# Automatic Image Captioning With PyTorch

Deepesh Garg  ·  *Follow*

5 min read  ·  Aug 20, 2020

▶ Listen          ⬆ Share          ••• More

> *"It's going to be interesting to see how society deals with artificial intelligence, but it will definitely be cool."*
>
> *- Colin Angle*

This is my first open source project . I was selected as a Participant for Open Source Contributions at **Student Code-in** . Actually, It was a two months programme where I was selected for contributions to a **Computer Vision** Project **: Image Captioning .** In this project, I design and train a CNN-RNN (Convolutional Neural Network — Recurrent Neural Network) model for automatically generating image captions. In this case, LSTM (Long Short Term Memory), is used which is a special kind of RNN that includes a memory cell, in order to maintain the information for a longer period of time.

The network is trained on the **Microsoft Common Objects in COntext** (MS COCO) dataset. The image captioning model is displayed below.
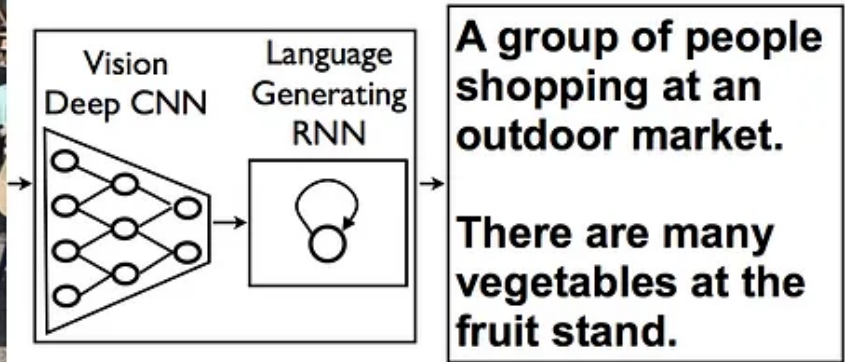
Image Source

## Dataset Used — MS COCO Dataset

The COCO dataset is one of the largest, publicly available image datasets and it is meant to represent realistic scenes. What I mean by this is that COCO does not overly pre-process images, instead these images come in a variety of shapes with a variety of objects and environment/lighting conditions that closely represent what you might get if you compiled images from many different cameras around the world.

To explore the dataset, you can check out the **dataset website**

## Captions

COCO is a richly labeled dataset; it comes with class labels, labels for segments of an image, and a set of captions for a given image . Here is an example :

Image Source — Udacity

## Visualize the Dataset

The Microsoft **C**ommon **O**bjects in **CO**ntext (MS COCO) dataset is a large-scale dataset for scene understanding. The dataset is commonly used to train and benchmark object detection, segmentation, and captioning algorithms.

```python
1   import os
2   import sys
3   sys.path.append('/opt/cocoapi/PythonAPI')
4   from pycocotools.coco import COCO
5
6   # initialize COCO API for instance annotations
7   dataDir = '/home/Project/Udacity-Computer-Vision-Nanodegree-Program/project_2_image_cap
8   dataType = 'val2014'
9   instances_annFile = os.path.join(dataDir, 'annotations/instances_{}.json'.format(dataTy
10  coco = COCO(instances_annFile)
11
12  # initialize COCO API for caption annotations
13  captions_annFile = os.path.join(dataDir, 'annotations/captions_{}.json'.format(dataType
14  coco_caps = COCO(captions_annFile)
15
16  # get image ids
17  ids = list(coco.anns.keys())
```

