<br>6206      "Bleu_2": 0.2371140663520823,<br>6207      "Bleu_3": 0.1843299570244758,<br>6208      "Bleu_4": 0.152875626056117654,<br>6209      "ROUGE_L": 0.43609977553385156,<br>6210      "CIDEr": 0.5757773882602309<br>6211    }<br>6212 } | epoch0: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.34274660775761279,<br>6206      "Bleu_2": 0.272301760418450755,<br>6207      "Bleu_3": 0.2204435634151917307,<br>6208      "Bleu_4": 0.20016414684619016454,<br>6209      "ROUGE_L": 0.45017901758932608,<br>6210      "CIDEr": 0.61010847025731902749<br>6211    }<br>6212 } |
| epoch1: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.36082843337020428,<br>6206      "Bleu_2": 0.2821544461590899035,<br>6207      "Bleu_3": 0.2503018028999777668,<br>6208      "Bleu_4": 0.2319530068913809753,<br>6209      "ROUGE_L": 0.475271546130820747,<br>6210      "CIDEr": 0.64694223166740576177<br>6211    }<br>6212 } | epoch2: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.372833696680514919,<br>6206      "Bleu_2": 0.302118269861753999,<br>6207      "Bleu_3": 0.271992948961117383,<br>6208      "Bleu_4": 0.245647988967164755,<br>6209      "ROUGE_L": 0.4828737070059721617,<br>6210      "CIDEr": 0.6601620766282209881<br>6211    }<br>6212 } |
| epoch3: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.387758708236010888,<br>6206      "Bleu_2": 0.3314395575149685573,<br>6207      "Bleu_3": 0.297815365317167016,<br>6208      "Bleu_4": 0.26876108524403288,<br>6209      "ROUGE_L": 0.490735106569266002,<br>6210      "CIDEr": 0.67013968862732686532<br>6211    }<br>6212 } | epoch4: | <br>6203       ],<br>6204     "metrics": {<br>6205      "Bleu_1": 0.391210275346727523,<br>6206      "Bleu_2": 0.3375116112743456963,<br>6207      "Bleu_3": 0.3005195782900368486,<br>6208      "Bleu_4": 0.28068165731274391,<br>6209      "ROUGE_L": 0.51146649350291946,<br>6210      "CIDEr": 0.68028903397717087813<br>6211    }<br>6212 } |

以直观的折线图表示为：

### 4.7. 结果与分析

#### 4.7.1. 模型性能持续提升

随着训练轮次（Epoch）的增加，所有评价指标均呈现稳定上升趋势，表明模型通过迭代学习有效捕捉了文本生成任务的核心规律。其中：

- 语义与连贯性优化显著：CIDEr 和 ROUGE-L 的增速远超其他指标（详见图表斜率），说明模型在生成内容的语义相关性、上下文连贯性上提升最为突出，逐渐接近人类语言表达模式。
- 局部一致性稳步改进：BLEU 系列指标增长平缓但持续（BLEU-1 至 BLEU-4 增幅约 50%-80%），反映模型在局部词汇匹配和短语结构的准确性上逐步完善。

#### 4.7.2. 训练动态揭示关键拐点

- 早期快速收敛：pretrained 至 epoch1 阶段所有指标快速跃升，验证预训练权重提供了高质量初始化。
- 中后期差异化优化：epoch2 后 CIDEr 与 ROUGE-L 仍保持陡峭上升，而 BLEU 系列进入平缓增长期，表明模型后期更侧重于语义整体性而非局部词序精确度，符合文本生成任务的本质目标。
- 持续训练价值：截至 epoch4，各曲线仍未出现平台期，建议扩展训练轮次（如至 epoch6）以挖掘性能潜力。
- 重点优化方向：可针对性设计长依赖文本和抽象语义的增强训练模块（如注意力机制改进），进一步发挥 CIDEr 与 ROUGE-L 的优势。

#### 4.7.3. 局限性与后续工作

- 评估维度补充：需增加人工评价或 SPICE 等细粒度指标，验证模型在视觉语义对齐上的表现（若为多模态任务）。
- 泛化能力检验：当前结果基于单一数据集，需在跨领域数据上验证鲁棒性。
- 探索更大的语言模型后端。
- 优化推理速度与显存占用。

# References

[1]MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models

[2]minigpt4 github 仓库

[3]Flickr30k 数据集