

In [6]:

```
1 %%html
2 <style>
3 table {display: block;}
4 td {
5     font-size: 18px
6 }
7 .rendered_html { font-size: 28px; }
8 *{ line-height: 200%; }
9 </style>
```

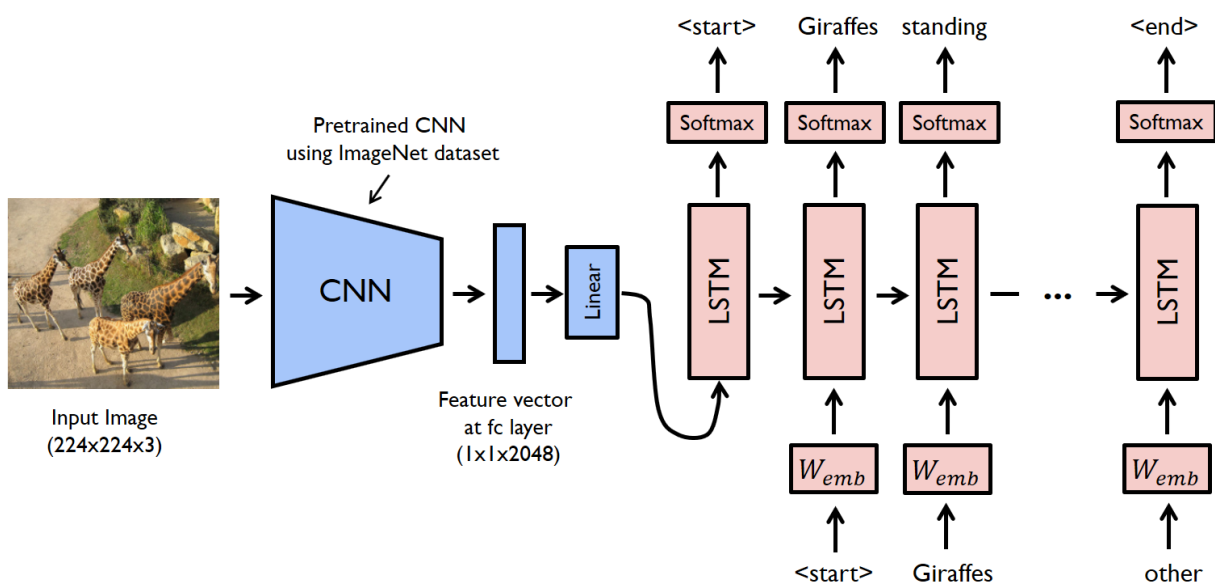
Welcome to the Natural Language Processing and the Web WS2022/23

Practical Class 7: Image Captioning

Adapted from the practice class Xintong Wang delivered last year

Introduction:

Image captioning aims to convert a given input image into a **natural language description**. The encoder-decoder framework is widely used for this task. The image encoder is a **convolutional neural network (CNN)**. The decoder is a **long short-term memory (LSTM)** network.



Training phase

For the encoder part, the pretrained CNN extracts the feature vector from a given input image. The feature vector is linearly transformed to have the same dimension as the input dimension of the LSTM network. For the decoder part, source and target texts are predefined. For example, if the image description is Giraffes standing next to each other, the source sequence is a list containing [`<start>`, Giraffes, standing, next, to, each, other] and the target sequence is a list containing [Giraffes, standing, next, to, each, other, `<end>`]. Using these source and target sequences and the feature vector, the LSTM decoder is trained as a [language model](#) conditioned on the feature vector.

Testing phase

In the test phase, the encoder part is almost same as the training phase. For the decoder part, there is a significant difference between the training phase and the test phase. In the test phase, the LSTM decoder can't see the image description. To deal with this problem, the LSTM decoder feeds back the previously generated word to the next input.!!!

Inference phase

In this phase, users could input any images and get the caption for each image using the model we already trained.

Goal for this practice class:

1. Understand the standard process to implement image captioning model.

2. Understand the pipeline for deep learning.
3. Understand the difference among training, testing, and inference.

[Assignment](#) for this class will be one question (5 scores) and two extra coding parts (10 scores).