**COCO_API.py** hosted with ❤ by **GitHub**      **view raw**

```python
1   import numpy as np
2   import skimage.io as io
3   import matplotlib.pyplot as plt
4   %matplotlib inline
5
6   # pick a random image and obtain the corresponding URL
7   ann_id = np.random.choice(ids)
8   img_id = coco.anns[ann_id]['image_id']
9   img = coco.loadImgs(img_id)[0]
10  url = img['coco_url']
11
12  # print URL and visualize corresponding image
13  print(url)
14  I = io.imread(url)
15  plt.axis('off')
16  plt.imshow(I)
17  plt.show()
18
19  # load and display captions
20  annIds = coco_caps.getAnnIds(imgIds=img['id']);
21  anns = coco_caps.loadAnns(annIds)
22  coco_caps.showAnns(anns)
```
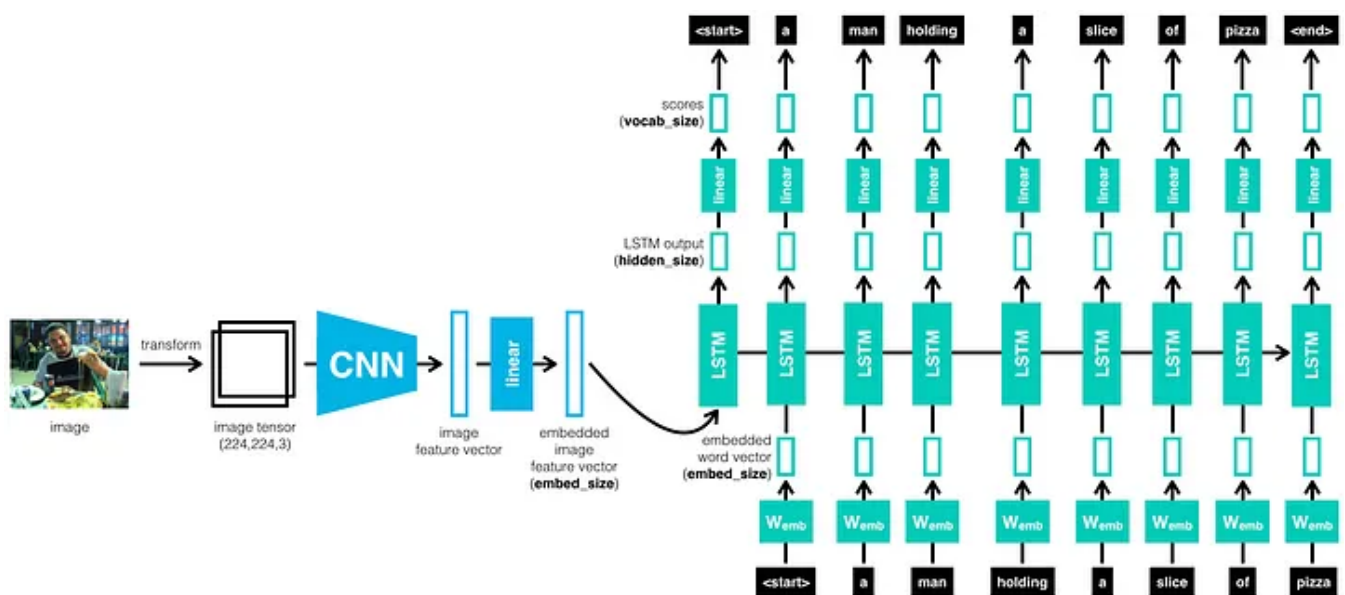
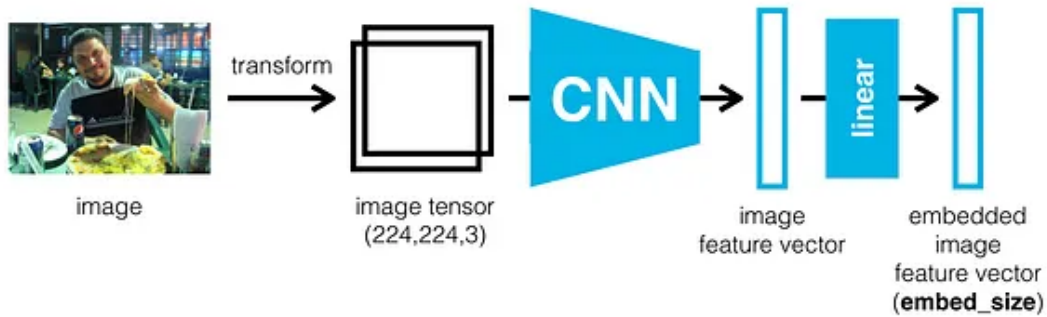**Visualize.py** hosted with ❤ by **GitHub**      **view raw**

```
A herd of animals grazing on a lush green field.
A field with many cows and they're all laying down
Cattle lying on the grass in a field while birds fly above them.
A herd of cattle graze in a grassy field.
Birds flying over cows in a green pasture.
```

