



به نام خدا دانشگاه تهران - دانشگاه صنعتی خواجه نصیرالدین طوسی تهران دانشگاه مهندسی برق و کامپیوتر

توصيف عكس

محمدجواد احمدي	نام و نام خانوادگی
4.115	شمارهٔ دانشجویی

قوصيف عكس معتمل المستعمل المست

فهرست مطالب

٣											پاسخ پرسش اول													
٣					•																نسمت ۱	پاسخ ف	1.1	
٣٣																					نسمت ۲	پاسخ ف	۲.۱	

توصيف عكس توصيف عكس

فهرست تصاوير

11	نمایش نمونهدادهٔ آموزش و ارزیابی	١
۱۵	نمونهای از تصویر/کپشن یک دسته از دادهها	۲
۲٧	نمودار خطا در حالت رزنت فریز	٣
۲۸	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ اول	*
۲۸	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ دوم	۵
4	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ سوم	۶
79	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ چهارم	٧
۳.	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ پنجم	٨
٣٣	نمودار خطای آموزش و تست در حالت رزنت فریز	٩
٣٩	نمودار خطا در حالت رزنت بدون فريز	١.
۴.	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ اول	11
۴.	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ دوم	17
41	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ سوم	١٣
41	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ چهارم	14
47	نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ پنجم	۱۵
41	نمودار خطای آموزش و تست در حالت رزنت بدون فریز	18
41	نمودار خطا در حالت رزنت فريز و بدون فريز	1
۵۴	نمودار خطا در حالت رزنت بدون فريز (Bidirectional LSTM)	١٨
۵۵	نمودار خطا در حالت رزنت فریز، بدون فریز و بدون فریز (Bidirectional LSTM)	19

توصيف عكس محمد عكس محمد المحمد المحمد

پرسش ۱. توصیف عکس

١ پاسخ پرسش اول

توضيح پوشهٔ كدهاى توصيف عكس

کدهای مربوط به این قسمت، علاوه بر پوشهٔ محلی کدها در این لینک گوگل کولب آورده شده است. مدلهای دخیره شده هم از طریق این لینک در دسترس هستند.

۱.۱ پاسخ قسمت ۱

در گام اول و برای راحتی، مجموعه داده را روی گوگل درایو بارگذاری می کنیم. هم چنین، فراخوانی این فایل بدون نیاز به Mount کردن و استفاده از دستورات زیر مشکل حل می شود و استفاده از دستورات زیر مشکل حل می شود (منبع). در ادامه هم دستوراتی برای خروج داده ها از حالت فشرده و رفتن به یک پوشهٔ مخصوص مجموعه داده نوشته ایم و در انتها فایل فشرده را حذف کرده ایم.

```
!pip install --upgrade --no-cache-dir gdown
!gdown 1DArAhv1ieTdmvmu0Up8_2u5j4iVgmqn7

!mkdir Dataset
!unzip /content/archive.zip -d Dataset/
```

در ادامه دستوراتی را برای بارگذاری داده ها از یک فایل متنی می نویسیم. ابتدا مسیر فایل متنی برا تعیین کرده و سپس با استفاده از تابع read_csv، داده های موجود در فایل متنی به صورت یک دیتافریم خوانده می شوند. در ادامه، تعداد تصاویر دارای شرح، تعداد تصاویر یکتا به صورت متنی چاپ می شوند. در انتها نیز، با استفاده از تابع head، پنج سطر اول داده های داده های داده می شوند. دستورات و نتیجهٔ خروجی به شرح زیر هستند:

```
import pandas as pd

# Set the path to the captions file

captions_path = 'captions.txt'

# Load the captions file into a pandas DataFrame

captions_df = pd.read_csv(captions_path)

# Count the number of images with captions

num_images = len(captions_df)

# Print the total number of images with captions
```

توصيف عكس ﴿ اللَّهُ ال

```
print(f"Total number of images with captions: {num_images}")
15 # Print the number of unique images
unique_images = len(captions_df['image'].unique())
17 print(f"Number of unique images with captions: {unique_images}")
19 # Print the number of unique captions
20 unique_captions = len(captions_df['caption'].unique())
21 print(f"Number of unique captions: {unique_captions}")
23 # Display the first 5 rows of the DataFrame
24 print("\nSample captions:")
25 captions_df.head()
27 Total number of images with captions: 40455
28 Number of unique images with captions: 8091
29 Number of unique captions: 40201
31 Sample captions:
32 image caption
0 1000268201_693b08cb0e.jpg A child in a pink dress is climbing up a set o...
1 1000268201_693b08cb0e.jpg A girl going into a wooden building .
2 1000268201_693b08cb0e.jpg A little girl climbing into a wooden playhouse .
36 3 1000268201_693b08cb0e.jpg A little girl climbing the stairs to her playh...
37 4 1000268201_693b08cb0e.jpg A little girl in a pink dress going into a woo...
```

در ادامه، دستوراتی شامل یک کلاس به نام Vocabulary می نویسیم که وظیفهٔ ساخت و استفاده از یک مجموعه واژگان را بر عهده دارد. به طور کلی، این کلاس دارای متدهایی برای ساخت واژگان، توکنسازی متن، عددیسازی متن و تشکیل دیکشنری واژگان است. در بخش ساخت واژگان، کلاس Vocabulary مقادیر اولیه را برای توکنهای ویژه (Special Tokens) تعریف می کند. به عنوان مثال PAD>< برای ید کردن دنبالهها، SOS>< برای نشان دادن شروع جمله، PAD>< برای نشان دادن یایان جمله و VNK>< برای توکنهای ناشناخته. سپس با استفاده از تابع tokenizer که با استفاده از یک توکنساز با عبارت منظم رشته را به توکنها تقسیم میکند و آنها را به حروف کوچک تبدیل میکند، و واژگان را براساس آستانه فراوانی میسازد. ورودی frequency_threshold مقدار آستانهٔ فراوانی است که توکنها باید برای اضافه شدن به واژگان داشته باشند. درواقع، آستانهٔ فراوانی مشخص میکند که توکن هایی که تعداد تکرار آنها بیشتر از این آستانه است، در واژگان قرار بگیرند. در هنگام ساخت واژگان، همچنین دیکشنری stoi و stoi بروزرسانی می شوند تا اندیس ها و توکن ها با هم متناظر شوند. متد numericalize با توجه به متن ورودی توکن سازی را انجام داده و توکن ها را به اندیس های متناظر در دیکشنری واژگان تبدیل می کند. اگر توکن در واژگان موجود نباشد، توکن Vocabulary (نامشخص) جایگزین می شود. در ادامه، یک نمونه از کلاس Vocabulary با آستانهٔ فراوانی ۱ ساخته می شود و با استفاده از یک لیست از حملات، واژگان ساخته می شود و دیکشنری واژگان چاپ می شود. سیس یک متن جدید عددیسازی شده و با استفاده از واژگان، نمایش داده می شود. در این نتیجه، دیکشنری واژگان شامل تمام توکن های موجود در جملات است و به هر توكن يك انديس اختصاص داده شده است. همچنين متن نمونه به صورت عددي نمايش داده شده است، به این صورت که هر توکن در متن با مقدار متناظر در دیکشنری واژگان حایگزین شده است. دستورات و نتیجهٔ احرای آن به شرح برنامهٔ ۱ است: م وصيف عکس <u>Garas</u> توصيف عکس

Program 1: Creating a Vocabulary

```
# Import the Vocabulary class
2 from nltk.tokenize import RegexpTokenizer
4 class Vocabulary:
    def __init__(self, frequency_threshold):
      # Special tokens and their corresponding indices
      self.itos = {0: "<PAD>", 1: "<SOS>", 2: "<EOS>", 3: "<UNK>"}
      self.stoi = {"<PAD>": 0, "<SOS>": 1, "<EOS>": 2, "<UNK>": 3}
      self.frequency_threshold = frequency_threshold
    def __len__(self):
      # Return the total number of tokens in the vocabulary
      return len(self.itos)
    def tokenizer(self, text):
      # Tokenize the text using a regular expression tokenizer
17
      tokenizer = RegexpTokenizer(r'\w+')
      return [token.lower() for token in tokenizer.tokenize(text)]
    def build_vocab(self, caption_list):
      frequencies = {}
      idx = 4
      # Iterate over each caption in the list
      for caption in caption_list:
26
        # Tokenize the caption
2.7
        for token in self.tokenizer(caption):
          # Update the token frequencies
          if token not in frequencies:
            frequencies[token] = 1
          else:
            frequencies[token] += 1
          # Check if the token frequency reaches the threshold
          if frequencies[token] == self.frequency_threshold:
            # Add the token to the vocabulary with a new index
            self.stoi[token] = idx
            self.itos[idx] = token
            idx += 1
    def numericalize(self, text):
      # Tokenize the text
43
      tokenized_text = self.tokenizer(text)
```

توصيف عكس ﴿ وَصِيفَ عَكْسَ السَّاعِينَ السَّاعِينَ السَّاعِينَ السَّاعِينَ السَّاعِينَ السَّاعِينَ السَّاعِينَ

```
# Convert tokens to their corresponding indices in the vocabulary
      return [self.stoi[token] if token in self.stoi else self.stoi["<UNK>"] for token in
      tokenized text]
49 # Create an instance of the Vocabulary class with a frequency threshold of 1
50 v = Vocabulary(frequency_threshold=1)
52 # Build the vocabulary using a list of captions
53 v.build_vocab(["I am Mohammad Javad Ahmadi, a student of Dr. Keller's Deep Learning course."])
55 # Print the vocabulary dictionary
56 print(f"Vocabulary dictionary: {v.stoi}")
58 # Numericalize a new text using the vocabulary
59 numericalized_text = v.numericalize("I am Mohammad Javad Ahmadi, a student of Dr. Keller's Deep
      Learning course.")
61 # Print the numericalized text
62 print(f"Numericalized text: {numericalized_text}")
64 Vocabulary dictionary: {'<PAD>': 0, '<SOS>': 1, '<EOS>': 2, '<UNK>': 3, 'i': 4, 'am': 5, '
      mohammad': 6, 'javad': 7, 'ahmadi': 8, 'a': 9, 'student': 10, 'of': 11, 'dr': 12, 'keller':
      13, 's': 14, 'deep': 15, 'learning': 16, 'course': 17}
65 Numericalized text: [4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17]
```

در ادامه، یک کلاس با نام FlickrDataset تعریف می کنیم که از کلاس torch.utils.data.Dataset ارثبری می کند و برای مدیریت مجموعه داده ها در فرمت دیتاست فلیکر استفاده می شود. ابتدا کتابخانه های لازم فراخوانی می شوند. سپس، كلاس Dataset از كتابخانهٔ torch.utils.data وارد مي شود كه يك كلاس پايه براي ايجاد ديتاستهاي شخصي سازي شده در پای تورچ است. همچنین برای کار با تصویر، کلاس Image از کتابخانهٔ PIL وارد می شود. برای کار با تنسورها و سایر موارد و ابزارهای لازم هم تورچ را وارد می کنیم. در ادامه، یک کلاس به نام FlickrDataset تعریف می کنیم که از کلاس Dataset ارثبری می کند. متد init سازندهٔ این کلاس است که در آن مقادیر اولیه متغیرها و ورودی ها تنظیم می شود. سیس، فایل کیشن را با استفاده از تابع read_csv از کتابخانه pandas می خوانیم و به صورت یک جدول داده ها (DataFrame) در متغیر df ذخیره می کنیم. بر اساس نوع داده، تقسیم بندی داده ها را انجام می دهیم. اگر نوع داده train باشد، از ابتدای دیتافریم df تا ۹۰ درصد آن را به عنوان دادهٔ آموزش انتخاب می کنیم. اگر نوع داده test باشد، ۱۰ درصد باقی مانده دیتافریم df را بهنوان دادهٔ ارزیابی انتخاب میکنیم. در تعاریف دقت داریم که بخشی را برای محاسبهٔ اندیس شروع داده های ارزیابی براساس ۱۰ درصد باقی مانده دیتافریم در نظر بگیریم و دادههای ارزیابی را از اندیس شروع تا انتهای دیتافریم انتخاب کنیم. همچنین، اندیسهای ردیفها را دوباره از ابتدا شمارهگذاری میکنیم. مسیر ریشه (root directory) دیتاست را در متغیر self.root_dir ذخیره میکنیم. تبدیلات و ترانسفورمیشن هایی که روی تصاویر باید اعمال شوند را در متغیر self.transform ذخیره میکنیم. این ورودی اختیاری است و در صورتی که ارائه نشود، مقدار None را دارد.: نام تصاویر را از ستون image در دیتافریم self.df استخراج کرده و در متغیر self.imgs ذخیره میکنیم. متنهای کپشن را هم از ستون caption در دیتافریم self.df استخراج کرده و در متغير self.captions ذخيره مي كنيم. يك نمونه از كلاس Vocabulary را با استفاده از مقدار requency_threshold را توصيف عكس ^{كمجمع} توصيف عكس