GPU and Dataset Preparation

GPU Preparation

[Colab](#) provides you 12 hours GPU usage, 12G.

First of all, make sure this tutorial and your assignment is running on GPU environment.

Otherwise, it would be pretty slow to train you model.

If your GPU environment is ready to use, use command `!nvidia-smi` to monitor as following.

```
!nvidia-smi
```

In [2]:

```
1 !nvidia-smi
```

Sat Nov 20 23:16:13 2021

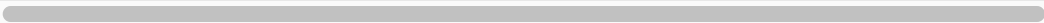
```
+-----+
+-----+
| NVIDIA-SMI 470.57.02      Driver Version: 470.57.02      CUDA Version: 1
1.4      |
+-----+-----+-----+
+-----+
| GPU  Name           Persistence-M| Bus-Id        Disp.A | Volatile Unco
rr. ECC |
| Fan  Temp   Perf   Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Com
pute M. |
|                                     |                    |
MIG M. |
+=====+=====+=====+
=====+
|   0   NVIDIA GeForce ...  On    | 00000000:04:00.0 Off |
N/A |
| 27%   25C    P8     18W / 250W |  4920MiB / 11019MiB |      0%
Default |
|                                     |                    |
N/A |
+-----+-----+-----+
+-----+
|   1   NVIDIA GeForce ...  On    | 00000000:05:00.0 Off |
N/A |
| 59%   76C    P2     226W / 250W |  9909MiB / 11178MiB |     98%
Default |
|                                     |                    |
N/A |
+-----+-----+-----+
+-----+
|   2   NVIDIA GeForce ...  On    | 00000000:08:00.0 Off |
N/A |
| 27%   27C    P8      1W / 250W |  6968MiB / 11019MiB |      0%
Default |
|                                     |                    |
N/A |
+-----+-----+-----+
+-----+
|   3   NVIDIA GeForce ...  On    | 00000000:09:00.0 Off |
N/A |
| 55%   71C    P2     264W / 250W |  9577MiB / 11178MiB |     66%
Default |
|                                     |                    |
N/A |
+-----+-----+-----+
+-----+
|   4   NVIDIA GeForce ...  On    | 00000000:84:00.0 Off |
N/A |
| 29%   25C    P8      8W / 250W |    1MiB / 11178MiB |      0%
Default |
|                                     |                    |
N/A |
+-----+-----+-----+
+-----+
|   5   NVIDIA TITAN Xp      On    | 00000000:85:00.0 Off |
```

```

N/A |
| 23%   19C   P8      8W / 250W |      787MiB / 12196MiB |      0%
Default |
|
N/A |
+-----+-----+-----+
-----+
|   6  NVIDIA GeForce ...  On   | 00000000:88:00.0 Off |
N/A |
| 29%   19C   P8      8W / 250W |      775MiB / 11178MiB |      0%
Default |
|
N/A |
+-----+-----+-----+
-----+
|   7  NVIDIA GeForce ...  On   | 00000000:89:00.0 Off |
N/A |
| 42%   79C   P2     246W / 250W |      7140MiB / 11019MiB |      96%
Default |
|
N/A |
+-----+-----+-----+
-----+

+-----+
-----+
| Processes:
|
| GPU   GI   CI           PID   Type   Process name                      GPU
Memory | ID   ID
ge      |
|=====
=====|
|    1  N/A  N/A     381721     C    python
9905MiB |
|    2  N/A  N/A     990296     C    /usr/bin/python3
6965MiB |
|    3  N/A  N/A    1079315     C    python3
9573MiB |
|    5  N/A  N/A     771734     C    ...cal/miniconda3/bin/python
783MiB  |
|    6  N/A  N/A    1148442     C    ...cal/miniconda3/bin/python
771MiB  |
|    7  N/A  N/A    2822267     C    python3
7137MiB |
+-----+
-----+

```



In this practical class, we will use PyTorch Framework.