## The CNN-RNN Architecture



### Encoder CNN

The encoder that I used was the pre-trained **ResNet**-50 architecture (with the final fully-connected layer removed) to extract features from a batch of pre-processed images. The output is then flattened to a vector, before being passed through a `Linear` layer to transform the feature vector to have the same size as the word embedding.

image    image tensor    image    embedded
   (224,224,3)    feature vector    image
                         feature vector
                         (embed_size)

```python
1    class EncoderCNN(nn.Module):
2        def __init__(self, embed_size):
3            super(EncoderCNN, self).__init__()
4            resnet = models.resnet50(pretrained=True)
5            for param in resnet.parameters():
6                param.requires_grad_(False)
7
8            modules = list(resnet.children())[:-1]
9            self.resnet = nn.Sequential(*modules)
10           self.embed = nn.Linear(resnet.fc.in_features, embed_size)
11           self.batch= nn.BatchNorm1d(embed_size,momentum = 0.01)
12           self.embed.weight.data.normal_(0., 0.02)
13           self.embed.bias.data.fill_(0)
14
15       def forward(self, images):
16           features = self.resnet(images)
17           features = features.view(features.size(0), -1)
18           features = self.batch(self.embed(features))
19           return features
```
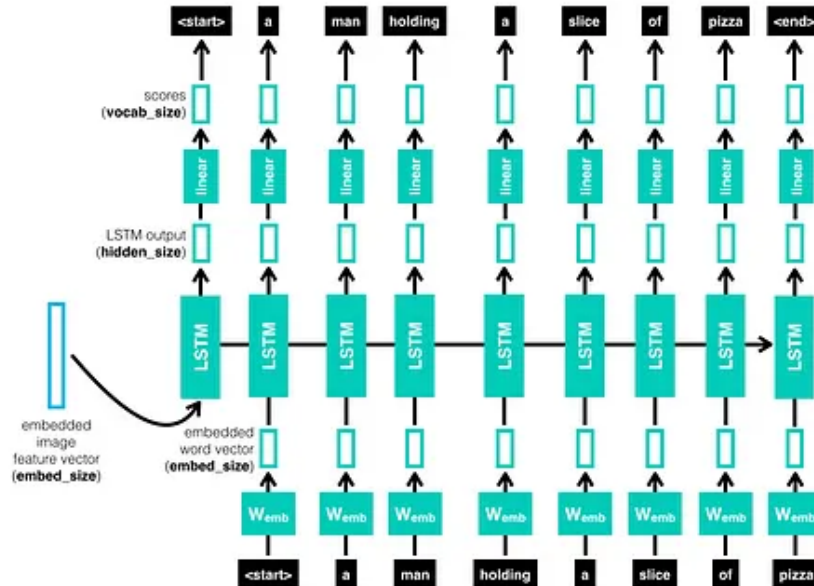
Encoder_CNN.py hosted with ❤ by GitHub            view raw

## Decoder RNN

The job of the RNN is to decode the process vector and turn it into a sequence of words. Thus, this portion of the network is often called a decoder. In this case, **LSTM (Long Short Term Memory),** is used which is a special kind of RNN that includes a memory cell, in order to maintain the information for a longer period of time.

```python
class DecoderRNN(nn.Module):
    def __init__(self, embed_size, hidden_size, vocab_size, num_layers=1):
        super(DecoderRNN, self).__init__()
        self.num_layers = num_layers
        self.hidden_size = hidden_size
        self.embed_size= embed_size
        self.drop_prob= 0.2
        self.vocabulary_size = vocab_size
        self.lstm = nn.LSTM(self.embed_size, self.hidden_size , self.num_layers,batch_
        self.dropout = nn.Dropout(self.drop_prob)
        self.embed = nn.Embedding(self.vocabulary_size, self.embed_size)
        self.linear = nn.Linear(hidden_size, self.vocabulary_size)
        self.embed.weight.data.uniform_(-0.1, 0.1)
        self.linear.weight.data.uniform_(-0.1, 0.1)
        self.linear.bias.data.fill_(0)

    def forward(self, features, captions):
        embeddings = self.embed(captions)
        features = features.unsqueeze(1)
        embeddings = torch.cat((features, embeddings[:, :-1,:]), dim=1)
        hiddens, c = self.lstm(embeddings)
        outputs = self.linear(hiddens)
        return outputs
```