به عنوان آستانهٔ تكرار استفاده كرده و آن را در متغير self.vocab ذخيره مي كنيم. با استفاده از نمونهٔ self.vocab، لغتنامه را بر اساس متنهای کپشن موجود در self.captions ایجاد میکنیم و آن را ساختاردهی میکنیم. در ادامه، متد len را برای بازگرداندن طول دیتاست که برابر با تعداد نمونهها در self.df است، تعریف میکنیم. متد getitem را هم برای بازگرداندن نمونههای دیتاست به صورت متناظر با اندیس تعریف می کنیم. متن کیشن مربوط به اندیس مورد نظر از متغیر self.captions را استخراج کرده و در متغیر caption ذخیره می کنیم. نام تصویر مربوط به اندیس مورد نظر از متغیر self.imgs را هم استخراج کرده و در متغیر image_name ذخیره می کنیم. ترکیب مسیر ریشه (self.root_dir) و نام تصویر image_name را به عنوان مسیر کامل تصویر در نظر گرفته و در متغیر image_path ذخیره میکنیم. سپس، با استفاده از تابع Image.open از کتابخانه PIL، تصویر را از مسیر image_path باز میکنیم و آن را به حالت RGB تبدیل میکنیم. تصویر در متغیر img ذخیره میشود. هم چنین اگر تبدیلاتی را مدنظر داشته باشیم، اِعمال می کنیم. بردارسازی کپشن را با توجه به نماد شروع SOS آغاز می کنیم و آن را در متغیر vectorized_caption ذخیره می کنیم. متن کپشن را به عدد تبدیل کرده و به متغیر vectorized_caption اضافه مى كنيم. تابع numericalize از لغتنامهٔ self.vocab استفاده مي كند. نماد پايان EOS را هم به بردارسازي كپشن اضافه می کنیم. تصویر img و کپشن بردارسازی شده یا vectorized_caption را به نوان تنسورها با استفاده از تابع torch. tensor بازمی گردانیم. در نهایت، این کد یک دیتاست شخصی سازی شده را برای استفاده در آموزش و ارزیابی مدلهای یادگیری عمیق ایجاد می کند. این دیتاست شامل تصاویر و کیشنهای متنی برای تصاویر است. همچنین، دیتاست شامل توابعی است که لغتنامه را براساس کیشنها ساخته و تصاویر و کیشنها را به صورت مناسب برای استفاده در مدلهای عمیق پیش پردازش مى كند. دستورات در برنامهٔ ۲ آورده شده اند.

Program 2: Creating a Custom Dataset

```
import pandas as pd
2 from torch.utils.data import Dataset
3 from PIL import Image
4 import torch
6 class FlickrDataset(Dataset):
   def __init__(self, root_dir, caption_file, transform=None, frequency_threshold=5, data_type='
     # Read the caption file into a DataFrame
     df = pd.read_csv(caption_file)
     # Split the dataset into train and test based on the data_type
     if data_type == 'train':
       # Select the first 90% of the DataFrame for training
       self.df = df.iloc[:int(0.9 * len(df))]
     elif data_type == 'test':
       # Select the remaining 10% of the DataFrame for testing
       test_start_index = int(0.9 * len(df))
       self.df = df.iloc[test_start_index:].reset_index(drop=True)
       # If data_type is neither 'train' nor 'test', do nothing
       pass
```

م عکس وصيف عکس ^{QARAS}

```
self.root_dir = root_dir
      self.transform = transform
      # Store the image names and captions
      self.imgs = self.df['image']
      self.captions = self.df['caption']
      # Initialize the vocabulary and build the vocabulary
      self.vocab = Vocabulary(frequency_threshold)
      self.vocab.build_vocab(self.captions.tolist())
    def __len__(self):
      # Return the length of the dataset
      return len(self.df)
    def __getitem__(self, index):
      # Get the caption and image information for the given index
     caption = self.captions[index]
     image_name = self.imgs[index]
     image_path = self.root_dir + '/' + image_name
     img = Image.open(image_path).convert("RGB")
      if self.transform is not None:
        # Apply the specified transformations to the image
       img = self.transform(img)
      # Convert the caption to a vectorized form
      vectorized_caption = [self.vocab.stoi["<SOS>"]]
51
      vectorized_caption += self.vocab.numericalize(caption)
      vectorized_caption.append(self.vocab.stoi["<EOS>"])
      # Return the image and its vectorized caption as tensors
      return img, torch.tensor(vectorized_caption)
```

سپس، دستوراتی برای تعریف توابع و کلاسهای مربوط به نمایش و پیش پردازش تصاویر و ایجاد نمونههای دیتاست می نویسیم. ایتدا کتابخانهٔ matplotlib را برای نمایش تصاویر وارد می کنیم. نام پای را هم برای برای کار با آرایهها و محاسبات عددی وارد می کنیم. تابع show_image را برای نمایش تصویر با استفاده از یک تانسور ورودی استفاده می کنیم. تنسور تصویر را به یک آرایهٔ نام پای تبدیل کرده و کارهایی از قبیل تغییر ابعاد آن به شکل (ارتفاع، عرض، کانال) انجام می دهیم. با استفاده از تابع matplotlib نام پای تبدیل کرده و کارهایی از قبیل تغییر ابعاد آن به شکل (ارتفاع، عرض، کانال) انجام می دهیم. با استفاده از تابع matplotlib تصویر را نمایش می دهیم. بخشی را برای نمایش کپشن و عنوان تصویر در نظر می گیریم و یک توقف هم برای به روزرسانی. تبدیل هایی که روی تصاویر باید اعمال شوند را تعریف می کنیم. برای این منظور از تابع Compose که ترکیبی از تبدیل هاست استفاده می کنیم. ابعاد تصویر را به ۲۲۴ در ۲۲۴ تبدیل می کنیم و تصاویر را به تنسور تلدیل می کنیم. دستورات در برنامهٔ ۴ آورده استفاده از کلاس مخصوصی که قبلاً تعریف کردیم، نمونههای آموزش و ارزیابی را معین می کنیم. دستورات در برنامهٔ ۴ آورده شده اند.

وصيف عكس توصيف عكس @ARASEME

Program 3: Make & Preprocess Dataset and Visualizing Data Samples

```
import matplotlib.pyplot as plt
2 import numpy as np
4 def show_image(image_tensor, title=None):
      Display an image represented as a tensor.
      Args:
          image_tensor (torch.Tensor): The input image tensor.
          title (str, optional): The title of the image. Defaults to None.
      # Convert the image tensor to a NumPy array and change the dimensions
      image_np = image_tensor.numpy().transpose((1, 2, 0))
      # Display the image using matplotlib
      plt.imshow(image_np)
      # Set the title of the image if provided
      if title is not None:
          plt.title(title)
      # Pause a bit to allow the plot to be updated
      plt.pause(0.01)
25 # Import the necessary libraries
26 import torchvision.transforms as T
28 # Define the transform to be applied to the images
29 transforms = T.Compose([
      T.Resize((224, 224)), # Resize the images to (224, 224)
      T.ToTensor() # Convert the images to tensors
32 1)
34 # Create instances of the FlickrDataset class for training and testing
35 train_dataset = FlickrDataset(
      root_dir="/content/Dataset/Images", # Path to the root directory of the images
      caption_file="captions.txt", # Path to the captions file
      transform=transforms, # Apply the defined transforms to the images
      data_type='train' # Specify the data type as 'train'
40 )
41
42 test_dataset = FlickrDataset(
      root_dir="/content/Dataset/Images", # Path to the root directory of the images
      caption_file="captions.txt", # Path to the captions file
```

```
transform=transforms, # Apply the defined transforms to the images

data_type='test' # Specify the data type as 'test'

7
```

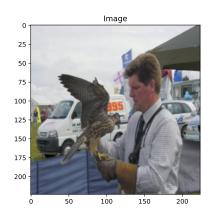
در ادامه دستوراتی را برای به نمایش گذاشتن تصویر و نمایش متن متناظر با آن تصویر مینویسیم. توضیح هر بخش کد ابتدا تصویر و شرحهای متناظر با تصویر صدم در دادههای آموزش از دیتاست train_dataset را بازیابی میکنیم. کد تصویر و کپشنهای متناظر کپشنهای متناظر می کنیم و میدهیم. در ادامه تصویر و عنوان آن را نمایش داده و کپشنهای متناظر با تصویر را به صورت توکنی چاپ می کنیم و سپس کد توکنهای موجود را به کلمات متناظر تبدیل می کنیم. دستورات مربوطه به این قسمت و نتایج مربوط به دو نمونه تصویر خروجی نشان داده شده در ؟؟ در شکل ۱ آورده شده است.

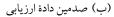
Program 4: Visualizing Data Samples

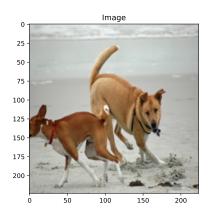
```
import matplotlib.backends.backend_pdf as pdf_backend
3 # Retrieve the image and captions for the 100th data point from the training dataset
4 image, captions = train_dataset[100]
6 # Display the image using the show_image function
7 show_image(image, "Image")
9 # Print the tokenized captions
print("Tokenized Captions:", captions)
12 # Convert the tokenized captions to their corresponding words
use | words = [train_dataset.vocab.itos[token] for token in captions.tolist()]
15 # Print the sentence
print("Sentence:")
17 print (words)
19 # Save the image as a PDF
20 pdf = pdf_backend.PdfPages('image100tr.pdf')
plt.savefig(pdf, format='pdf')
22 pdf.close()
24 Train:
25 Tokenized Captions: tensor([ 1, 48, 318, 1895, 10, 29, 15, 20, 49,
     9.
           50, 2])
28 ['<SOS>', 'two', 'different', 'breeds', 'of', 'brown', 'and', 'white', 'dogs', 'play', 'on', 'the
     ', 'beach', '<EOS>']
30 Test:
I Tokenized Captions: tensor([ 1, 12, 15, 8, 40, 3, 33, 48, 4, 3, 6, 16, 504,
      21)
```

توصيف عكس ﴿ اللَّهُ ال

```
32 Sentence:
33 ['<SOS>', 'the', 'man', 'is', 'wearing', '<UNK>', 'and', 'holding', 'a', '<UNK>', 'in', 'his', '
hands', '<EOS>']
```







(آ) صدمین دادهٔ آموزش

شکل ۱: نمایش نمونه دادهٔ آموزش و ارزیابی.

در ادامه، یک کلاس به نام Apppadd پیادهسازی می کنیم که به عنوان تابع تجمیع در دیتالودر استفاده می شود تا پدینگ (Padding) را برای جملات در هر دسته بندی اعمال کند. متد init این کلاس، دو پارامتر pad_idx و pad_idx را بریافت می کند. pad_idx مشخص می کند که آیا بعد دسته ها باید دریافت می کند. که آیا بعد دسته ها باید در اولین بُعد باشد یا نه. در ادامه پدینگ لازم روی جملات در دسته اِعمال می شود. ابتدا تصاویر را از دسته استخراج می کنیم و با استفاده از و با استفاده از unsqueeze با سیس تمام جملات را از دسته استخراج می کنیم و با استفاده از با استفاده از پدینگ توکن مشخص شده (pad_idx) پد می کنیم. در نهایت، تصاویر و جملات پدشده را به عنوان خروجی بازمی گردانیم. با استفاده از دستورات آورده شده در برنامهٔ ۵ می توانیم تصاویر و جملات را با پدینگ مناسب برای آموزش مدل به صورت دسته بندی شده (دسته ای) دریافت کنیم.