Make sure PyTorch is installed in your environment. To

install the appropriate PyTorch, check cuda version first.

```
# !nvcc --version
```

OR in your local linux machine as:

```
!cat /usr/local/cuda/version.txt
```

In []:

```
1 !cat /usr/local/cuda/version.txt
```

This result tells us that cuda that is running on our server is version 11

But please ensure to check it yourself in your Google Colab env. Because your running environment may different from where the notebook is running here.

[PyTorch Install Link: \(https://pytorch.org/get-started/locally/\)](https://pytorch.org/get-started/locally/)

If your environment is the same with this notebook, then the install command will be:

```
!pip install torch==1.10.0+cu111  
torchvision==0.11.1+cu111 torchaudio==0.10.0+cu111 -f https://download.pytorch.org/whl/cu111/torch_stable.html
```

Sidenote: If you want to use conda package management, you should install miniconda or anaconda first. Otherwise, you should use pip.

In []:

```
1 !pip install torch==1.10.0+cu111 torchvision==0.11.1+cu111 torchaudio==0.10.0+cu111
```

How to check if pytorch is installed successfully?

```
import torch
```

If you could import it, then you successfully install it.

Note: the package of PyTorch is called torch.

In [7]:

```
1 import torch
```

Dataset Preparation:

For this practice class we will use [flickr8k](#) to train our models.

flickr 8k (~1G) is a captioning dataset. You can download it from Kaggle using the link

<https://www.kaggle.com/adityajn105/flickr8k/activity>
(<https://www.kaggle.com/adityajn105/flickr8k/activity>).

But we will share the dataset for the practice class in our datas server.

As Colab will delete all your uploading files when 12hs limit comes. Now, we use the [drive.mount](#) method to make our notebook reach our datasets without uploading.

What you should do:

Download the dataset from [here](http://ltdata1.informatik.uni-hamburg.de/flickr8k/flickr8k.zip) (<http://ltdata1.informatik.uni-hamburg.de/flickr8k/flickr8k.zip>). Then uploading it to your own google drive. Run the code below, using the token your own google drive gives back to you.

Note: Please do not share your notebook, because it would be potential risks. We don't know who can reach your files.

In []:

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

In [9]:

```
1 !pwd
```

/raid/seid/par4sem/tmp

In [8]:

```
1 !ls /content/drive/MyDrive/
```

ch7.ipynb

Now we have [flickr8k](#) folder under our directory

Get the dataset from our server

The dataset is available in our server. Download it to your colab directory as follows `wget -`

`http://ltdatal.informatik.uni-hamburg.de/flickr8k/flickr8k.zip`

In [10]:

```
1 !wget - http://ltdatal.informatik.uni-hamburg.de/flickr8k/flickr8k.zip
```

```
--2021-11-20 23:28:42-- http://-/ (http://-/)  
Resolving - (-)... failed: Name or service not known.  
wget: unable to resolve host address '-'  
  
--2021-11-20 23:28:42-- http://ltdatal.informatik.uni-hamburg.de/flic  
kr8k/flickr8k.zip (http://ltdatal.informatik.uni-hamburg.de/flickr8k/f  
lickr8k.zip)  
Resolving ltdatal.informatik.uni-hamburg.de (ltdatal.informatik.uni-ha  
mburg.de)... 134.100.15.200  
Connecting to ltdatal.informatik.uni-hamburg.de (ltdatal.informatik.un  
i-hamburg.de)|134.100.15.200|:80... connected.  
HTTP request sent, awaiting response... 200 OK  
Length: 1112971163 (1.0G) [application/zip]  
Saving to: 'flickr8k.zip'  
  
flickr8k.zip      100%[=====>]    1.04G  51.7MB/s   in  
16s  
  
2021-11-20 23:28:58 (67.2 MB/s) - 'flickr8k.zip' saved [1112971163/111  
2971163]  
  
FINISHED --2021-11-20 23:28:58--  
Total wall clock time: 16s  
Downloaded: 1 files, 1.0G in 16s (67.2 MB/s)
```

