

- ✓ 1) Initial Imports and loading the utils function. The dataset is used is [Flickr 8k](#) from kaggle. Custom dataset and dataloader is implemented in [this](#) notebook.

```
#location of the data
data_location = "../input/flickr8k"
!ls $data_location
```

Images captions.txt

```
#reading the text data
import pandas as pd
caption_file = data_location + '/captions.txt'
df = pd.read_csv(caption_file)
print("There are {} image to captions".format(len(df)))
df.head(7)
```

There are 40455 image to captions

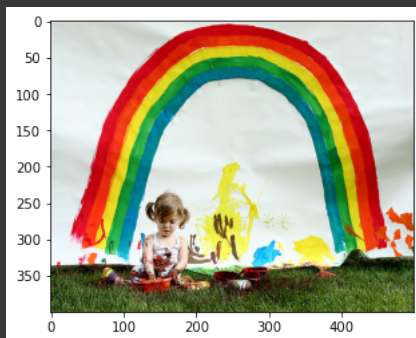
	image	caption
0	1000268201_693b08cb0e.jpg	A child in a pink dress is climbing up a set o...
1	1000268201_693b08cb0e.jpg	A girl going into a wooden building .
2	1000268201_693b08cb0e.jpg	A little girl climbing into a wooden playhouse .
3	1000268201_693b08cb0e.jpg	A little girl climbing the stairs to her playh...
4	1000268201_693b08cb0e.jpg	A little girl in a pink dress going into a woo...
5	1001773457_577c3a7d70.jpg	A black dog and a spotted dog are fighting
6	1001773457_577c3a7d70.jpg	A black dog and a tri-colored dog playing with...

```
import matplotlib.pyplot as plt
import matplotlib.image as mpimg

#select any index from the whole dataset
#single image has 5 captions
#so, select indx as: 1,6,11,16...
data_idx = 11

#eg path to be plot: ../input/flickr8k/Images/1000268201_693b08cb0e.jpg
image_path = data_location+"/Images/"+df.iloc[data_idx,0]
img=mpimg.imread(image_path)
plt.imshow(img)
plt.show()

#image consits of 5 captions,
#showing all 5 captions of the image of the given idx
for i in range(data_idx,data_idx+5):
    print("Caption:",df.iloc[i,1])
```



Caption: A little girl is sitting in front of a large painted rainbow .  
 Caption: A small girl in the grass plays with fingerpaints in front of a white canvas  
 Caption: There is a girl with pigtails sitting in front of a rainbow painting .  
 Caption: Young girl with pigtails painting outside in the grass .  
 Caption: A man lays on a bench while his dog sits by his

```
#imports
import os
from collections import Counter
import spacy
import torch
from torch.nn.utils.rnn import pad_sequence
from torch.utils.data import DataLoader, Dataset
import torchvision.transforms as T
```

```
from PIL import Image
```

```
#using spacy for the better text tokenization
spacy_eng = spacy.load("en")
```

```
#example
text = "This is a good place to find a city"
[token.text.lower() for token in spacy_eng.tokenizer(text)]
```

```
['this', 'is', 'a', 'good', 'place', 'to', 'find', 'a', 'city']
```

```
class Vocabulary:
    def __init__(self, freq_threshold):
        #setting the pre-reserved tokens int to string tokens
        self.itos = {0:"<PAD>", 1:"<SOS>", 2:"<EOS>", 3:"<UNK>"}

        #string to int tokens
        #its reverse dict self.itos
        self.stoi = {v:k for k,v in self.itos.items()}

        self.freq_threshold = freq_threshold

    def __len__(self): return len(self.itos)

    @staticmethod
    def tokenize(text):
        return [token.text.lower() for token in spacy_eng.tokenizer(text)]

    def build_vocab(self, sentence_list):
        frequencies = Counter()
        idx = 4

        for sentence in sentence_list:
            for word in self.tokenize(sentence):
                frequencies[word] += 1

            #add the word to the vocab if it reaches minum frequency threshold
            if frequencies[word] == self.freq_threshold:
                self.stoi[word] = idx
                self.itos[idx] = word
                idx += 1

    def numericalize(self, text):
        """ For each word in the text corresponding index token for that word form the vocab built as list """
        tokenized_text = self.tokenize(text)
        return [ self.stoi[token] if token in self.stoi else self.stoi["<UNK>"] for token in tokenized_text ]
```

```
#testing the vicab class
v = Vocabulary(freq_threshold=1)
```

```
v.build_vocab(["This is a good place to find a city"])
print(v.stoi)
print(v.numericalize("This is a good place to find a city here!!"))
```

```
{'<PAD>': 0, '<SOS>': 1, '<EOS>': 2, '<UNK>': 3, 'this': 4, 'is': 5, 'a': 6, 'good': 7, 'place': 8, 'to': 9, 'find': 10, 'city': 11}
[4, 5, 6, 7, 8, 9, 10, 6, 11, 3, 3, 3]
```

```

class FlickrDataset(Dataset):
    """
    FlickrDataset
    """
    def __init__(self, root_dir, captions_file, transform=None, freq_threshold=5):
        self.root_dir = root_dir
        self.df = pd.read_csv(caption_file)
        self.transform = transform

        #Get image and caption colum from the dataframe
        self.imgs = self.df["image"]
        self.captions = self.df["caption"]

        #Initialize vocabulary and build vocab
        self.vocab = Vocabulary(freq_threshold)
        self.vocab.build_vocab(self.captions.tolist())

    def __len__(self):
        return len(self.df)

    def __getitem__(self, idx):
        caption = self.captions[idx]
        img_name = self.imgs[idx]
        img_location = os.path.join(self.root_dir, img_name)
        img = Image.open(img_location).convert("RGB")

        #apply the transform to the image
        if self.transform is not None:
            img = self.transform(img)

        #numericalize the caption text
        caption_vec = []
        caption_vec += [self.vocab.stoi["<SOS>"]]
        caption_vec += self.vocab.numericalize(caption)
        caption_vec += [self.vocab.stoi["<EOS>"]]

        return img, torch.tensor(caption_vec)