**Decoder_RNN.py** hosted with ❤ by **GitHub**                                                      **view raw**

## Caption Pre-Processing

The captions also need to be pre-processed and prepped for training. In this example, for generating captions, I aimed to create a model that predicts the next

token of a sentence from previous tokens, So I turned the caption associated with any image into a list of tokenized words, before casting it to a PyTorch tensor that we can use to train the network.

## Tokenizing Captions

First, we iterate through all of the training captions and create a dictionary that maps all unique words to a numerical index. So, every word we come across will have a corresponding integer value that can be found in this dictionary. The words in this dictionary are referred to as vocabulary.

```python
1   sample_caption = 'A person doing a trick on a rail while riding a skateboard.'
2   import nltk
3
4   sample_tokens = nltk.tokenize.word_tokenize(str(sample_caption).lower())
5   print(sample_tokens)
6   sample_caption = []
7
8   start_word = data_loader.dataset.vocab.start_word
9   print('Special start word:', start_word)
10  sample_caption.append(data_loader.dataset.vocab(start_word))
11  print(sample_caption)
```

**Tokenize.py** hosted with ❤ by **GitHub**                                    **view raw**

```python
1   # Preview the word2idx dictionary.
2   dict(list(data_loader.dataset.vocab.word2idx.items())[:10])
```

**Preview.py** hosted with ❤ by **GitHub**                                     **view raw**

## Output

```
{'<start>': 1,
'<end>': 0,
'<unk>': 2,
'a': 3,
'and': 6,
'clean': 5,
'decorated': 8,
'empty': 9,
'very': 4,
'well': 7}
```

```
1    # Modify the minimum word count threshold.
2    vocab_threshold = 6
3
4    # Obtain the data loader.
5    data_loader = get_loader(transform=transform_train,
6                             mode='train',
7                             batch_size=batch_size,
8                             vocab_threshold=vocab_threshold,
9                             vocab_from_file=False)
```

**Vocab.py** hosted with ❤ by **GitHub**                                                    **view raw**

```
1    # Print the total number of keys in the word2idx dictionary.
2    print('Total number of tokens in vocabulary:', len(data_loader.dataset.vocab))
```

**Total_Tokens.py** hosted with ❤ by **GitHub**                                            **view raw**

## Output

```
Total number of tokens in vocabulary: 8099
```

## Conversion Of Word To Vectors

The words first must be turned into a numerical representation so that a network can use normal loss functions and optimizers to calculate the difference between a predicted word and ground truth word (from a known, training caption) . So, we typically turn a sequence of words into a sequence of numerical values; a vector of numbers where each number maps to a specific word in our vocabulary.



## Training The Model

We have two model components, i.e. encoder and decoder, we train them jointly by passing the output of the encoder, which is the latent space vector, to the decoder, which, in turn, is the recurrent neural network.