unzip the file to a folder

```
unzip flickr8k.zip -d flickr8k
```

In [11]:

```
1 !unzip flickr8k.zip -d flickr8k
```

Archive: flickr8k.zip

```

inflating: flickr8k/Images/1000268201_693b08cb0e.jpg
inflating: flickr8k/Images/1001773457_577c3a7d70.jpg
inflating: flickr8k/Images/1002674143_1b742ab4b8.jpg
inflating: flickr8k/Images/1003163366_44323f5815.jpg
inflating: flickr8k/Images/1007129816_e794419615.jpg
inflating: flickr8k/Images/1007320043_627395c3d8.jpg
inflating: flickr8k/Images/1009434119_febe49276a.jpg
inflating: flickr8k/Images/1012212859_01547e3f17.jpg
inflating: flickr8k/Images/1015118661_980735411b.jpg
inflating: flickr8k/Images/1015584366_dfcec3c85a.jpg
inflating: flickr8k/Images/101654506_8eb26cfb60.jpg
inflating: flickr8k/Images/101669240_b2d3e7f17b.jpg
inflating: flickr8k/Images/1016887272_03199f49c4.jpg
inflating: flickr8k/Images/1019077836_6fc9b15408.jpg
inflating: flickr8k/Images/1019604187_d087bf9a5f.jpg
inflating: flickr8k/Images/1020651753_06077ec457.jpg
inflating: flickr8k/Images/1022454332_6af2c1449a.jpg
inflating: flickr8k/Images/1022454428_b6b660a67b.jpg

```

move the file to your Gdrive

In []:

```
1 !mv flickr8k /content/drive/MyDrive/
```

Data Loader

There will be one folder containing all the images and one text file under your directory.

Images: 8091 pictures.

Caption: description for each image

The start point for any deep learning models is data processing. Be patience, because most of the bugs will happen in the future if you make mistakes in this step.

In Pytorch, the interface for training to use is called [DataLoader](#).

In [38]:

```
1 # Check the content
2 !ls /content/drive/MyDrive/flickr8k/
```

captions.txt Images

In [39]:

```
1 # How many images we have
2 !ls /content/drive/MyDrive/flickr8k/Images | wc -l
```

8091

In [40]:

```
1 # See how the description file looks like, image_name, description
2 !tail /content/drive/MyDrive/flickr8k/captions.txt
```

```
997338199_7343367d7f.jpg,A person stands near golden walls .
997338199_7343367d7f.jpg,a woman behind a scrolled wall is writing
997338199_7343367d7f.jpg,A woman standing near a decorated wall writes
.
997338199_7343367d7f.jpg,The walls are covered in gold and patterns .
997338199_7343367d7f.jpg,"Woman writing on a pad in room with gold , d
ecorated walls ."
997722733_0cb5439472.jpg,A man in a pink shirt climbs a rock face
997722733_0cb5439472.jpg,A man is rock climbing high in the air .
997722733_0cb5439472.jpg,A person in a red shirt climbing up a rock fa
ce covered in assist handles .
997722733_0cb5439472.jpg,A rock climber in a red shirt .
997722733_0cb5439472.jpg,A rock climber practices on a rock climbing w
all .
```


In [41]:

```
1 # lets see how the image for the file '997722733_0cb5439472.jpg' looks like, for
2 from PIL import Image
3 Image.open('/content/drive/MyDrive/flickr8k/Images/997722733_0cb5439472.jpg')
```

Out[41]:



IMPORTANT THING We should know in this step!

We want to convert text to numerical values

1. We need a Vocabulary mapping each word to a index
2. We need to setup a Pytorch dataset to load the data
3. Setup padding of every batch (all examples should be of same seq_len and setup dataloader)

We will look at `itos` --> `index_to_sentence` and `stoi`-->`sentence_to_index` here.