Program 5: Implement Padding for Sentences in Each Batch

```
class Apppadd:

"""

Collate function to apply padding to captions with DataLoader.

"""

def __init__(self, pad_idx, batch_first=False):

"""

Initialize the CapsCollate class.

Args:

pad_idx (int): The index of the padding token.

batch_first (bool): Whether to return the batch dimension as the first dimension.

"""
```

توصيف عكس محمده عكس هما المحمدة المحمد

```
self.pad_idx = pad_idx
          self.batch_first = batch_first
      def __call__(self, batch):
          Apply padding to the captions in the batch.
          Args:
              batch (list): List of tuples containing image and caption pairs.
          Returns:
              torch. Tensor: Batch of images.
              torch. Tensor: Batch of padded captions.
          0.00
          # Extract images from the batch and unsqueeze them
          imgs = [item[0].unsqueeze(0) for item in batch]
          imgs = torch.cat(imgs, dim=0)
          # Extract captions from the batch
          captions = [item[1] for item in batch]
          # Pad the captions using pad_sequence
          captions = pad_sequence(captions, batch_first=self.batch_first, padding_value=self.
      pad_idx)
37
          return imgs, captions
```

در ادامه دو متغیر BATCH_SIZE و NUM_WORKERS را تعریف می کنیم که به ترتیب برای تعیین اندازهٔ دسته و تعداد فرآیندهای WORKER در دیتالودر استفاده می شوند. به عنوان توضیحی از برخی متغیرهای داخل دیتالودر؛ از num_workers برای تعیین WORKER در دیتالودر استفاده می شود. به این صورت که، زمانی که داده ها تعداد پردازشهای موازی که برای بارگذاری دادهها به کار گرفته می شود با توجه برای آموزش شبکه باید از دیسک بارگذاری شوند، تعداد ورکرهایی که برای بارگذاری دادهها باید به کار گرفته می شود با توجه به مواردی قابل تعیین است: تعداد ورکرها باید کمتر از تعداد کل هستههای پردازنده سیستم باشد تا برنامه به صورت صحیح و با عملکرد بهینه اجرا شود. تعداد ورکرها باید به اندازهٔ کافی بزرگ باشد تا بتوانند با توجه به مواردی مانند حجم داده و سرعت بارگذاری آنها، دادهها را به صورت موازی بارگذاری و تحویل دهند. بنابراین، با افزایش این پارامتر و متغیر، تعداد ورکرها بیشتر پردازنده و میزان منابع دیگر مانند حافظهی سیستمی را در نظر گرفته شود. سپس، متغیر pad_idx تعریف شده است که مقدار پردازنده و میزان منابع دیگر مانند حافظهی سیستمی را در نظر گرفته شود. سپس، متغیر train_dataset بن توکن به عنوان نشانهای پردازید و میزان نمابه در مدل استفاده می شود. سپس، با استفاده از تابع دیتالودر، یک شئ دیتالودر برای بارگیری مجموعه داده آموزشی (train_dataset vocab) است. این توکن به عنوان نشانهای شود، در اینجا مجموعه دادهٔ آموزشی (train_dataset) است. که کدام مجموعه دادهٔ آموزشی (train_dataset) است که تعداد فرآیندهای ورکر برای بارگیری داده را تعیین میکند. که کمدام مجموعه دادهٔ آموزشی (کر برای بارگیری داده را تعیین میکند. که کمدام محموعه دادهٔ آموزشی (کمدر برای بارگیری دوره به صورت تصادفی مخلوط شوند یا خیر. collate_fi تابع دیتالودر، یک پسخص میکند که کمدام محموعه دادهٔ آموزشی در اینجا مجموعه دادهٔ آموزشی و خور برای بارگیری داده را تعیین میکند. در اینجا محموعه دادهٔ آموزشی در اینجا مخموعه دادهٔ آموزشی در اینجا مخموعه دادهٔ آموزشی در اینجا محموعه دادهٔ آموزشی در اینجا مخموعه دادهٔ آموزشی در اینجا مخموعه دادهٔ آموزشی در اینجا مخلوط شوند یا خیر. اینجا داده دا اینجا در اینجا مخلوط شوند یا خیر. در اینجا در ا

قوصيف عكس هجمه عكس هجمه المستقلم المست

از کلاس Apppadd استفاده شده است. در نهایت، با استفاده از این شئ دیتالودر می توان داده های آموزشی را به صورت دسته ای با سایز مشخص، با پدینگ کپشن، تصادفی شده و با تعداد فرآیندهای ورکر مناسب بارگیری کرد. دستورات مربوطه به شرح زیر است.

```
import multiprocessing
3 NUM_WORKERS = multiprocessing.cpu_count()
4 print("Maximum number of workers:", NUM_WORKERS)
6 BATCH_SIZE = 4
7 \text{ NUM_WORKERS} = 2
9 # Token to represent the padding
pad_idx = train_dataset.vocab.stoi["<PAD>"]
12 # Create a DataLoader object to load the training dataset
13 data_loader = DataLoader(
      dataset=train_dataset,
                                           # Specify the dataset to load
      batch size=BATCH SIZE,
                                           # Set the batch size
     num_workers=NUM_WORKERS,
                                           # Set the number of worker processes for data loading
                                           # Shuffle the data for each epoch
     collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True) # Specify the collate function for
      padding captions
```

در ادامه دستوراتی برای نمایش تصاویر و کپشنهای یک دستهٔ تصادفی می نویسیم. ابتدا با استفاده از تابع iterator از دیتالودر به نام dataiter ایجاد می کنیم. این iterator قادر به بازیابی دستههای متوالی از دادههاست. با استفاده از استفاده از ابع batch بازیابی می شود و در متغیر batch قرار می گیرد. سپس، با استفاده از ستفاده از captions از دادهها از dataiter بازیابی می شود و در متغیرهای قرار می گیرد. سپس، با استفاده از captions بازیابی می شود و در متغیرهای مربوط به خود (متغیرهای refective images) قرار می دهیم. متغیر aمیدهیم. متغیرهای و کپشنها را به ترتیب در متغیرهای و BATCH_SIZE و تعداد تصاویر موجود در images می دهیم. می دهیم. میشود. سپس، با یک حلقهٔ for روی اندیسهای تا Septime batch و تعداد تصاویر موجود در دسته را نمایش می دهیم. برای هر تصویر و توضیحات مربوطه را از متغیرهای images و می کنیم. اندیس توکن EOS را پیدا کرده و کپشن را از کپشن عددی شده با استفاده از train_dataset.vocab.itos استخراج می کنیم. اندیس توکن EOS را پیدا کرده و کپشن را برش می دهیم تا کپشن تکمیلی حذف شود. لیست کلمات برچسب کپشن را به یک رشته تبدیل می کنیم و یک شکل جدید برای نمایش تصویر ایجاد می کنیم. سپس، تصویر را همراه با برچسب کپشن در شکل نمایش می دهیم و عرض تصویر را محاسبه کرده و ارتفاع جعبه توضیح را تعیین می کنیم. یک جعبهٔ رنگی با برچسب کپشن را اضافه می کنیم. شکل را به عنوان یک فایل پدری دی دی اف ذخیره می کنیم و آن را نمایش می دهیم. دستورات مربوطه به شرح زیر و نمونه ای از اجرا در شکل ۲ نمایش داده شده

```
import matplotlib.pyplot as plt
import numpy as np

# Generating the iterator from the dataloader
dataiter = iter(data_loader)
```

توصيف عكس هجم وصيف عكس هجم المجمع المجمع

```
7 # Getting the next batch
8 batch = next(dataiter)
10 # Unpacking the batch
images, captions = batch
13 # Determine the effective batch size
14 effective_batch_size = min(BATCH_SIZE, len(images))
16 # Showing information of each image in the batch
for i in range(effective_batch_size):
      img, cap = images[i], captions[i]
      # Extracting the caption label from the numericalized caption
      caption_label = [train_dataset.vocab.itos[token] for token in cap.tolist()]
      # Finding the index of '<EOS>' token to truncate the caption
      eos_index = caption_label.index('<EOS>')
24
      caption_label = caption_label[1:eos_index]
25
      # Joining the caption label words into a single string
27
      caption_label = ' '.join(caption_label)
28
      # Create a new figure
30
      fig = plt.figure()
31
      # Displaying the image with the caption label
33
      plt.imshow(np.transpose(img.numpy(), (1, 2, 0)))
      plt.axis('off')
      # Calculating the width of the image
      img_width = img.shape[1]
38
      # Calculating the height of the caption box
      caption_height = int(img_width / 50)
41
      # Adding a colored box with the caption label
43
      plt.text(0, -10, caption_label, bbox=dict(facecolor='white', edgecolor='black', boxstyle='
      round'),
               fontsize=8, color='black', ha='left', va='top')
45
46
      # Save the plot as a PDF file
      plt.savefig(f'captionedimage{i}.pdf', format='pdf')
48
```

۱۵ ® معکم توصیف عکس

Display the plot plt.show()





a black dog is chasing a tan and white dog on a sunny day



(آ) دادهٔ اول



(د) دادهٔ چهارم



(ج) دادهٔ سوم

شکل ۲: نمونهای از تصویر /کیشن یک دسته از دادهها.

حال به سراغ پیادهسازی مدل میرویم و قبل از هرچه تکلیف دادهها را با توجه به تمامی کلاسها و توابعی که تاکنون نوشتهایم مشخص می کنیم. برای این منظور دستوراتی می نویسیم که در آن متغیر pad_idx نشان می دهد که در نوشتار عددی کپشنها، علامت PAD با چه عددی نشان داده شود. متغیر transforms را هم برای ایجاد یک ترکیب از تبدیلات که بر تصاویر باید اعمال شوند تعریف میکنیم. این تحویلها شامل تغییر اندازهٔ پیکسلب تصاویر، بریدن تصادفی تصاویر، تبدیل تصاویر به تنسور، و نر مال سازی تصاویر با توجه به اعداد موردنیاز برای رزنت پیش آموزش دیده هستند. سپس، یک نمونه از کلاس FlickrDataset با استفاده از پارامترهای مناسب ایجاد میکنیم. این نمونه برای آموزش استفاده می شود و ترکیب تبدیل های مختلف روی تصاویر اعمال می شود. در ادامه، یک نمونه از کلاس دیتالودر برای مجموعه آموزش ایجاد می شود. این نمونه از دیتالودر تعداد تصاویر در هر دسته را برابر با موارد از پیش مشخص شده تعیین می کند. همچنین، داده ها در هر دورهٔ آموزش قبل از اعمال دسته بندی، مخلوط می شوند تا ترتیب داده ها تاثیری بر فر آیند آموزش نداشته باشد. همچنین، تابع collate_fn برای اعمال پرکردن (padding) به کپشنها استفاده می شود. دستورات به شرح زیر است: المجاهدة توصيف عكس GARASIM

```
BATCH_SIZE = 256
2 NUM_WORKERS = 2
3 pad_idx = train_dataset.vocab.stoi["<PAD>"]
5 # Define the transformations to be applied, including resizing, random cropping,
6 # converting to tensor, and normalization using ResNet statistics
7 transforms = T.Compose([
      T.Resize(256),
     T.RandomCrop(224),
     T. ToTensor(),
      T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225))
12 1)
14 # Create an instance of the FlickrDataset for training
15 train_dataset = FlickrDataset(
     root_dir="/content/Dataset/Images",
     caption_file="captions.txt",
     transform=transforms,
     data_type='train'
20 )
22 # Create a DataLoader for the training dataset
23 train_loader = DataLoader(
     dataset=train_dataset,
    batch_size=BATCH_SIZE,
    num_workers=NUM_WORKERS,
      shuffle=True,
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
29 )
```

در قسمت تعریف مدل یک سری کلاسها و مدلها را تعریف می کنیم که در ترکیب با یکدیگر برای تولید کپشن برای تصاویر استفاده می شوند. سه کلاس LSTM ،ResNet تعریف شدهاند که در ترکیب با هم برای عملیات توصیف تصاویر استفاده می شوند. هر یک از این کلاس را تعیین می کنند.

● کلاس ResNet یک مدل از شبکه عصبی ResNet را تعریف می کند. این مدل دارای یک لایه ResNet است که از تصاویر ویژگیهای آنها را استخراج می کند. لایهٔ Embedding خروجی لایهٔ تماماً متصل آخر مدل ResNet-18 را به یک فضای تعبیه شده با اندازهٔ مشخص تبدیل می کند. در این کلاس، می توانیم تصمیم بگیریم که آیا بخشهای مختلف مدل ResNet قابل آموزش باشند یا خیر. در متد اولیهٔ این کلاس، مدل ResNet-18 پیش آموزش دیده بارگذاری می شود. سپس پارامتر train_resnet که نشان می دهد آیا باید مدل قابل آموزش باشد یا خیر، استفاده می شود تا تصمیم گیری شود که آیا پارامترهای مدل قابل آموزش باشند یا خیر. سپس، لایههای مورد نیاز از مدل تا لایهٔ تماماً متصل آخر با استفاده از [1-:]((Internet.children استخراج می شوند و به عنوان لایهٔ تبدیل بردار ویژگی ها به بردار تعبیه شده استفاده می شود. در نهایت، لایه های relu (تابع فعال سازی ReLU) و dropout (لایه tropout) هم در کار وارد می شوند.