```

```

#defining the transform to be applied
transforms = T.Compose([
    T.Resize((224,224)),
    T.ToTensor()
])

```

```

def show_image(inp, title=None):
    """Imshow for Tensor."""
    inp = inp.numpy().transpose((1, 2, 0))
    plt.imshow(inp)
    if title is not None:
        plt.title(title)
    plt.pause(0.001) # pause a bit so that plots are updated

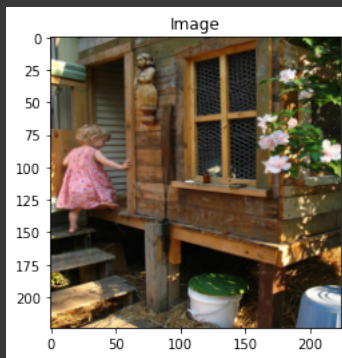
```

```

#testing the dataset class
dataset = FlickrDataset(
    root_dir = data_location+"/Images",
    captions_file = data_location+"/captions.txt",
    transform=transforms
)

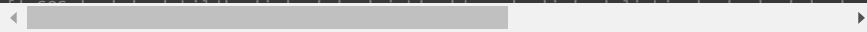
img, caps = dataset[0]
show_image(img, "Image")
print("Token:", caps)
print("Sentence:")
print([dataset.vocab.itos[token] for token in caps.tolist()])

```



Token: tensor([ 1, 4, 28, 8, 4, 195, 151, 17, 32, 67, 4, 353, 11, 711, 8, 24, 3, 496, 5, 2])

Sentence:



```
class CapsCollate:
    """
    Collate to apply the padding to the captions with dataloader
    """
    def __init__(self,pad_idx,batch_first=False):
        self.pad_idx = pad_idx
        self.batch_first = batch_first

    def __call__(self,batch):
        imgs = [item[0].unsqueeze(0) for item in batch]
        imgs = torch.cat(imgs,dim=0)

        targets = [item[1] for item in batch]
        targets = pad_sequence(targets, batch_first=self.batch_first, padding_value=self.pad_idx)
        return imgs,targets
```

```
#writing the dataloader
#setting the constants
BATCH_SIZE = 64
NUM_WORKER = 1

#token to represent the padding
pad_idx = dataset.vocab.stoi["<PAD>"]

data_loader = DataLoader(
    dataset=dataset,
    batch_size=BATCH_SIZE,
    num_workers=NUM_WORKER,
    shuffle=True,
    collate_fn=CapsCollate(pad_idx=pad_idx,batch_first=True)
)
```