No. Of Epochs = 1

Batch Size = 32

```python
1   import torch
2   import torch.nn as nn
3   from torchvision import transforms
4   import sys
5   sys.path.append('/opt/cocoapi/PythonAPI')
6   from pycocotools.coco import COCO
7   from data_loader import get_loader
8   from model import EncoderCNN, DecoderRNN
9   import math
10
11
12  ## TODO #1: Select appropriate values for the Python variables below.
13  batch_size = 32          # batch size
14  vocab_threshold = 6        # minimum word count threshold
15  vocab_from_file = True    # if True, load existing vocab file
16  embed_size = 512           # dimensionality of image and word embeddings
17  hidden_size = 512          # number of features in hidden state of the RNN decoder
18  num_epochs = 1             # number of training epochs (1 for testing)
19  save_every = 1             # determines frequency of saving model weights
20  print_every = 200          # determines window for printing average loss
21  log_file = 'training_log.txt'      # name of file with saved training loss and perplex
22
23  # (Optional) TODO #2: Amend the image transform below.
24  transform_train = transforms.Compose([
25      transforms.Resize(256),                       # smaller edge of image resized to
26      transforms.RandomCrop(224),                   # get 224x224 crop from random loc
27      transforms.RandomHorizontalFlip(),            # horizontally flip image with pro
28      transforms.ToTensor(),                        # convert the PIL Image to a tenso
29      transforms.Normalize((0.485, 0.456, 0.406),    # normalize image for pre-trained
30                           (0.229, 0.224, 0.225))])
31
32  # Build data loader.
33  data_loader = get_loader(transform=transform_train,
34                           mode='train',
35                           batch_size=batch_size,
36                           vocab_threshold=vocab_threshold,
37                           vocab_from_file=vocab_from_file)
38
39  # The size of the vocabulary.
40  vocab_size = len(data_loader.dataset.vocab)
41
42  # Initialize the encoder and decoder.
43  encoder = EncoderCNN(embed_size)
44  decoder = DecoderRNN(embed_size, hidden_size, vocab_size)
45
46  # Move models to GPU if CUDA is available.
47  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
48  encoder.to(device)
```

```
48    encoder.to(device)
49    decoder.to(device)
50
51    # Define the loss function.
52    criterion = nn.CrossEntropyLoss().cuda() if torch.cuda.is_available() else nn.CrossEntr
53
54    # TODO #3: Specify the learnable parameters of the model.
55    params = list(decoder.parameters()) + list(encoder.embed.parameters()) + list(encoder.k
56
57    # TODO #4: Define the optimizer.
58    optimizer = torch.optim.Adam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08)
59    # optimizer = torch.optim.Adam(params, lr=0.01, betas=(0.9, 0.999), eps=1e-08)
60    # optimizer = torch.optim.RMSprop(params, lr=0.01, alpha=0.99, eps=1e-08)
61
62    # Set the total number of training steps per epoch.
63    total_step = math.ceil(len(data_loader.dataset.caption_lengths) / data_loader.batch_sam
```

To figure out how well our model is doing, we can look at how the training loss and perplexity evolve during training — and for the purposes of this project, we can amend the hyperparameters based on this information. However, this will not tell you if your model is overfitting to the training data, and, unfortunately, overfitting is a problem that is commonly encountered when training image captioning models. For this project, you need not worry about overfitting. This project does not have strict requirements regarding the performance of your model, and you just need to demonstrate that your model has learned something when you generate captions on the test data.

## Prediction Function

The **get_prediction** function was used to loop over images in the test dataset and print model's predicted caption.

```
1    def get_prediction():
2        orig_image, image = next(iter(data_loader))
3        plt.imshow(np.squeeze(orig_image))
4        plt.title('Sample Image')
5        plt.show()
6        image = image.to(device)
7        features = encoder(image).unsqueeze(1)
8        output = decoder.sample(features)
9        sentence = clean_sentence(output)
10       print(sentence)
```

Prediction.py hosted with ❤ by **GitHub**                                          view raw

## Predicted Results



**A large elephant standing next to a tree .**

**A person holding a cell phone in their hands .**

## More Predictions

This is my complete open source project on **GitHub** .

**References**

**1.** Show, Attend and Tell: Neural Image Caption Generation with Visual Attention( https://arxiv.org/pdf/1502.03044.pdf)

**2.** https://github.com/sauravraghuvanshi/Udacity-Computer-Vision-Nanodegree-Program/tree/master/project_2_image_captioning_project

# Happy Learning !

Image Captioning   Deep Learning   Neural Networks   Pytorch

# Written by Deepesh Garg

6 Followers

CSE Undergrad 🎓 || ML Enthusiast 💻 || Coder 👨‍💻 || Optimistic About New And Developing AI Technologies 💯

---

## More from Deepesh Garg



👤 Deepesh Garg

## My Student Code-In, 2020 Journey

Technology innovation is starting to explode and having open-source material out there really helps this explosion. You get students and...