توصيف عكس مالا من المالات الما

● کلاس LSTM یک مدل شبکه عصبی بازگشتی LSTM را تعریف میکند. این مدل هم دارای یک لایهٔ LSTM ویژگی است که کلمات کپشن را به یک فضای تعبیه شده تبدیل میکند. سپس، با استفاده از لایهٔ LSTM، بردارهای ویژگی تصویر و بردارهای تعبیه شدهٔ کلمات کپشن با یکدیگر ترکیب می شوند و خروجی را تولید میکنند. این خروجی را می توان برای پیش بینی کلمات کپشن بعدی استفاده کرد. در متد اولیهٔ این کلاس، لایه تعبیهٔ کلمات (embedding) با اندازهٔ ورودی embed_size و اندازهٔ وضعیت مخفی با اندازهٔ ورودی embed_size و اندازهٔ وضعیت مخفی با اندازهٔ می شود. در انتها، لایهٔ خطی نیز برای پیش بینی کلمهٔ بعدی استفاده می شود. در انتها، لایهٔ کتوریف می شود.

● کلاس CNNtoRNN یک مدل ترکیبی از ResNet و ResNet را تعریف می کند. این مدل شامل دو بخش اصلی است: بخش انکودر و بخش دیکودر. بخش انکودر شامل مدل رزنت می شود که وظیفهٔ استخراج ویژگی های تصویر را دارد. بخش دیکودر مدل بازگشتی را شامل می شود که با استفاده از ویژگی های استخراج شده توسط انکودر و کلمات کپشن، خروجی های پیش بینی را تولید می کند. در متد اولیهٔ این کلاس، یک نمونه از مدل رزنت به عنوان انکودر و یک نمونه از مدل بازگشتی LSTM به عنوان دیکودر ایجاد می شود. در متد forward این کلاس، تصاویر از مدل انکودر گذر داده شده و و ویژگی های تصویر استخراج می شوند. سپس ویژگی ها و کپشن ها از مدل دیکودر گذر داده شده و خروجی ها محاسبه می شوند.

در مجموع، توابع init هر کلاس وظیفهٔ مقداردهی اولیه پارامترها و تعریف لایهها و ماژولهای مورد نیاز مدل را دارند. توابع forward هر کلاس وظیفهٔ انجام محاسبات لازم برای انتقال به جلوی مدل را دارند. در کلاس رزنت، ویژگیهای استخراجشده توسط مدل به فضای تعبیهٔ خروجی مدل تبدیل میشوند. در کلاس بازگشتی هم، بردارهای تعبیهٔ کلمات و ویژگیهای تصویر با یکدیگر ترکیب کرده و خروجی را تولید میکنند. در کلاس CNNtoRNN، تصاویر از طریق بخش انکودر یا مدل رزنت عبور میکنند و ویژگیهای استخراجشده به بخش دیکودر یا مدل بازگشتی منتقل میشوند تا خروجیهای پیشبینی را تولید کنند. در ادامه، کلاس بازگشتی دارای تابع generate_caption است که با استفاده از تصاویر ویژگی و نمونههای مخفی آغازین، شروع به تولید کپشنهای متنی میکند. تا زمانی که علامت EOS که نشان دهندهٔ پایان کپشن است تولید نشود یا تعداد کلمات کپشن تولیدشده به حداکثر تعداد کلمات مجاز نرسد، کلمات پیشبینی شده اضافه و به عنوان ورودی به مدل داده می شود. در نهایت، نمایش رشته ای کلمات کپشن تولید شده به عنوان خروجی تابع ارائه می شود. با ارائهٔ این توضیحات مدل در برنامهٔ ۶ تعریف شده است.

Program 6: Model

```
#ResNet Model
class ResNet(nn.Module):

def __init__(self, embed_size, train_resnet=False):

"""

Initialize the ResNet model with a specified embedding size.

Args:

embed_size (int): Size of the embedding output.

train_resnet (bool): Whether to train the ResNet backbone or not.

"""

super(ResNet, self).__init__()

# Load the pretrained ResNet-18 model
```

قوصيف عكس ^{©ARAS®} توصيف عكس

```
resnet = models.resnet18(pretrained=True)
          # Set the requires_grad flag of the ResNet parameters
          # to control whether they are trainable or not
          for param in resnet.parameters():
              param.requires_grad_(train_resnet)
          # Extract the modules of the ResNet model up to the last fully connected layer
          modules = list(resnet.children())[:-1]
          # Create the ResNet backbone with the extracted modules
          self.resnet = nn.Sequential(*modules)
          # Create the embedding layer
          self.embed = nn.Linear(resnet.fc.in_features, embed_size)
          # Activation function and dropout layer
          self.relu = nn.ReLU()
          self.dropout = nn.Dropout(0.5)
32
      def forward(self, images):
          ....
          Forward pass of the ResNet model.
          Args:
              images (tensor): Input images tensor.
          Returns:
41
              features (tensor): Embedded features tensor.
          features = self.dropout(self.relu(self.resnet(images)))
          features = features.view(features.size(0), -1)
          features = self.embed(features)
          return features
49 # LSTM Network
50 class LSTM(nn.Module):
      def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5):
51
          Initialize the LSTM model with specified sizes and parameters.
          Args:
              embed_size (int): Size of the word embedding.
              hidden_size (int): Size of the hidden state of the LSTM.
              vocab_size (int): Size of the vocabulary.
```

توصيف عكس هم في المحمد المحمد

```
num_layers (int): Number of layers in the LSTM (default: 1).
               drop_prob (float): Dropout probability (default: 0.5).
          super().__init__()
62
          # Word embedding layer
          self.embedding = nn.Embedding(vocab_size, embed_size)
          # LSTM layer
          self.lstm = nn.LSTM(embed_size, hidden_size, num_layers=num_layers, batch_first=True)
          # Linear layer for prediction
          self.linear = nn.Linear(hidden_size, vocab_size)
71
          # Dropout layer
          self.dropout = nn.Dropout(drop_prob)
      def forward(self, features, captions):
          Forward pass of the LSTM model.
80
          Args:
               features (tensor): Image features tensor.
81
               captions (tensor): Captions tensor.
83
          Returns:
               x (tensor): Output tensor.
          # Vectorize the caption by passing it through the embedding layer
          embeds = self.dropout(self.embedding(captions[:, :-1]))
          # Concatenate the features and captions
          x = torch.cat((features.unsqueeze(1), embeds), dim=1)
91
          # Pass through the LSTM layer
          x, _ = self.lstm(x)
94
          # Pass through the linear layer for prediction
          x = self.linear(x)
97
          return x
99
100
      def generate_caption(self, inputs, hidden=None, max_len=20, vocab=None):
101
102
          Generate captions given the image features.
```

توصيف عكس GARASHI

```
104
105
           Args:
               inputs (tensor): Input tensor of image features.
               hidden (tuple): Hidden state of the LSTM (default: None).
107
               max_len (int): Maximum length of the generated caption (default: 20).
               vocab (Vocab): Vocabulary object (default: None).
110
           Returns:
               caption (list): Generated caption as a list of words.
           # Inference part
114
           # Given the image features, generate the captions
116
           batch_size = inputs.size(0)
118
           captions = []
120
           for i in range(max_len):
               output, hidden = self.lstm(inputs, hidden)
               output = self.linear(output)
               output = output.view(batch_size, -1)
               # Select the word with the highest value
126
               predicted_word_idx = output.argmax(dim=1)
128
               # Save the generated word
129
               captions.append(predicted_word_idx.item())
132
               # End if <EOS> is detected
               if vocab.itos[predicted_word_idx.item()] == "<EOS>":
                   break
134
               # Send the generated word as the next caption
136
               inputs = self.embedding(predicted_word_idx.unsqueeze(0))
138
           # Convert the vocabulary indices to words and return the sentence
139
           return [vocab.itos[idx] for idx in captions]
142 # Utilizing the powerful fusion of ResNet and LSTM for image captioning
  class CNNtoRNN(nn.Module):
       def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5,
       train_resnet=False):
           Initialize the CNNtoRNN model with specified sizes and parameters.
146
```

توصيف عكس كالمحمدة توصيف عكس

```
Args:
148
               embed_size (int): Size of the word embedding.
149
               hidden_size (int): Size of the hidden state of the LSTM.
               vocab_size (int): Size of the vocabulary.
               num_layers (int): Number of layers in the LSTM (default: 1).
               drop_prob (float): Dropout probability (default: 0.5).
               train_resnet (bool): Whether to train the ResNet backbone or not (default: False).
           super().__init__()
156
           # Encoder (ResNet)
           self.encoder = ResNet(embed_size, train_resnet)
159
           # Decoder (LSTM)
161
           self.decoder = LSTM(embed_size, hidden_size, vocab_size, num_layers, drop_prob)
       def forward(self, images, captions):
164
           Forward pass of the CNNtoRNN model.
166
167
           Args:
               images (tensor): Input images tensor.
               captions (tensor): Captions tensor.
           Returns:
               outputs (tensor): Output tensor.
           # Pass the images through the encoder (ResNet) to get features
175
           features = self.encoder(images)
           # Pass the features and captions through the decoder (LSTM) to get outputs
178
           outputs = self.decoder(features, captions)
180
           return outputs
```

در ادامه دستوراتی را برای تعریف یک حلقهٔ آموزشی برای آموزش مدل شبکه عصبی می نویسیم و سپس نموداری از تغییرات خطا را در طول آموزش رسم می کند. ابتدا تعداد دورههای آموزشی را تعریف می کنیم. سپس، یک لیست خالی ایجاد می کنیم تا خطاهای آموزش ذخیره شوند. یک حلقه به ازای تعداد دورههای آموزشی معین ایجاد می کنیم. متغیر Tunning_loss را ایجاد و آن را با مقدار صفر مقداردهی اولیه می کنیم. این متغیر برای محاسبه میانگین خطا در هر دوره استفاده می شود. سپس، یک حلقهٔ داخلی روی دادههای آموزش ایجاد می کنیم و تصاویر و کپشنها را به دستگاه مشخص شده (مثلاً GPU) منتقل می کنیم. پس از آن، مدل را در حالت آموزش قرار می دهیم و گرادیانهای موجود در بهینه ساز را صفر می کنیم. این کار قبل از محاسبه گرادیان برای دستهٔ فعلی ضروری است. با اعمال مدل روی تصاویر و کپشنها، خروجیهای پیش بینی را تولید می کنیم. خطای دسته های فعلی را با استفاده از تابع خطا محاسبه می کنیم. و می مینیم بارامترها با استفاده از الگوریتم پس انتشار خطا محاسبه می کنیم. یکسان تغییر شکل می دهیم. در ادامه، گرادیان را برای همه پارامترها با استفاده از الگوریتم پس انتشار خطا محاسبه می کنیم.

توصيف عكس ٩٨٣٨٥٥٤ توصيف عكس

بهینهساز، گامهای بهینهسازی را برای بهروزرسانی پارامترها براساس محاسبات انجامشده در مرحله قبل انجام می دهد. خطا برای دستهٔ فعلی به running_loss اضافه می شود تا در نهایت میانگین خطا برای هر دوره محاسبه شود. بنابراین، میانگین خطا برای هر دوره محاسبه می شود، با این توضیح که idx نشان دهندهٔ شمارهٔ دستهٔ فعلی در داده های آموزش است. در نهایت مقادیر شمارهٔ دوره، خطا و... چاپ می شوند. میانگین خطای آموزش را به لیست train_loss اضافه می کنیم تا بعداً برای رسم نمودار استفاده شود. مدل را پس از هر دوره در فایلی با نام Model.pth ذخیره می کنیم و نمودارهای مورد نیاز را نمایش می دهیم. مجموعهٔ این دستورات در زیر آورده شده اند:

```
# Set the number of training epochs
2 \text{ num\_epochs} = 100
4 # Create an empty list to store the training loss
5 train_loss = []
7 # Training loop
8 for epoch in range(num_epochs):
     running_loss = 0
     # Iterate over the training data loader
     for idx, (image, captions) in enumerate(iter(train_loader)):
          # Move the image and captions to the specified device
          image, captions = image.to(device), captions.to(device)
          # Set the model to train mode
          model.train()
          # Zero the gradients in the optimizer
          optimizer.zero_grad()
          # Feed forward
          outputs = model(image, captions)
          # Calculate the batch loss
          loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
          # Backward pass
         loss.backward()
          # Update the parameters in the optimizer
          optimizer.step()
          # Accumulate the running loss
          running_loss += loss.item()
      # Calculate the average loss for the epoch
      average_loss = running_loss / (idx + 1)
```