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

```
#getting the next batch
batch = next(dataiter)
```

```
#unpacking the batch
images, captions = batch
print(len(captions))
#showing info of image in single batch
print(captions.shape)
```

```
batch = next(dataiter)
```

```
#unpacking the batch
images, captions = batch
print(len(captions[0]))
print(len(captions[1]))
print(len(captions[2]))
```

```
64
torch.Size([64, 25])
29
29
29
```

```
#location of the training data
data_location = "../input/flickr8k"
#copy dataloader
!cp ../input/data-loader/data_loader.py .

#imports
import numpy as np
import torch
from torch.utils.data import DataLoader, Dataset
import torchvision.transforms as T

#custom imports
# from data_loader import FlickrDataset, get_data_loader

cp: cannot stat '../input/data-loader/data_loader.py': No such file or directory
```

## 2) Implementing the Helper function to plot the Tensor image

```
#show the tensor image
import matplotlib.pyplot as plt
def show_image(img, title=None):
    """Imshow for Tensor."""

    #unnormailize
    img[0] = img[0] * 0.229
    img[1] = img[1] * 0.224
    img[2] = img[2] * 0.225
    img[0] += 0.485
    img[1] += 0.456
    img[2] += 0.406

    img = img.numpy().transpose((1, 2, 0))

    plt.imshow(img)
    if title is not None:
        plt.title(title)
    plt.pause(0.001) # pause a bit so that plots are updated
```

```
#Initiate the Dataset and Dataloader

#setting the constants
data_location = "../input/flickr8k"
BATCH_SIZE = 256
# BATCH_SIZE = 6
NUM_WORKER = 4

#defining the transform to be applied
transforms = T.Compose([
    T.Resize(226),
    T.RandomCrop(224),
    T.ToTensor(),
    T.Normalize((0.485, 0.456, 0.406), (0.229, 0.224, 0.225))
])

#testing the dataset class
dataset = FlickrDataset(
    root_dir = data_location+"/Images",
    captions_file = data_location+"/captions.txt",
    transform=transforms
)

#writing the dataloader
data_loader = DataLoader(
    dataset=dataset,
    batch_size=BATCH_SIZE,
    num_workers=NUM_WORKER,
    shuffle=True,
    # batch_first=False
    collate_fn=CapsCollate(pad_idx=pad_idx, batch_first=True)
)

#vocab_size
vocab_size = len(dataset.vocab)

device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
device
```

```
device(type='cuda', index=0)
```