In [14]:

```
1 # use spaCy for preprocessing of the description text
2 !python -m spacy download en_core_web_sm
```

Collecting en-core-web-sm==3.0.0

Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.0.0/en_core_web_sm-3.0.0-py3-none-any.whl (https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.0.0/en_core_web_sm-3.0.0-py3-none-any.whl) (13.7 MB)

Requirement already satisfied: spacy<3.1.0,>=3.0.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from en-core-web-sm==3.0.0) (3.0.3)
Requirement already satisfied: wasabi<1.1.0,>=0.8.1 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (0.8.2)
Requirement already satisfied: typer<0.4.0,>=0.3.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (0.3.2)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.0.5)
Requirement already satisfied: packaging>=20.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (20.4)
Requirement already satisfied: catalogue<2.1.0,>=2.0.1 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.0.1)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.0.1)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.0.5)
Requirement already satisfied: numpy>=1.15.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.18.5)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (4.47.0)
Requirement already satisfied: blis<0.8.0,>=0.4.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (0.7.4)
Requirement already satisfied: setuptools in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (49.1.0.post20200710)
Requirement already satisfied: thinc<8.1.0,>=8.0.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (8.0.1)
Requirement already satisfied: importlib-metadata>=0.20 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.10.1)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.25.1)

Requirement already satisfied: pydantic<1.8.0,>=1.7.1 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.7.3)

Requirement already satisfied: srsly<3.0.0,>=2.4.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.4.0)

Requirement already satisfied: Jinja2 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.11.2)

Requirement already satisfied: pathy in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (0.4.0)

Requirement already satisfied: typing-extensions>=3.7.4 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.7.4.3)

Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.0.5)

Requirement already satisfied: zipp>=0.5 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from importlib-metadata>=0.20->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.1.0)

Requirement already satisfied: six in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from packaging>=20.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.15.0)

Requirement already satisfied: pyparsing>=2.0.2 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from packaging>=20.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.4.7)

Requirement already satisfied: idna<3,>=2.5 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2.10)

Requirement already satisfied: urllib3<1.27,>=1.21.1 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.25.9)

Requirement already satisfied: chardet<5,>=3.0.2 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.0.4)

Requirement already satisfied: certifi>=2017.4.17 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (2020.12.5)

Requirement already satisfied: click<7.2.0,>=7.1.1 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from typer<0.4.0,>=0.3.0->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (7.1.2)

Requirement already satisfied: MarkupSafe>=0.23 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from Jinja2->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (1.1.1)

Requirement already satisfied: smart-open<4.0.0,>=2.2.0 in /srv/home/yimam/anaconda3/lib/python3.7/site-packages (from pathy->spacy<3.1.0,>=3.0.0->en-core-web-sm==3.0.0) (3.0.0)

Installing collected packages: en-core-web-sm

Successfully installed en-core-web-sm-3.0.0

WARNING: You are using pip version 21.0.1; however, version 21.3.1 is available.

You should consider upgrading via the '/srv/home/yimam/anaconda3/bin/python -m pip install --upgrade pip' command.

✓ Download and installation successful

You can now load the package via `spacy.load('en_core_web_sm')`

In [15]:

```
1 import os # when loading file paths
2 import pandas as pd # for lookup in annotation file
3 import spacy # for tokenizer
4 import torch
5 from torch.nn.utils.rnn import pad_sequence # pad batch
6 from torch.utils.data import DataLoader, Dataset
7 from PIL import Image # Load img
8 import torchvision.transforms as transforms
9
10
11 # Download with: python -m spacy download en
12 # spacy_eng = spacy.load("en")
13
14 # if not in colab, use the following
15 # Download with: python -m spacy download en_core_web_sm
16 spacy_eng = spacy.load("en_core_web_sm")
```

1. Build vocabulary and numericalize

Tokenizer:

I attend NLP4Web Course.

--- Tokenizer ---

```
['i', 'attend', 'nlp4web', 'class']
```

Build Vocab.: Counting the frequency. If meet the [freq_threshold](#) then add to our vocabulary.