توصيف عكس GARASIM

```
# Print the epoch number and the average loss
40
      print(f'Epoch: {epoch+1} - Train Loss: {average_loss}')
41
42
43
      # Append the average loss to the train_loss list
      train_loss.append(average_loss)
      # Save the model after each epoch
      torch.save(model.cpu().state_dict(), 'Model.pth')
47
      # Move the model back to the specified device
      model.cuda()
52 import matplotlib.pyplot as plt
54 # Set the figure size and dpi for better quality
55 plt.figure(figsize=(8, 6), dpi=80)
57 # Plot the training loss
58 plt.plot(train_loss, label='Training Loss')
60 # Set the plot title and axis labels
61 plt.title("Training Loss per Epoch")
62 plt.xlabel("Epoch")
63 plt.ylabel("Loss")
65 # Customize the grid and ticks
66 plt.grid(True, linestyle='--', linewidth=0.5)
67 plt.xticks(range(len(train_loss)))
68 plt.yticks()
70 # Add a legend
71 plt.legend()
73 # Save the plot as a PDF file
74 plt.savefig("loss_plot.pdf", format='pdf', bbox_inches='tight')
76 # Show the plot
77 plt.show()
```

حال برای قسمت ارزیابی مطابق توضیحات داده شده در قسمت آموزش، دیتالودر قسمت تست را هم به همان ترتیب تشکیل می دهیم و سپس مدل را روی تعدادی از دادگان تست می کنیم و نتیجه را به صورتی مناسب نمایش می دهیم. دستورات را به صورتی نوشته ایم که یک تصویر تصادفی از مجموعهٔ ارزیابی انتخاب کند و کپشنهای پیش بینی شده و واقعی را روی تصویر و در جعبه هایی با رنگهای متفاوت نشان دهد. دستورات به شرح زیر است:

```
import matplotlib.pyplot as plt
```

توصيف عكس ٩٨٣٨٥٥٤٤ توصيف عكس

```
def show_image_with_captions(image, predicted_caption, real_caption):
      """Display an image with predicted and real captions."""
      # Denormalize the image tensor
      image[0] = image[0] * 0.229
      image[1] = image[1] * 0.224
      image[2] = image[2] * 0.225
      image[0] += 0.485
10
      image[1] += 0.456
11
      image[2] += 0.406
      # Convert the image tensor to a numpy array and transpose the dimensions
14
      image = image.numpy().transpose((1, 2, 0))
15
      # Display the image
      plt.imshow(image)
18
      # Add predicted caption box
20
      plt.text(
21
          0, -20, predicted_caption, color='white', backgroundcolor='blue',
          fontsize=12, verticalalignment='top', bbox=dict(facecolor='blue', alpha=0.8, edgecolor='
      white', pad=5)
      # Add real caption box
      plt.text(
          0, -2, real_caption, color='black', backgroundcolor='green',
          fontsize=12, verticalalignment='top', bbox=dict(facecolor='green', alpha=0.8, edgecolor='
      white', pad=5)
      plt.axis('off')
      plt.tight_layout()
35 # Define the test data transformations
36 test_transforms = T.Compose([
      T.Resize((224, 224)), # Resize the images to the specified size
      T.ToTensor(), # Convert the images to tensors
      T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225)) # Normalize the image tensors
40 ])
42 # Create the test dataset
43 test_dataset = FlickrDataset(
44 root_dir="Images", # Root directory of the dataset
```

توصيف عكس GARASHI

```
caption_file="captions.txt", # File containing the captions
      transform=test_transforms, # Apply the specified transformations to the images
      frequency_threshold=1, # Frequency threshold for filtering captions
      data_type='test' # Specify the type of data (in this case, 'test')
49 )
51 # Create the test data loader
52 test_loader = DataLoader(
      dataset=test_dataset, # Use the created test dataset
      batch_size=BATCH_SIZE, # Number of samples per batch
      num_workers=NUM_WORKERS, # Number of worker threads for data loading
      shuffle=True, # Shuffle the data for each epoch
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True) # Function to collate and preprocess
      the data
58 )
60 # Get a batch of images and captions from the test loader
images, captions = next(iter(test_loader))
# Iterate over a single image and its captions
64 for i in range(1):
      # Set the model to evaluation mode
      model.eval()
      # Get a single test image and create a clone
68
      test_image = torch.clone(images)[i].unsqueeze(0)
      # Disable gradient calculation during inference
71
      with torch.no_grad():
          # Encode the test image using the model's encoder
          features = model.encoder(test_image[0:1].to(device))
          # Generate captions for the test image using the model's decoder
          predicted_captions = model.decoder.generate_caption(features.unsqueeze(0), vocab=
      train_dataset.vocab)
          # Get the ground truth caption for the test image
          cap = torch.clone(captions)[i]
          caption_label = [test_dataset.vocab.itos[token] for token in cap.tolist()]
81
          # Find the index of the end-of-sequence token '<EOS>' in the caption
          eos_index = caption_label.index('<EOS>')
          # Extract the caption tokens from the start to the '<EOS>' token
          caption_label = caption_label[1:eos_index]
```

توصيف عكس GARASIM

```
# Convert the caption tokens to a string
caption_label = ' '.join(caption_label)

# Create the predicted caption string with proper formatting
predicted_caption = "Predicted Caption: " + ' '.join(predicted_captions[1:len(
predicted_captions)-1])

# Create the actual caption string with proper formatting
real_caption = "Actual Caption: " + caption_label

# Display the image with the predicted and real captions
show_image_with_captions(test_image[0], predicted_caption, real_caption)

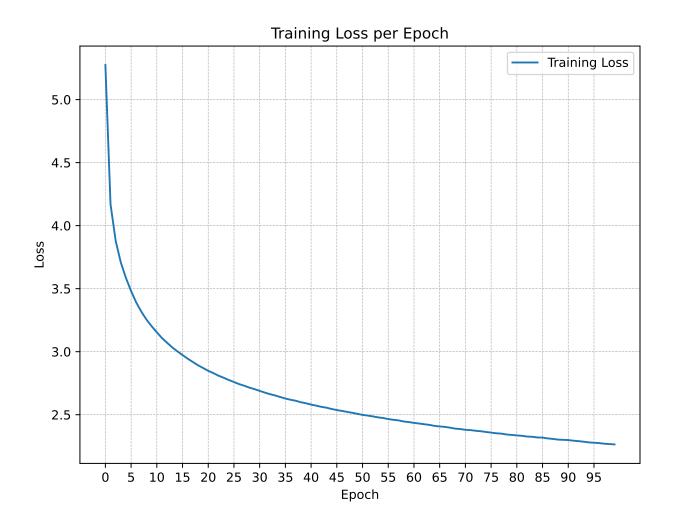
# Save the plot as a PDF file
plt.savefig("testimagecaptions.pdf", format='pdf', bbox_inches='tight')
```

با اجرای این فرآیند و دستورات نمودار خطا به صورتی است که در شکل ۳ نمایش داده شده و نتایج سه نمونهٔ تصادفی از مجموعهٔ ارزیابی به صورتی است که در شکل ۴ تا شکل ۸ نمایش داده شده است. همان طور که از نتایج مشخص است مدل هم چنان جا برای آموزش دارد اما به دلیل کمبود زمان (تمرین ها را تک نفره انجام می دهم) از ادامهٔ آموزش صرف نظر شده است. نتایج کپشن های پیش بینی شده روی داده های ارزیابی هم نشان می دهد که آموزش مدل قابل قبول بوده است و کپشن ها نزدیکی بالایی با آن چه در تصویر جریان دارد دارند. هر چند بعضی جاها لغزش های بسیار کوچکی هم دیده می شود.

```
Epoch: 1 - Train Loss: 5.2742267722016445
2 Epoch: 2 - Train Loss: 4.169356649572199
3 Epoch: 3 - Train Loss: 3.8797206595227434
4 Epoch: 4 - Train Loss: 3.7080819373364213
5 Epoch: 5 - Train Loss: 3.584198301488703
6 Epoch: 6 - Train Loss: 3.481250404477953
7 Epoch: 7 - Train Loss: 3.3901297129117527
8 Epoch: 8 - Train Loss: 3.3170354816463443
9 Epoch: 9 - Train Loss: 3.254032086659145
10 Epoch: 10 - Train Loss: 3.2015588150157797
14 Epoch: 90 - Train Loss: 2.3006849122214152
15 Epoch: 91 - Train Loss: 2.2989422974886593
16 Epoch: 92 - Train Loss: 2.293771370307549
17 Epoch: 93 - Train Loss: 2.2902511366597422
18 Epoch: 94 - Train Loss: 2.2857611779566414
19 Epoch: 95 - Train Loss: 2.2802960872650146
20 Epoch: 96 - Train Loss: 2.277562533225213
21 Epoch: 97 - Train Loss: 2.274853866417091
22 Epoch: 98 - Train Loss: 2.269936553248159
```

توصيف عكس GARASIM

```
23 Epoch: 99 - Train Loss: 2.267233173330347
24 Epoch: 100 - Train Loss: 2.26415508443659
```



شكل ٣: نمودار خطا در حالت رزنت فريز.

برای نمایش خطای تست، از آنجا که دادن دادههای تست در هنگام آموزش به نوعی تقلب محسوب میشود، دستورات زیر را نوشتیم تا بخشی از دادههای آموزش را بهعنوان دادهٔ اعتبارسنجی درنظر بگیریم و بتوانیم خطای تست را گزارش و رسم کنیم.

```
import torch.nn as nn
import torch.optim as optim
import torchvision.transforms as T
from torch.utils.data import DataLoader
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

# Set the batch size, number of workers, and pad index
BATCH_SIZE = 256
```

توصيف عكس ٩٨٣٨٥٥٤٤ توصيف عكس

Predicted Caption: a young boy is jumping into a pool

Actual Caption: a little girl with inflatable armbands is jumping into the pool



شکل ۴: نتایج پیش بینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ اول.

Predicted Caption: a black and white dog is running through a grassy field

Actual Caption: two dogs are wrestling in the park



شکل ۵: نتایج پیشبینی روی داده های ارزیابی در حالت رزنت فریز - نمونهٔ دوم.

NUM_WORKERS = 10

pad_idx = train_dataset.vocab.stoi["<PAD>"]

توصيف عكس GARASE

Predicted Caption: a young boy is playing with a white <UNK>

Actual Caption: a curly haired boy jumps up on a bed



شکل ۶: نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ سوم.

Predicted Caption: a man is standing on a beach with a surfboard in the background Actual Caption: a child is splashing in the water at the beach



شکل ۷: نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت فریز - نمونهٔ چهارم.

```
# Define the transformations to be applied, including resizing, random cropping,

# converting to tensor, and normalization using ResNet statistics

# transforms = T.Compose([
```

توصيف عكس محمد توصيف عكس

Predicted Caption: a little girl in a pink dress is sitting on a wooden bench

Actual Caption: two young girls lay on the carpet next to wooden blocks



شکل ۸: نتایج پیشبینی روی داده های ارزیابی در حالت رزنت فریز - نمونهٔ پنجم.

```
T.Resize(256),
      T.RandomCrop(224),
      T.ToTensor(),
      T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225))
21 ])
23 # Create an instance of the FlickrDataset for training
24 train_dataset = FlickrDataset(
     root_dir="/content/Dataset/Images",
     caption_file="captions.txt",
    transform=transforms,
     data_type='train'
31 # Split the dataset into training and validation sets
32 train_dataset, val_dataset = train_test_split(train_dataset, test_size=0.1, random_state=42)
34 # Create a DataLoader for the training dataset
35 train_loader = DataLoader(
     dataset=train_dataset,
     batch_size=BATCH_SIZE,
     num_workers=NUM_WORKERS,
    shuffle=True,
```

TO Eagle 20 To Garasing To Ga

```
collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
41 )
43 # Create a DataLoader for the validation dataset
44 val_loader = DataLoader(
      dataset=val_dataset,
      batch_size=BATCH_SIZE,
      num_workers=NUM_WORKERS,
      shuffle=False,
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
50 )
52 # Set the number of training epochs
num_epochs = 20
55 # Create empty lists to store the training and validation loss
56 train_loss = []
57 val_loss = []
59 # Training loop
60 for epoch in range(num_epochs):
      running_train_loss = 0
      running_val_loss = 0
      # Training
64
      model.train()
      for idx, (image, captions) in enumerate(train_loader):
          image, captions = image.to(device), captions.to(device)
          optimizer.zero_grad()
          outputs = model(image, captions)
          loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
          loss.backward()
71
          optimizer.step()
          running_train_loss += loss.item()
      average_train_loss = running_train_loss / len(train_loader)
      train_loss.append(average_train_loss)
75
      # Validation
      model.eval()
      with torch.no_grad():
          for idx, (image, captions) in enumerate(val_loader):
81
              image, captions = image.to(device), captions.to(device)
              outputs = model(image, captions)
              loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
              running_val_loss += loss.item()
```