### 3) Defining the Model Architecture

Model is seq2seq model. In the **encoder** pretrained ResNet model is used to extract the features. Decoder, is the implementation of the Bahdanau Attention Decoder. In the decoder model **LSTM cell**.

```
import torch
import numpy as np
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

import torchvision.models as models
from torch.utils.data import DataLoader, Dataset
import torchvision.transforms as T

class EncoderCNN(nn.Module):
    def __init__(self):
        super(EncoderCNN, self).__init__()
        resnet = models.resnet50(pretrained=True)
        for param in resnet.parameters():
            param.requires_grad_(False)

        modules = list(resnet.children())[:-2]
        self.resnet = nn.Sequential(*modules)

    def forward(self, images):
        features = self.resnet(images) # (batch_size, 2048, 7, 7)
        features = features.permute(0, 2, 3, 1) # (batch_size, 7, 7, 2048)
        features = features.view(features.size(0), -1, features.size(-1)) # (batch_size, 49, 2048)
        return features
```

```
#Bahdanau Attention
class Attention(nn.Module):
    def __init__(self, encoder_dim, decoder_dim, attention_dim):
        super(Attention, self).__init__()

        self.attention_dim = attention_dim

        self.W = nn.Linear(decoder_dim, attention_dim)
        self.U = nn.Linear(encoder_dim, attention_dim)

        self.A = nn.Linear(attention_dim, 1)

    def forward(self, features, hidden_state):
        u_hs = self.U(features) # (batch_size, num_layers, attention_dim)
        w_ah = self.W(hidden_state) # (batch_size, attention_dim)

        combined_states = torch.tanh(u_hs + w_ah.unsqueeze(1)) # (batch_size, num_layers, attention_dim)

        attention_scores = self.A(combined_states) # (batch_size, num_layers, 1)
        attention_scores = attention_scores.squeeze(2) # (batch_size, num_layers)

        alpha = F.softmax(attention_scores, dim=1) # (batch_size, num_layers)

        attention_weights = features * alpha.unsqueeze(2) # (batch_size, num_layers, features_dim)
        attention_weights = attention_weights.sum(dim=1) # (batch_size, num_layers)

        return alpha, attention_weights
```

```

#Attention Decoder
class DecoderRNN(nn.Module):
    def __init__(self, embed_size, vocab_size, attention_dim, encoder_dim, decoder_dim, drop_prob=0.3):
        super().__init__()

        #save the model param
        self.vocab_size = vocab_size
        self.attention_dim = attention_dim
        self.decoder_dim = decoder_dim

        self.embedding = nn.Embedding(vocab_size, embed_size)
        self.attention = Attention(encoder_dim, decoder_dim, attention_dim)

        self.init_h = nn.Linear(encoder_dim, decoder_dim)
        self.init_c = nn.Linear(encoder_dim, decoder_dim)
        self.lstm_cell = nn.LSTMCell(embed_size+encoder_dim, decoder_dim, bias=True)
        self.f_beta = nn.Linear(decoder_dim, encoder_dim)

        self.fcn = nn.Linear(decoder_dim, vocab_size)
        self.drop = nn.Dropout(drop_prob)

    def forward(self, features, captions):

        #vectorize the caption
        embeds = self.embedding(captions)

        # Initialize LSTM state
        h, c = self.init_hidden_state(features) # (batch_size, decoder_dim)

        #get the seq length to iterate
        seq_length = len(captions[0])-1 #Exclude the last one
        batch_size = captions.size(0)
        num_features = features.size(1)

        preds = torch.zeros(batch_size, seq_length, self.vocab_size).to(device)
        alphas = torch.zeros(batch_size, seq_length, num_features).to(device)

        for s in range(seq_length):
            alpha, context = self.attention(features, h)
            lstm_input = torch.cat((embeds[:, s], context), dim=1)
            h, c = self.lstm_cell(lstm_input, (h, c))

            output = self.fcn(self.drop(h))

            preds[:, s] = output
            alphas[:, s] = alpha

        return preds, alphas

    def generate_caption(self, features, max_len=20, vocab=None):
        # Inference part
        # Given the image features generate the captions

        batch_size = features.size(0)
        h, c = self.init_hidden_state(features) # (batch_size, decoder_dim)

        alphas = []

        #starting input
        word = torch.tensor(vocab.stoi['<SOS>']).view(1, -1).to(device)
        embeds = self.embedding(word)

        captions = []

        for i in range(max_len):
            alpha, context = self.attention(features, h)

            #store the apla score
            alphas.append(alpha.cpu().detach().numpy())

            lstm_input = torch.cat((embeds[:, 0], context), dim=1)
            h, c = self.lstm_cell(lstm_input, (h, c))
            output = self.fcn(self.drop(h))
            output = output.view(batch_size, -1)

```

```

        #select the word with most val
        predicted_word_idx = output.argmax(dim=1)

        #save the generated word
        captions.append(predicted_word_idx.item())

        #end if <EOS detected>
        if vocab.itos[predicted_word_idx.item()] == "<EOS>":
            break

        #send generated word as the next caption
        embeds = self.embedding(predicted_word_idx.unsqueeze(0))

    #covert the vocab idx to words and return sentence
    return [vocab.itos[idx] for idx in captions],alphas

def init_hidden_state(self, encoder_out):
    mean_encoder_out = encoder_out.mean(dim=1)
    h = self.init_h(mean_encoder_out) # (batch_size, decoder_dim)
    c = self.init_c(mean_encoder_out)
    return h, c

```

```

class EncoderDecoder(nn.Module):
    def __init__(self, embed_size, vocab_size, attention_dim, encoder_dim, decoder_dim, drop_prob=0.3):
        super().__init__()
        self.encoder = EncoderCNN()
        self.decoder = DecoderRNN(
            embed_size=embed_size,
            vocab_size = len(dataset.vocab),
            attention_dim=attention_dim,
            encoder_dim=encoder_dim,
            decoder_dim=decoder_dim
        )

    def forward(self, images, captions):
        features = self.encoder(images)
        outputs = self.decoder(features, captions)
        return outputs

```

#### 4) Setting Hyperparameter and Init the model

```

#Hyperparams
embed_size=300
vocab_size = len(dataset.vocab)
attention_dim=256
encoder_dim=2048
decoder_dim=512
learning_rate = 3e-4

```

```

#init model
model = EncoderDecoder(
    embed_size=300,
    vocab_size = len(dataset.vocab),
    attention_dim=256,
    encoder_dim=2048,
    decoder_dim=512
).to(device)

criterion = nn.CrossEntropyLoss(ignore_index=dataset.vocab.stoi["<PAD>"])
optimizer = optim.Adam(model.parameters(), lr=learning_rate)