Index plus 1 (start from 4, because we have pad, sos, eos, unk tokens)

In [16]:

```
1 class Vocabulary:
2     def __init__(self, freq_threshold):
3         self.itos = {0: "<PAD>", 1: "<SOS>", 2: "<EOS>", 3: "<UNK>"}
4         self.stoi = {"<PAD>": 0, "<SOS>": 1, "<EOS>": 2, "<UNK>": 3}
5         self.freq_threshold = freq_threshold
6
7     def __len__(self):
8         return len(self.itos)
9
10    @staticmethod
11    def tokenizer_eng(text):
12        return [tok.text.lower() for tok in spacy_eng.tokenizer(text)]
13
14    def build_vocabulary(self, sentence_list):
15        frequencies = {}
16        idx = 4
17
18        for sentence in sentence_list:
19            for word in self.tokenizer_eng(sentence):
20                if word not in frequencies:
21                    frequencies[word] = 1
22
23                else:
24                    frequencies[word] += 1
25
26                if frequencies[word] == self.freq_threshold:
27                    self.stoi[word] = idx
28                    self.itos[idx] = word
29                    idx += 1
30
31    def numericalize(self, text):
32        tokenized_text = self.tokenizer_eng(text)
33
34        return [
35            self.stoi[token] if token in self.stoi else self.stoi["<UNK>"]
```

```
36         for token in tokenized_text
37     ]
```

In [43]:

```
1 #spaCy based tokenizer
2 examplesent = ["the old fox jumps over the lazy dog."]
3 spacy_eng.tokenizer(examplesent[0])
```

Out[43]:

the old fox jumps over the lazy dog.

In [44]:

```
1 vocab = Vocabulary(1)
2 vocab.build_vocabulary(examplesent)
3 print("index to sentence==>", vocab.itos)
4 print("sentence to index==>", vocab.stoi)
```

```
index to sentence==> {0: '<PAD>', 1: '<SOS>', 2: '<EOS>', 3: '<UNK>',
4: 'the', 5: 'old', 6: 'fox', 7: 'jumps', 8: 'over', 9: 'lazy', 10: 'd
og', 11: '.'}
sentence to index==> {'<PAD>': 0, '<SOS>': 1, '<EOS>': 2, '<UNK>': 3,
'the': 4, 'old': 5, 'fox': 6, 'jumps': 7, 'over': 8, 'lazy': 9, 'dog':
10, '.': 11}
```

In [45]:

```
1 vocab.stoi["<PAD>"], vocab.stoi["<SOS>"], vocab.stoi["<EOS>"], vocab.stoi["<UNK>"]
```

Out[45]:

(0, 1, 2, 3)

In [46]:

```
1 vocab.numericalize("fox and diffirent dog")
```

Out[46]:

[6, 3, 3, 10]

2. Make Flickr Dataset/DataLoader, image and caption

Image part is more easy than the caption part

In [17]:

```
1 class FlickrDataset(Dataset):
2     def __init__(self, root_dir, captions_file, transform=None, freq_threshold=5
3         self.root_dir = root_dir
4         self.df = pd.read_csv(captions_file)
5         self.transform = transform
6
7         # Get img, caption columns
8         self.imgs = self.df["image"]
9         self.captions = self.df["caption"]
10
11        # Initialize vocabulary and build vocab
12        self.vocab = Vocabulary(freq_threshold)
13        self.vocab.build_vocabulary(self.captions.tolist())
14
15    def __len__(self):
16        return len(self.df)
17
18    def __getitem__(self, index):
19        caption = self.captions[index]
20        img_id = self.imgs[index]
21        img = Image.open(os.path.join(self.root_dir, img_id)).convert("RGB")
22
23        if self.transform is not None:
24            img = self.transform(img)
25
26        numericalized_caption = [self.vocab.stoi["<SOS>"]]
27        numericalized_caption += self.vocab.numericalize(caption)
28        numericalized_caption.append(self.vocab.stoi["<EOS>"])
29
30        return img, torch.tensor(numericalized_caption)
```

We want to load data in the way of batch not just a single pair.