توصيف عكس GARASE

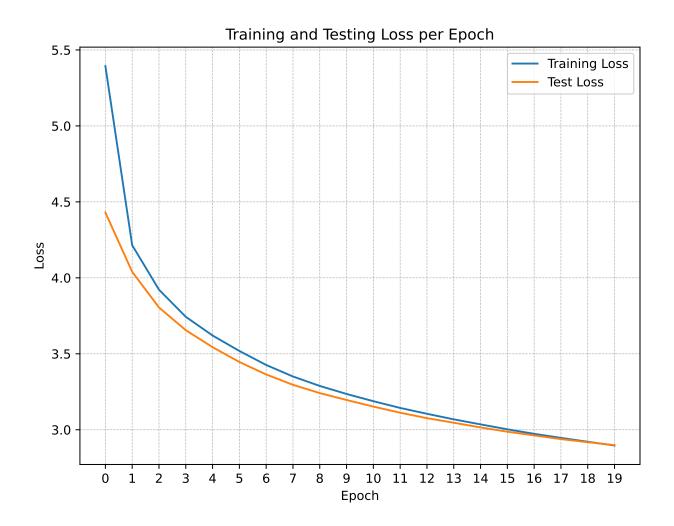
```
average_val_loss = running_val_loss / len(val_loader)
      val_loss.append(average_val_loss)
      print(f'Epoch: {epoch+1} - Train Loss: {average_train_loss:.4f} - Val Loss: {average_val_loss
      :.4f}')
      # Save the model after each epoch
      torch.save(model.cpu().state_dict(), 'ModelwithVal.pth')
      model.cuda()
94 # Plotting the training and validation loss
95 plt.figure(figsize=(8, 6), dpi=80)
96 plt.plot(train_loss, label='Training Loss')
97 plt.plot(val_loss, label='Validation Loss')
98 plt.title("Training and Validation Loss per Epoch")
99 plt.xlabel("Epoch")
100 plt.ylabel("Loss")
plt.grid(True, linestyle='--', linewidth=0.5)
plt.xticks(range(num_epochs))
plt.legend()
plt.savefig("loss_plot.pdf", format='pdf', bbox_inches='tight')
```

نتایج به صورتی است که در شکل ۱۶ آمده است. اینباز برای عدم تکرار فرآیند قبل و صرفاً برای نمایش نتایج تنها برای ۲۰ دوره آموزش دادیم. نتیجه روی سه دادهٔ تست تصادفی هم در مورد ارزیابی قرار گرفته که در دفتر چه کد گوگل کولب قابل مشاهده است و این جا برای جلوگیری از ازد حام بیش تر نیاورده ایم.

```
Epoch: 1 - Train Loss: 5.3937 - Val Loss: 4.4301
2 Epoch: 2 - Train Loss: 4.2150 - Val Loss: 4.0397
3 Epoch: 3 - Train Loss: 3.9220 - Val Loss: 3.8049
4 Epoch: 4 - Train Loss: 3.7434 - Val Loss: 3.6553
5 Epoch: 5 - Train Loss: 3.6204 - Val Loss: 3.5436
6 Epoch: 6 - Train Loss: 3.5186 - Val Loss: 3.4460
7 Epoch: 7 - Train Loss: 3.4264 - Val Loss: 3.3641
8 Epoch: 8 - Train Loss: 3.3506 - Val Loss: 3.2957
9 Epoch: 9 - Train Loss: 3.2884 - Val Loss: 3.2414
10 Epoch: 10 - Train Loss: 3.2355 - Val Loss: 3.1959
H Epoch: 11 - Train Loss: 3.1880 - Val Loss: 3.1525
12 Epoch: 12 - Train Loss: 3.1436 - Val Loss: 3.1122
13 Epoch: 13 - Train Loss: 3.1050 - Val Loss: 3.0762
14 Epoch: 14 - Train Loss: 3.0685 - Val Loss: 3.0463
15 Epoch: 15 - Train Loss: 3.0353 - Val Loss: 3.0153
16 Epoch: 16 - Train Loss: 3.0030 - Val Loss: 2.9872
17 Epoch: 17 - Train Loss: 2.9733 - Val Loss: 2.9627
18 Epoch: 18 - Train Loss: 2.9463 - Val Loss: 2.9385
19 Epoch: 19 - Train Loss: 2.9210 - Val Loss: 2.9180
```

وصيف عكس ١٩٩٤ وصيف عكس ١٩٣٤

```
20 Epoch: 20 - Train Loss: 2.8963 - Val Loss: 2.8987
```



شکل ۹: نمودار خطای آموزش و تست در حالت رزنت فریز.

۲.۱ پاسخ قسمت ۲

برای پاسخ به این قسمت از تکرار مجدد موارد گفته شده در پاسخ قسمت ۱ صرف نظر می کنیم. تنها قسمتی که نیاز به تغییر دارد قسمت تعریف مدل است. در قسمت قبل قابلیت آموزش لایه های رزنت را خاموش نگه داشته بودیم و در این قسمت، مطابق با خواستهٔ سوال آن را فعال می کنیم. دستور تغییریافته و مجموعهٔ دستورات تغییریافتهٔ مربوط به این بخش در زیر آورده شده است (train_resnet=True):

```
#ResNet Model
class ResNet(nn.Module):
def __init__(self, embed_size, train_resnet=False):

#ResNet Model
Init__(self, embed_size, train_resnet=False):

#ResNet model with a specified embedding size.
```

توصيف عكس @ARAS توصيف عكس

```
Args:
              embed_size (int): Size of the embedding output.
              train_resnet (bool): Whether to train the ResNet backbone or not.
          super(ResNet, self).__init__()
11
          # Load the pretrained ResNet-18 model
          resnet = models.resnet18(pretrained=True)
14
          # Set the requires_grad flag of the ResNet parameters
          # to control whether they are trainable or not
          for param in resnet.parameters():
              param.requires_grad_(train_resnet)
          # Extract the modules of the ResNet model up to the last fully connected layer
          modules = list(resnet.children())[:-1]
          # Create the ResNet backbone with the extracted modules
          self.resnet = nn.Sequential(*modules)
          # Create the embedding layer
          self.embed = nn.Linear(resnet.fc.in_features, embed_size)
          # Activation function and dropout layer
          self.relu = nn.ReLU()
31
          self.dropout = nn.Dropout(0.5)
34
      def forward(self, images):
          Forward pass of the ResNet model.
          Args:
              images (tensor): Input images tensor.
          Returns:
41
              features (tensor): Embedded features tensor.
43
          features = self.dropout(self.relu(self.resnet(images)))
44
          features = features.view(features.size(0), -1)
          features = self.embed(features)
47
          return features
49 # LSTM Network
50 class LSTM(nn.Module):
```

```
def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5):
52
          Initialize the LSTM model with specified sizes and parameters.
          Args:
              embed_size (int): Size of the word embedding.
              hidden_size (int): Size of the hidden state of the LSTM.
              vocab_size (int): Size of the vocabulary.
              num_layers (int): Number of layers in the LSTM (default: 1).
              drop_prob (float): Dropout probability (default: 0.5).
61
          super().__init__()
          # Word embedding layer
          self.embedding = nn.Embedding(vocab_size, embed_size)
          # LSTM layer
67
          self.lstm = nn.LSTM(embed_size, hidden_size, num_layers=num_layers, batch_first=True)
          # Linear layer for prediction
70
          self.linear = nn.Linear(hidden_size, vocab_size)
          # Dropout layer
          self.dropout = nn.Dropout(drop_prob)
      def forward(self, features, captions):
          Forward pass of the LSTM model.
          Args:
              features (tensor): Image features tensor.
81
              captions (tensor): Captions tensor.
83
          Returns:
              x (tensor): Output tensor.
          # Vectorize the caption by passing it through the embedding layer
          embeds = self.dropout(self.embedding(captions[:, :-1]))
          # Concatenate the features and captions
          x = torch.cat((features.unsqueeze(1), embeds), dim=1)
91
92
          # Pass through the LSTM layer
          x, _ = self.lstm(x)
```

توصيف عكس @ARAS توصيف عكس

```
# Pass through the linear layer for prediction
           x = self.linear(x)
           return x
99
       def generate_caption(self, inputs, hidden=None, max_len=20, vocab=None):
101
102
           Generate captions given the image features.
104
105
           Args:
               inputs (tensor): Input tensor of image features.
106
               hidden (tuple): Hidden state of the LSTM (default: None).
107
               max_len (int): Maximum length of the generated caption (default: 20).
               vocab (Vocab): Vocabulary object (default: None).
109
110
           Returns:
               caption (list): Generated caption as a list of words.
           # Inference part
114
           # Given the image features, generate the captions
           batch_size = inputs.size(0)
117
118
           captions = []
           for i in range(max_len):
               output, hidden = self.lstm(inputs, hidden)
               output = self.linear(output)
               output = output.view(batch_size, -1)
               # Select the word with the highest value
126
               predicted_word_idx = output.argmax(dim=1)
128
               # Save the generated word
129
               captions.append(predicted_word_idx.item())
130
               # End if <EOS> is detected
               if vocab.itos[predicted_word_idx.item()] == "<EOS>":
                   break
134
               # Send the generated word as the next caption
136
               inputs = self.embedding(predicted_word_idx.unsqueeze(0))
138
           # Convert the vocabulary indices to words and return the sentence
139
           return [vocab.itos[idx] for idx in captions]
```

قوصيف عكس <u>همجمعه</u> توصيف عكس

```
142 # Utilizing the powerful fusion of ResNet and LSTM for image captioning
  class CNNtoRNN(nn.Module):
       def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5,
       train_resnet=False):
           Initialize the CNNtoRNN model with specified sizes and parameters.
146
           Args:
148
               embed_size (int): Size of the word embedding.
149
               hidden_size (int): Size of the hidden state of the LSTM.
               vocab_size (int): Size of the vocabulary.
               num_layers (int): Number of layers in the LSTM (default: 1).
               drop_prob (float): Dropout probability (default: 0.5).
               train_resnet (bool): Whether to train the ResNet backbone or not (default: False).
154
           super().__init__()
156
           # Encoder (ResNet)
158
           self.encoder = ResNet(embed_size, train_resnet)
           # Decoder (LSTM)
161
           self.decoder = LSTM(embed_size, hidden_size, vocab_size, num_layers, drop_prob)
162
       def forward(self, images, captions):
164
           Forward pass of the CNNtoRNN model.
167
           Args:
               images (tensor): Input images tensor.
               captions (tensor): Captions tensor.
170
           Returns:
               outputs (tensor): Output tensor.
174
           # Pass the images through the encoder (ResNet) to get features
           features = self.encoder(images)
           # Pass the features and captions through the decoder (LSTM) to get outputs
178
           outputs = self.decoder(features, captions)
179
180
           return outputs
182
183
```