4 min read · Sep 1, 2020

👏 52          💬                                            🔖         •••

---

See all from Deepesh Garg

---

# Recommended from Medium



👤 Ranjithkumar Panjabikesan, Enterprise Architect

## Build a Image Captioning App leveraging Hugging face, BLIP Image Captioning Model and Gradio

In this article, we will look at how we can harness the combined power of Hugging face, Salesforce BLIP Image captioning models, Gradio and...

7 min read · Oct 16, 2023

👏 30      💬                                                    🔖⁺        •••

Everton Gomede, PhD

# Progressive Growing in Generative Adversarial Networks

Introduction

6 min read · Oct 18, 2023

## Lists

### Natural Language Processing
1200 stories · 668 saves

### Practical Guides to Machine Learning
10 stories · 1058 saves

### data science and AI
40 stories · 72 saves

### Staff Picks
582 stories · 750 saves

$$Discriminator \; \frac{D(G(z)}{}$$

Mohammad Daraeee

## Generative Adversarial Networks (GANs)

link to codes of this article:
https://github.com/Mhddaraaa/start/tree/main/Generative_Adversarial_Networks(GANs)

6 min read · Jan 25, 2024

S Subash Palvel

## Image Captioning using CNN-RNN Architectures

Image captioning is the task of generating a textual description for an image. It combines computer vision and natural language processing...

2 min read  ·  Sep 13, 2023



Dream AI

## Multi-Label Video Classification using PyTorch Lightning Flash

Author: Rafay Farhan at DreamAI Software (Pvt) Ltd
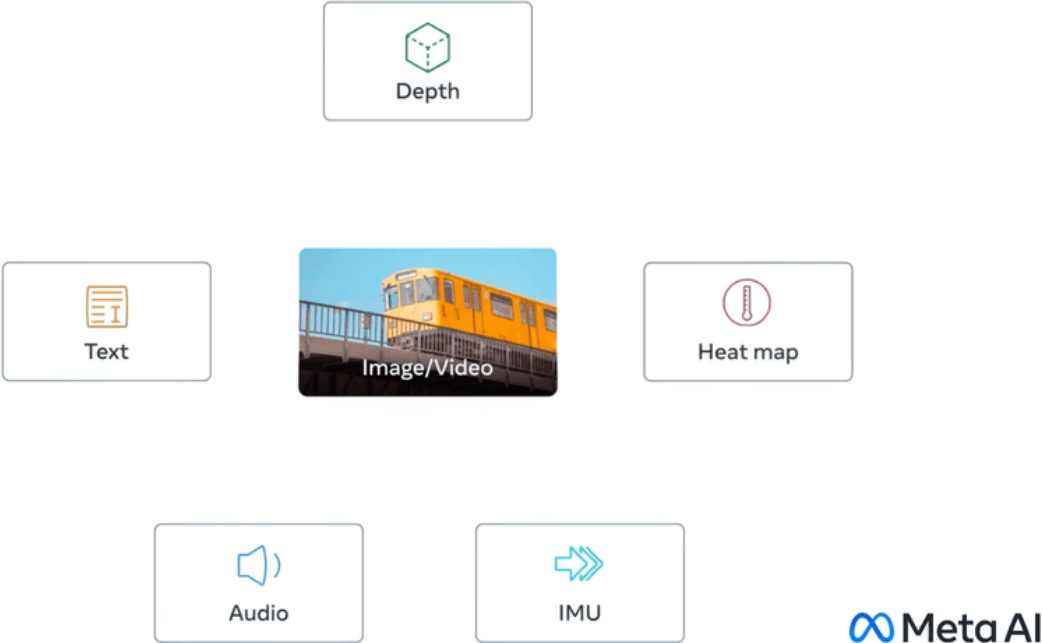
14 min read  ·  Sep 13, 2023

168    1

Fahim Rustamy, PhD in Towards Data Science

## CLIP Model and The Importance of Multimodal Embeddings

CLIP, which stands for Contrastive Language-Image Pretraining, is a deep learning model developed by OpenAI in 2021. CLIP's embeddings for...

10 min read · Dec 11, 2023

310          1

**See more recommendations**