In [18]:

```
1 class MyCollate:
2     def __init__(self, pad_idx):
3         self.pad_idx = pad_idx
4
5     def __call__(self, batch):
6         imgs = [item[0].unsqueeze(0) for item in batch]
7         imgs = torch.cat(imgs, dim=0)
8         targets = [item[1] for item in batch]
9         targets = pad_sequence(targets, batch_first=False, padding_value=self.pa
10
11         return imgs, targets
```

Dataloader is the interface for PyTorch to reach data after processing

In [19]:

```
1 def get_loader(root_folder, annotation_file, transform, batch_size=32, num_workers=
2
3     dataset = FlickrDataset(root_folder, annotation_file, transform=transform)
4
5     pad_idx = dataset.vocab.stoi[ "<PAD>" ]
6
7     loader = DataLoader(
8         dataset=dataset,
9         batch_size=batch_size,
10        num_workers=num_workers,
11        shuffle=shuffle,
12        pin_memory=pin_memory,
13        collate_fn=MyCollate(pad_idx=pad_idx),
14    )
15
16    return loader, dataset
```

In [20]:

```
1 transform = transforms.Compose(
2     [transforms.Resize((224, 224)),
3     transforms.ToTensor(),])
```

In [24]:

```
1 loader, dataset = get_loader(  
2     #  "./flickr8k/Images/", "./flickr8k/captions.txt", transform=transform  
3     "/content/drive/MyDrive/flickr8k/Images/", /content/drive/MyDrive/flickr8k/  
4  
5     )  
6  
7 for idx, (imgs, captions) in enumerate(loader):  
8     print(imgs.shape)  
9     print(captions.shape)  
10    break
```

```
torch.Size([32, 3, 224, 224])  
torch.Size([20, 32])
```

One more thing we should finish here before we start our model part.

In order to have a look at the current performance, we sample some cases from our dataset and generate related captions comparing with the ground truth.

What you should do: (Assignment 2)

- 1. Make a test_examples folder and put five images/rename in this folder**
- 2. Copy each caption to the code below**

In [49]:

```
1 import torch
2 import torchvision.transforms as transforms
3 from PIL import Image
```

In [50]:

```
1 def print_examples(model, device, dataset):
2     transform = transforms.Compose(
3         [
4             transforms.Resize((299, 299)),
5             transforms.ToTensor(),
6             transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
7         ]
8     )
9
10    model.eval()
11    test_img1 = transform(Image.open("test_examples/dog.jpg").convert("RGB")).ur
12        0
13    )
14    print("Example 1 CORRECT: Dog on a beach by the ocean")
15    print(
16        "Example 1 OUTPUT: "
17        + " ".join(model.caption_image(test_img1.to(device), dataset.vocab))
18    )
19    test_img2 = transform(
20        Image.open("test_examples/child.jpg").convert("RGB")
21    ).unsqueeze(0)
22    print("Example 2 CORRECT: Child holding red frisbee outdoors")
23    print(
24        "Example 2 OUTPUT: "
25        + " ".join(model.caption_image(test_img2.to(device), dataset.vocab))
26    )
27    test_img3 = transform(Image.open("test_examples/bus.png").convert("RGB")).ur
28        0
29    )
30    print("Example 3 CORRECT: Bus driving by parked cars")
31    print(
32        "Example 3 OUTPUT: "
33        + " ".join(model.caption_image(test_img3.to(device), dataset.vocab))
34    )
35    test_img4 = transform(
```

```

36         Image.open("test_examples/boat.png").convert("RGB")
37     ).unsqueeze(0)
38     print("Example 4 CORRECT: A small boat in the ocean")
39     print(
40         "Example 4 OUTPUT: "
41         + " ".join(model.caption_image(test_img4.to(device), dataset.vocab))
42     )
43     test_img5 = transform(
44         Image.open("test_examples/horse.png").convert("RGB")
45     ).unsqueeze(0)
46     print("Example 5 CORRECT: A cowboy riding a horse in the desert")
47     print(
48         "Example 5 OUTPUT: "
49         + " ".join(model.caption_image(test_img5.to(device), dataset.vocab))
50     )
51     model.train()

```

Save and loading checkpoints

In [51]:

```

1  def save_checkpoint(state, filename="my_checkpoint.pth.tar"):
2      print("=> Saving checkpoint")
3      torch.save(state, filename)
4
5
6  def load_checkpoint(checkpoint, model, optimizer):
7      print("=> Loading checkpoint")
8      model.load_state_dict(checkpoint["state_dict"])
9      optimizer.load_state_dict(checkpoint["optimizer"])
10     step = checkpoint["step"]
11     return step

```

Model

In [52]:

```
1 import torch
2 import torch.nn as nn
3 import statistics
4 import torchvision.models as models
```

For the image understanding part, we name it as [EncoderCNN](#).