توصيف عكس ٩٨٣٨٥٥٥ توصيف عكس

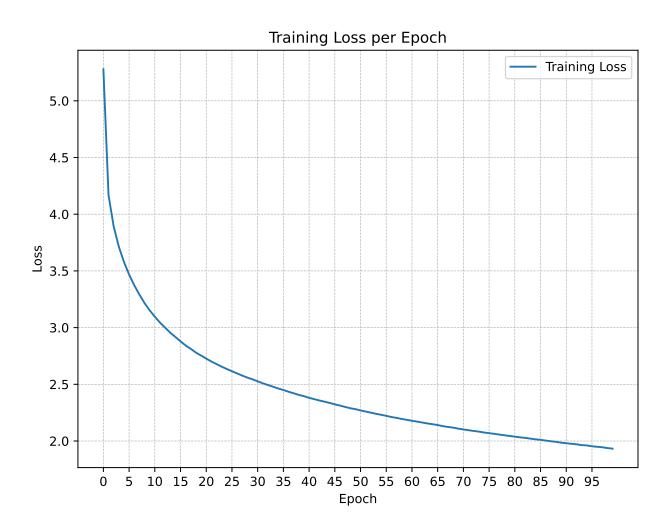
با اجرای این فرآیند و دستورات نمودار خطا به صورتی است که در شکل ۱۰ نمایش داده شده و نتایج سه نمونهٔ تصادفی از مجموعهٔ ارزیابی به صورتی است که در شکل ۱۱ تا شکل ۱۵ نمایش داده شده است. همان طور که از نتایج مشخص است مدل هم چنان جا برای آموزش دارد اما به دلیل کمبود زمان (تمرین ها را تک نفره انجام می دهم) از ادامهٔ آموزش صرف نظر شده است. می توان از () torch. cuda. amp. GradScaler هم برای افزایش سرعت و بهره وری در فرآیند آموزش استفاده کرد. به هر حال، نتایج کپشن های پیش بینی شده روی داده های ارزیابی هم نشان می دهد که آموزش مدل قابل قبول بوده است و کپشن ها نزدیکی بالایی با آن چه در تصویر جریان دارد دارند. هر چند بعضی جاها لغزش های بسیار کوچکی هم دیده می شود.

```
Epoch: 1 - Train Loss: 5.2784986095828605
2 Epoch: 2 - Train Loss: 4.1686154712330215
3 Epoch: 3 - Train Loss: 3.892264472854721
4 Epoch: 4 - Train Loss: 3.71403486245162
5 Epoch: 5 - Train Loss: 3.5789838320725447
6 Epoch: 6 - Train Loss: 3.4679061566199456
7 Epoch: 7 - Train Loss: 3.3738862034324164
8 Epoch: 8 - Train Loss: 3.2915008068084717
9 Epoch: 9 - Train Loss: 3.2161747125478892
10 Epoch: 10 - Train Loss: 3.152081139437802
14 Epoch: 90 - Train Loss: 1.985385377090294
15 Epoch: 91 - Train Loss: 1.9803814221095373
16 Epoch: 92 - Train Loss: 1.9754978493377044
17 Epoch: 93 - Train Loss: 1.9709438619080124
18 Epoch: 94 - Train Loss: 1.9638588253434721
19 Epoch: 95 - Train Loss: 1.961058676659644
20 Epoch: 96 - Train Loss: 1.9537817581550225
21 Epoch: 97 - Train Loss: 1.9489550223717322
22 Epoch: 98 - Train Loss: 1.9450463039891703
23 Epoch: 99 - Train Loss: 1.937939615516396
24 Epoch: 100 - Train Loss: 1.9327541773135846
```

برای نمایش خطای تست، از آنجا که دادن دادههای تست در هنگام آموزش به نوعی تقلب محسوب میشود، دستورات زیر را نوشتیم تا بخشی از دادههای آموزش را بهعنوان دادهٔ اعتبارسنجی درنظر بگیریم و بتوانیم خطای تست را گزارش و رسم کنیم.

```
import torch
import torch.nn as nn
import torch.optim as optim
import torchvision.transforms as T
from torch.utils.data import DataLoader
```

قوصيف عكس <u>محمدة عكس همومية المحمدة ا</u>



شكل ۱۰: نمودار خطا در حالت رزنت بدون فريز.

```
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

# Set the batch size, number of workers, and pad index

BATCH_SIZE = 256

NUM_WORKERS = 10

pad_idx = train_dataset.vocab.stoi["<PAD>"]

# Define the transformations to be applied, including resizing, random cropping,

# converting to tensor, and normalization using ResNet statistics

transforms = T.Compose([
    T.Resize(256),
    T.RandomCrop(224),
    T.ToTensor(),
    T.ToTensor(),
    T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225))

1 ])
```

توصيف عكس توصيف عكس

Predicted Caption: a brown dog is jumping up in the air to catch a ball

Actual Caption: a brown dog jumps as he looks at a soccer ball



شکل ۱۱: نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ اول.

Predicted Caption: people are walking on a city street

Actual Caption: two elderly women are walking past a younger woman on a public path



شکل ۱۲: نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ دوم.

```
# Create an instance of the FlickrDataset for training
train_dataset = FlickrDataset(
    root_dir="/content/Dataset/Images",
    caption_file="captions.txt",
    transform=transforms,
    data_type='train'
```

توصيف عكس @ARAS توصيف عكس

Predicted Caption: a young boy wearing a blue shirt is running on a blue plastic toy
Actual Caption: a young boy paying in a small pool with water



شکل ۱۳: نتایج پیش بینی روی داده های ارزیابی در حالت رزنت بدون فریز - نمونهٔ سوم.

Predicted Caption: a group of children playing in a fountain

Actual Caption: two children are in costume in a dance studio



شکل ۱۴: نتایج پیش بینی روی داده های ارزیابی در حالت رزنت بدون فریز - نمونهٔ چهارم.

```
29 )
30
31 # Split the dataset into training and validation sets
```

توصيف عكس كما المحالية عكس المحالية الم

Predicted Caption: a man in a red jacket is walking on a rocky path

Actual Caption: man on bicycle trek high in the mountains



شکل ۱۵: نتایج پیشبینی روی دادههای ارزیابی در حالت رزنت بدون فریز - نمونهٔ پنجم.

```
32 train_dataset, val_dataset = train_test_split(train_dataset, test_size=0.1, random_state=42)
34 # Create a DataLoader for the training dataset
35 train_loader = DataLoader(
     dataset=train_dataset,
     batch_size=BATCH_SIZE,
     num_workers=NUM_WORKERS,
     shuffle=True,
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
41 )
43 # Create a DataLoader for the validation dataset
44 val_loader = DataLoader(
     dataset=val_dataset,
     batch_size=BATCH_SIZE,
     num_workers=NUM_WORKERS,
     shuffle=False,
     collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
50 )
52 # Set the number of training epochs
```

توصيف عكس ١٩٣٤ توصيف عكس

```
num_epochs = 20
55 # Create empty lists to store the training and validation loss
56 train_loss = []
57 val_loss = []
59 # Training loop
60 for epoch in range(num_epochs):
      running_train_loss = 0
      running_val_loss = 0
      # Training
      model.train()
      for idx, (image, captions) in enumerate(train_loader):
          image, captions = image.to(device), captions.to(device)
          optimizer.zero_grad()
          outputs = model(image, captions)
69
          loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
          loss.backward()
71
          optimizer.step()
          running_train_loss += loss.item()
      average_train_loss = running_train_loss / len(train_loader)
      train_loss.append(average_train_loss)
75
      # Validation
      model.eval()
      with torch.no_grad():
          for idx, (image, captions) in enumerate(val_loader):
              image, captions = image.to(device), captions.to(device)
              outputs = model(image, captions)
              loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
              running_val_loss += loss.item()
      average_val_loss = running_val_loss / len(val_loader)
85
      val_loss.append(average_val_loss)
      print(f'Epoch: {epoch+1} - Train Loss: {average_train_loss:.4f} - Val Loss: {average_val_loss
      :.4f}')
      # Save the model after each epoch
      torch.save(model.cpu().state_dict(), 'ModelwithVal.pth')
      model.cuda()
94 # Plotting the training and validation loss
95 plt.figure(figsize=(8, 6), dpi=80)
96 plt.plot(train_loss, label='Training Loss')
```

توصيف عكس كمال المحالات المحال

```
97 plt.plot(val_loss, label='Validation Loss')
98 plt.title("Training and Validation Loss per Epoch")
99 plt.xlabel("Epoch")
100 plt.ylabel("Loss")
101 plt.grid(True, linestyle='--', linewidth=0.5)
102 plt.xticks(range(num_epochs))
103 plt.legend()
104 plt.savefig("loss_plot.pdf", format='pdf', bbox_inches='tight')
105 plt.show()
```

نتایج به صورتی است که در شکل ۱۶ آمده است. اینباز برای عدم تکرار فرآیند قبل و صرفاً برای نمایش نتایج تنها برای ۲۰ دوره آموزش دادیم. نتیجه روی سه دادهٔ تست تصادفی هم در مورد ارزیابی قرار گرفته که در دفتر چه کد گوگل کولب قابل مشاهده است و این جا برای جلوگیری از ازد حام بیش تر نیاورده ایم.

```
import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torchvision.transforms as T
5 from torch.utils.data import DataLoader
6 from sklearn.model_selection import train_test_split
7 import matplotlib.pyplot as plt
9 # Set the batch size, number of workers, and pad index
10 BATCH_SIZE = 256
11 NUM_WORKERS = 10
pad_idx = train_dataset.vocab.stoi["<PAD>"]
# Define the transformations to be applied, including resizing, random cropping,
# converting to tensor, and normalization using ResNet statistics
16 transforms = T.Compose([
     T.Resize(256),
     T.RandomCrop(224),
     T. ToTensor(),
      T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225))
21 ])
23 # Create an instance of the FlickrDataset for training
24 train_dataset = FlickrDataset(
     root_dir="/content/Dataset/Images",
     caption_file="captions.txt",
     transform=transforms,
     data_type='train'
29 )
31 # Split the dataset into training and validation sets
32 train_dataset, val_dataset = train_test_split(train_dataset, test_size=0.1, random_state=42)
```

توصيف عکس
GARASHI
توصيف عکس

```
34 # Create a DataLoader for the training dataset
35 train_loader = DataLoader(
      dataset=train_dataset,
      batch_size=BATCH_SIZE,
      num_workers=NUM_WORKERS,
      shuffle=True,
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
41 )
42
43 # Create a DataLoader for the validation dataset
44 val_loader = DataLoader(
      dataset=val_dataset,
      batch_size=BATCH_SIZE,
      num_workers=NUM_WORKERS,
      shuffle=False,
      collate_fn=Apppadd(pad_idx=pad_idx, batch_first=True)
50 )
52 # Set the number of training epochs
num_epochs = 20
55 # Create empty lists to store the training and validation loss
56 train_loss = []
57 val_loss = []
59 # Training loop
60 for epoch in range(num_epochs):
      running_train_loss = 0
      running_val_loss = 0
      # Training
      model.train()
      for idx, (image, captions) in enumerate(train_loader):
          image, captions = image.to(device), captions.to(device)
          optimizer.zero_grad()
68
          outputs = model(image, captions)
          loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
          loss.backward()
71
          optimizer.step()
          running_train_loss += loss.item()
      average_train_loss = running_train_loss / len(train_loader)
      train_loss.append(average_train_loss)
75
      # Validation
```

توصيف عكس GARASIM

```
model.eval()
      with torch.no_grad():
          for idx, (image, captions) in enumerate(val_loader):
               image, captions = image.to(device), captions.to(device)
81
               outputs = model(image, captions)
              loss = criterion(outputs.view(-1, vocab_size), captions.view(-1))
              running_val_loss += loss.item()
       average_val_loss = running_val_loss / len(val_loader)
      val_loss.append(average_val_loss)
      print(f'Epoch: {epoch+1} - Train Loss: {average_train_loss:.4f} - Val Loss: {average_val_loss
      :.4f}')
      # Save the model after each epoch
      torch.save(model.cpu().state_dict(), 'ModelwithValNoFreeze.pth')
       model.cuda()
94 # Plotting the training and validation loss
95 plt.figure(figsize=(8, 6), dpi=80)
96 plt.plot(train_loss, label='Training Loss')
97 plt.plot(val_loss, label='Test Loss')
98 plt.title("Training and Testing Loss per Epoch")
99 plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.grid(True, linestyle='--', linewidth=0.5)
plt.xticks(range(num_epochs))
103 plt.legend()
plt.savefig("lossplotwithval2.pdf", format='pdf', bbox_inches='tight')
105 plt.show()
```

در ادامه از دستورات زیر برای مقایسهٔ قسمتهای اول و دوم و نمایش نمودار خطا برای هر دوی آنها استفاده می کنیم و نتایج به صورتی است که در شکل ۱۷ نشان داده شده است.

```
# Import the necessary libraries
import matplotlib.pyplot as plt

# Set the figure size
plt.figure(figsize=(10, 6))

# Set the plot title and axis labels
plt.title("Training Loss per Epoch")
plt.xlabel("Epoch")
plt.ylabel("Loss")

# Plot the training loss curves
plt.plot(train_loss, label='Without Fine-tuning (Freeze)')
```

توصيف عكس ٩٩٨٥ توصيف عكس



شکل ۱۶: نمودار خطای آموزش و تست در حالت رزنت بدون فریز.

```
plt.plot(NoFreeze_train_loss, label='With Fine-tuning (No Freeze)')

# Add a legend to the plot
plt.legend()

# Add grid lines to the plot
plt.grid(True)

# Adjust the layout for a more compact and professional appearance
plt.tight_layout()

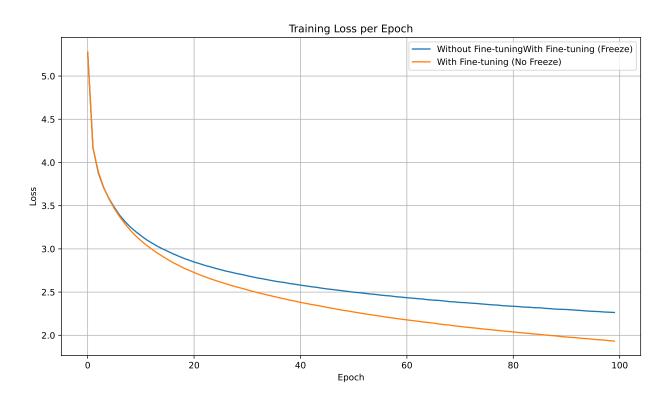
# Add a legend
plt.legend()

# Add a legend
plt.legend()

# Save the plot as a PDF file
plt.savefig("comparelossplot2.pdf", format='pdf', bbox_inches='tight')
```

وصيف عكس ما المحمد المح

```
30
31 # Display the plot
32 plt.show()
```



شكل ١٧: نمودار خطا در حالت رزنت فريز و بدون فريز.