```



```
#helper function to save the model
def save_model(model,num_epochs):
    model_state = {
        'num_epochs':num_epochs,
        'embed_size':embed_size,
        'vocab_size':len(dataset.vocab),
        'attention_dim':attention_dim,
        'encoder_dim':encoder_dim,
        'decoder_dim':decoder_dim,
        'state_dict':model.state_dict()
    }

    torch.save(model_state,'attention_model_state.pth')
```

## 5) Training Job from above configs

```
num_epochs = 10
print_every = 100

for epoch in range(1,num_epochs+1):
    for idx, (image, captions) in enumerate(iter(data_loader)):
        image,captions = image.to(device),captions.to(device)

        # Zero the gradients.
        optimizer.zero_grad()

        # Feed forward
        outputs,attentions = model(image, captions)

        # Calculate the batch loss.
        targets = captions[:,1:]
        loss = criterion(outputs.view(-1, vocab_size), targets.reshape(-1))

        # Backward pass.
        loss.backward()

        # Update the parameters in the optimizer.
        optimizer.step()

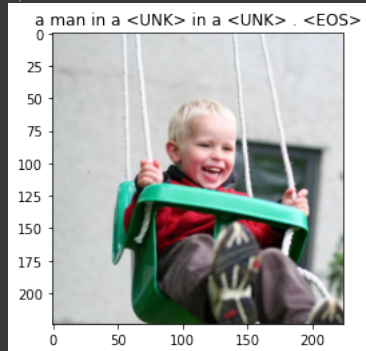
        if (idx+1)%print_every == 0:
            print("Epoch: {} loss: {:.5f}".format(epoch,loss.item()))

        #generate the caption
        model.eval()
        with torch.no_grad():
            dataiter = iter(data_loader)
            img,_ = next(dataiter)
            features = model.encoder(img[0:1].to(device))
            caps,alphas = model.decoder.generate_caption(features,vocab=dataset.vocab)
            caption = ' '.join(caps)
            show_image(img[0],title=caption)

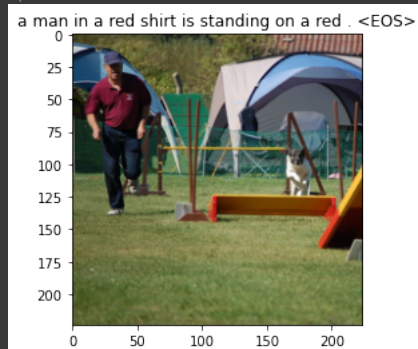
        model.train()

#save the latest model
save_model(model,epoch)
```

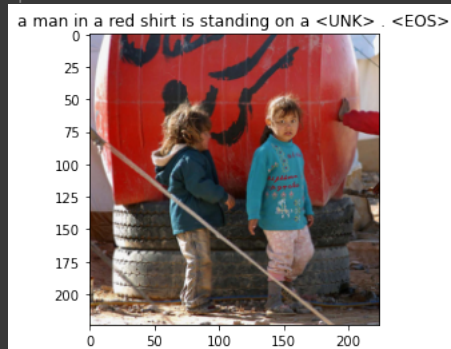
Epoch: 1 loss: 4.24822



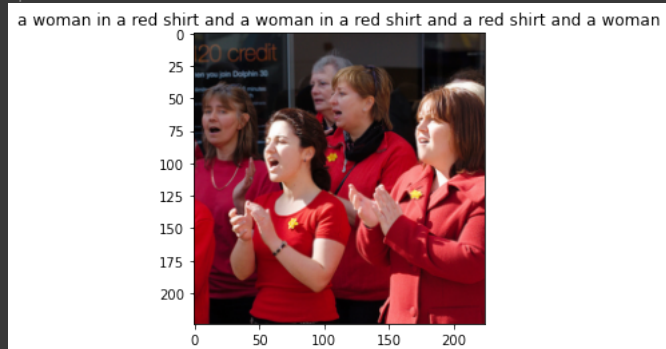
Epoch: 2 loss: 3.56758



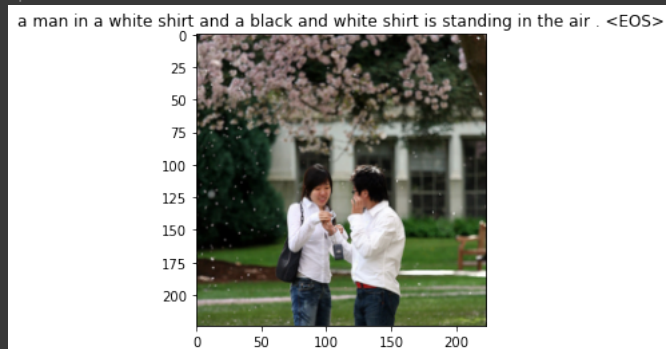
Epoch: 3 loss: 3.31058



Epoch: 4 loss: 3.07021

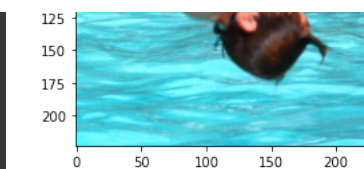


Epoch: 5 loss: 2.98379

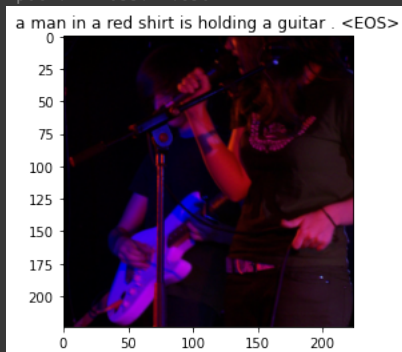


Epoch: 6 loss: 2.89330

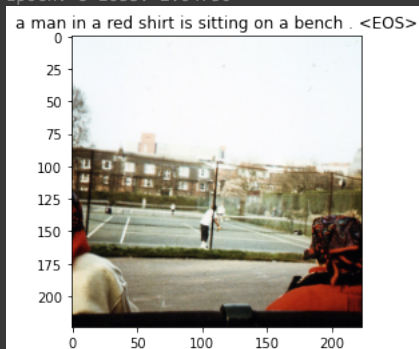




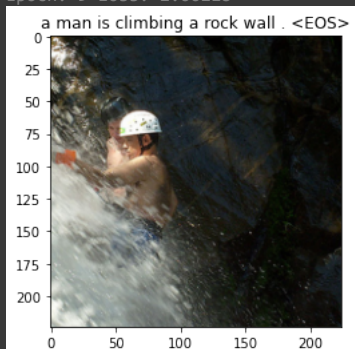
Epoch: 7 loss: 2.83097



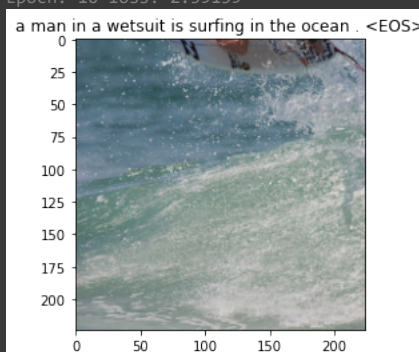
Epoch: 8 loss: 2.64736



Epoch: 9 loss: 2.66218



Epoch: 10 loss: 2.59159



## 6 Visualizing the attentions

Defining helper functions

- Given the image generate captions and attention scores
- Plot the attention scores in the image

```

#generate caption
def get_caps_from(features_tensors):
    #generate the caption
    model.eval()
    with torch.no_grad():
        features = model.encoder(features_tensors.to(device))
        caps,alphas = model.decoder.generate_caption(features,vocab=dataset.vocab)
        caption = ' '.join(caps)
        show_image(features_tensors[0],title=caption)

    return caps,alphas

#Show attention
def plot_attention(img, result, attention_plot):
    #untransform
    img[0] = img[0] * 0.229
    img[1] = img[1] * 0.224
    img[2] = img[2] * 0.225
    img[0] += 0.485
    img[1] += 0.456
    img[2] += 0.406

    img = img.numpy().transpose((1, 2, 0))
    temp_image = img

    fig = plt.figure(figsize=(15, 15))

    len_result = len(result)
    for l in range(len_result):
        temp_att = attention_plot[l].reshape(7,7)

        ax = fig.add_subplot(len_result//2,len_result//2, l+1)
        ax.set_title(result[l])
        img = ax.imshow(temp_image)
        ax.imshow(temp_att, cmap='gray', alpha=0.7, extent=img.get_extent())

    plt.tight_layout()
    plt.show()

```

```

#show any 1
dataiter = iter(data_loader)
images,_ = next(dataiter)

img = images[0].detach().clone()
img1 = images[0].detach().clone()
caps,alphas = get_caps_from(img.unsqueeze(0))

plot_attention(img1, caps, alphas)

```



```
#show any 1
dataiter = iter(data_loader)
images,_ = next(dataiter)

img = images[0].detach().clone()
img1 = images[0].detach().clone()
caps,alphas = get_caps_from(img.unsqueeze(0))

plot_attention(img1, caps, alphas)
```