به عنوان یک حرکت اضافی همان حالت قبلی (پاسخ قسمت ۲) را با استفاده از LSTM دوطرفه هم آزمودیم (bidirectional=True):

```
#ResNet Model
class ResNet(nn.Module):

def __init__(self, embed_size, train_resnet=False):
    """

Initialize the ResNet model with a specified embedding size.

Args:
    embed_size (int): Size of the embedding output.
    train_resnet (bool): Whether to train the ResNet backbone or not.

"""
super(ResNet, self).__init__()

# Load the pretrained ResNet-18 model
resnet = models.resnet18(pretrained=True)

# Set the requires_grad flag of the ResNet parameters
# to control whether they are trainable or not
```

F9 Tegnie 32 Teg

```
for param in resnet.parameters():
              param.requires_grad_(train_resnet)
          # Extract the modules of the ResNet model up to the last fully connected layer
          modules = list(resnet.children())[:-1]
          # Create the ResNet backbone with the extracted modules
          self.resnet = nn.Sequential(*modules)
          # Create the embedding layer
          self.embed = nn.Linear(resnet.fc.in_features, embed_size)
          # Activation function and dropout layer
          self.relu = nn.ReLU()
31
          self.dropout = nn.Dropout(0.5)
      def forward(self, images):
34
          Forward pass of the ResNet model.
          Args:
              images (tensor): Input images tensor.
40
          Returns:
              features (tensor): Embedded features tensor.
          features = self.dropout(self.relu(self.resnet(images)))
          features = features.view(features.size(0), -1)
          features = self.embed(features)
          return features
49 # LSTM Network
50 class LSTM(nn.Module):
      def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5):
          Initialize the LSTM model with specified sizes and parameters.
          Args:
              embed_size (int): Size of the word embedding.
              hidden_size (int): Size of the hidden state of the LSTM.
              vocab_size (int): Size of the vocabulary.
              num_layers (int): Number of layers in the LSTM (default: 1).
              drop_prob (float): Dropout probability (default: 0.5).
          super().__init__()
```

۵۰ توصيف عکس

```
# Word embedding layer
           self.embedding = nn.Embedding(vocab_size, embed_size)
           # LSTM layer
           self.lstm = nn.LSTM(embed_size, hidden_size, num_layers=num_layers, batch_first=True,
       bidirectional = True)
           # Linear layer for prediction
           self.linear = nn.Linear(hidden_size*2, vocab_size)
71
           # Dropout layer
           self.dropout = nn.Dropout(drop_prob)
      def forward(self, features, captions):
           Forward pass of the LSTM model.
           Args:
               features (tensor): Image features tensor.
81
               captions (tensor): Captions tensor.
           Returns:
               x (tensor): Output tensor.
           # Vectorize the caption by passing it through the embedding layer
           embeds = self.dropout(self.embedding(captions[:, :-1]))
           # Concatenate the features and captions
          x = torch.cat((features.unsqueeze(1), embeds), dim=1)
91
           # Pass through the LSTM layer
          x, _ = self.lstm(x)
           # Pass through the linear layer for prediction
          x = self.linear(x)
97
          return x
100
      def generate_caption(self, inputs, hidden=None, max_len=20, vocab=None):
102
103
           Generate captions given the image features.
           Args:
105
               inputs (tensor): Input tensor of image features.
```

قوصيف عكس معمد توصيف عكس هم المعمد ال

```
hidden (tuple): Hidden state of the LSTM (default: None).
107
               max_len (int): Maximum length of the generated caption (default: 20).
108
               vocab (Vocab): Vocabulary object (default: None).
           Returns:
               caption (list): Generated caption as a list of words.
           # Inference part
           # Given the image features, generate the captions
116
           batch_size = inputs.size(0)
118
           captions = []
119
120
           for i in range(max_len):
               output, hidden = self.lstm(inputs, hidden)
               output = self.linear(output)
               output = output.view(batch_size, -1)
125
               # Select the word with the highest value
126
               predicted_word_idx = output.argmax(dim=1)
128
               # Save the generated word
129
               captions.append(predicted_word_idx.item())
               # End if <EOS> is detected
               if vocab.itos[predicted_word_idx.item()] == "<EOS>":
                   break
134
135
               # Send the generated word as the next caption
136
               inputs = self.embedding(predicted_word_idx.unsqueeze(0))
           # Convert the vocabulary indices to words and return the sentence
139
           return [vocab.itos[idx] for idx in captions]
140
142 # Utilizing the powerful fusion of ResNet and LSTM for image captioning
  class CNNtoRNN(nn.Module):
       def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1, drop_prob=0.5,
       train_resnet=False):
           Initialize the CNNtoRNN model with specified sizes and parameters.
146
147
               embed_size (int): Size of the word embedding.
149
               hidden_size (int): Size of the hidden state of the LSTM.
```

```
vocab_size (int): Size of the vocabulary.
               num_layers (int): Number of layers in the LSTM (default: 1).
               drop_prob (float): Dropout probability (default: 0.5).
               train_resnet (bool): Whether to train the ResNet backbone or not (default: False).
154
           super().__init__()
156
           # Encoder (ResNet)
           self.encoder = ResNet(embed_size, train_resnet)
159
           # Decoder (LSTM)
           self.decoder = LSTM(embed_size, hidden_size, vocab_size, num_layers, drop_prob)
162
164
      def forward(self, images, captions):
           Forward pass of the CNNtoRNN model.
167
           Args:
               images (tensor): Input images tensor.
169
               captions (tensor): Captions tensor.
           Returns:
               outputs (tensor): Output tensor.
           # Pass the images through the encoder (ResNet) to get features
           features = self.encoder(images)
           # Pass the features and captions through the decoder (LSTM) to get outputs
178
           outputs = self.decoder(features, captions)
181
           return outputs
183
# Initialize the model, loss function, and optimizer
186 NoFreeze_model_bi = CNNtoRNN(embed_size, hidden_size, vocab_size, num_layers, train_resnet=True).
       to(device)
187 criterion = nn.CrossEntropyLoss(ignore_index=train_dataset.vocab.stoi["<PAD>"])
188 optimizer = optim.Adam(NoFreeze_model_bi.parameters(), lr=learning_rate)
```

با اجرای این فرآیند و دستورات نمودار خطا به صورتی است که در شکل ۱۸ نمایش داده شده و همان طور که از نتایج مشخص است هرچند مدل همچنان جا برای آموزش دارد؛ اما نتایج در تعداد دوره های کم تر به مراتب بهتر از حالات قبلی شده است. در این پیاده سازی از مقیاس کردن گرادیان ها یا () torch.cuda.amp.GradScaler هم استفاده شده است تا سرعت و بهره وری فرآیند آموزش بالاتر رود. نتایج کپشن های پیش بینی شده روی داده های ارزیابی اما نشان می دهد که هرچند نتایج آموزش مدل قابل قبول نبوده و به نوعی اورفیت رخ داده است.

وصيف عكس ١٩٩٤ وصيف عكس ١٩٩٤

```
Epoch: 1 - Train Loss: 4.764396175637946
2 Epoch: 2 - Train Loss: 2.955980390935511
3 Epoch: 3 - Train Loss: 2.192415646739773
4 Epoch: 4 - Train Loss: 1.7028564723221573
5 Epoch: 5 - Train Loss: 1.3591378785513497
6 Epoch: 6 - Train Loss: 1.106173901707976
7 Epoch: 7 - Train Loss: 0.9153213192532946
8 Epoch: 8 - Train Loss: 0.7661450617796891
9 Epoch: 9 - Train Loss: 0.6474878558745751
10 Epoch: 10 - Train Loss: 0.5507056039946896
13 .
14 Epoch: 20 - Train Loss: 0.1382495445700792
15 Epoch: 21 - Train Loss: 0.12208282926699499
16 Epoch: 22 - Train Loss: 0.10844990913267735
17 Epoch: 23 - Train Loss: 0.09631544533934626
18 Epoch: 24 - Train Loss: 0.08601176170827625
19 Epoch: 25 - Train Loss: 0.07654622593126097
20 Epoch: 26 - Train Loss: 0.06880965792424196
21 Epoch: 27 - Train Loss: 0.06132213951095001
22 Epoch: 28 - Train Loss: 0.05519910268746056
23 Epoch: 29 - Train Loss: 0.04966954329422304
24 Epoch: 30 - Train Loss: 0.045182802746971173
```

در ادامه از دستورات زیر برای مقایسهٔ قسمتهای اول و دوم و این قسمت و نمایش نمودار خطا برای هر سهی اینها استفاده میکنیم و نتایج بهصورتی است که در شکل ۱۹ نشان داده شده است.

```
# Import the necessary libraries
import matplotlib.pyplot as plt

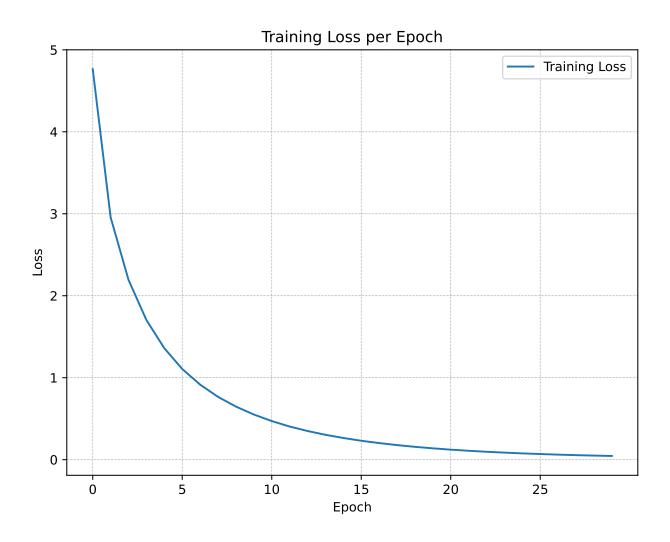
# Set the figure size
plt.figure(figsize=(10, 6))

# Set the plot title and axis labels
plt.title("Training Loss per Epoch")
plt.xlabel("Epoch")
plt.ylabel("Loss")

# Plot the training loss curves
plt.plot(train_loss, label='Without Fine-tuning (Freeze)')
plt.plot(NoFreeze_train_loss, label='With Fine-tuning (No Freeze) + Bidirectional LSTM')

# Add a legend to the plot
plt.legend()
```

وصيف عكس ١٩٥٤ وصيف عكس ١٩٥٤



شكل ۱۸: نمودار خطا در حالت رزنت بدون فريز (Bidirectional LSTM).

```
# Add grid lines to the plot
plt.grid(True)

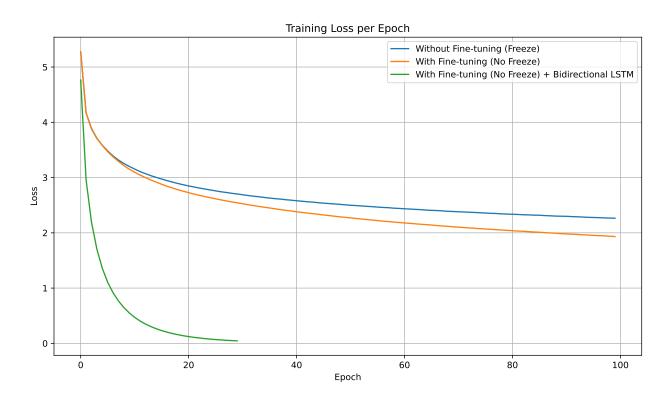
# Adjust the layout for a more compact and professional appearance
plt.tight_layout()

# Add a legend
plt.legend()

# Save the plot as a PDF file
plt.savefig("comparelossplot3.pdf", format='pdf', bbox_inches='tight')

# Display the plot
plt.show()
```

توصيف عکس ^{©ARAS™} توصيف عکس



شكل ۱۹: نمودار خطا در حالت رزنت فريز، بدون فريز و بدون فريز (Bidirectional LSTM).