Here we use pretrained inception model, also you could change it to VGG or ResNet.

Two parameters: [embed_size](#), [train_CNN](#)

In [53]:

```
1 class EncoderCNN(nn.Module):
2     def __init__(self, embed_size, train_CNN=False):
3         super(EncoderCNN, self).__init__()
4         self.train_CNN = train_CNN
5         self.inception = models.inception_v3(pretrained=True, aux_logits=False)
6         self.inception.fc = nn.Linear(self.inception.fc.in_features, embed_size)
7         self.relu = nn.ReLU()
8         self.times = []
9         self.dropout = nn.Dropout(0.5)
10
11     def forward(self, images):
12         features = self.inception(images)
13         return self.dropout(self.relu(features))
```

For the caption generation part, we name it as
DecoderRNN.

Here we use LSTM

Three parameters: **vocab_size**, **embed_size**,
num_layers.

In [54]:

```
1 class DecoderRNN(nn.Module):
2     def __init__(self, embed_size, hidden_size, vocab_size, num_layers):
3         super(DecoderRNN, self).__init__()
4         self.embed = nn.Embedding(vocab_size, embed_size)
5         self.lstm = nn.LSTM(embed_size, hidden_size, num_layers)
6         self.linear = nn.Linear(hidden_size, vocab_size)
7         self.dropout = nn.Dropout(0.5)
8
9     def forward(self, features, captions):
10         embeddings = self.dropout(self.embed(captions))
11         embeddings = torch.cat((features.unsqueeze(0), embeddings), dim=0)
12         hiddens, _ = self.lstm(embeddings)
13         outputs = self.linear(hiddens)
14         return outputs
```

Now we have encoder and decoder, what we would do next is to obtain them in a way of end to end learning.

Pay attention to caption_image function:

```
max_length, '<EOS>', itos
```

In [55]:

```
1 class CNNtoRNN(nn.Module):
2     def __init__(self, embed_size, hidden_size, vocab_size, num_layers):
3         super(CNNtoRNN, self).__init__()
4         self.encoderCNN = EncoderCNN(embed_size)
5         self.decoderRNN = DecoderRNN(embed_size, hidden_size, vocab_size, num_la
6
7     def forward(self, images, captions):
8         features = self.encoderCNN(images)
9         outputs = self.decoderRNN(features, captions)
10        return outputs
11
12    def caption_image(self, image, vocabulary, max_length=50):
13        result_caption = []
14
15        with torch.no_grad():
16            x = self.encoderCNN(image).unsqueeze(0)
17            states = None
18
19            for _ in range(max_length):
20                hiddens, states = self.decoderRNN.lstm(x, states)
21                output = self.decoderRNN.linear(hiddens.squeeze(0))
22                predicted = output.argmax(1)
23                result_caption.append(predicted.item())
24                x = self.decoderRNN.embed(predicted).unsqueeze(0)
25
26                if vocabulary.itos[predicted.item()] == "<EOS>":
27                    break
28
29        return [vocabulary.itos[idx] for idx in result_caption]
```

Training

In [56]:

```
1 import torch
2 from tqdm import tqdm
3 import torch.nn as nn
4 import torch.optim as optim
5 import torchvision.transforms as transforms
6 from torch.utils.tensorboard import SummaryWriter
```

First of all, we need creat a model and set parameters

In [57]:

```
1 transform = transforms.Compose(
2     [
3         transforms.Resize((356, 356)),
4         transforms.RandomCrop((299, 299)),
5         transforms.ToTensor(),
6         transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
7     ]
8 )
9
10 train_loader, dataset = get_loader(
11     # root_folder="./flickr8k/Images",
12     # annotation_file="./flickr8k/captions.txt",
13     root_folder="/content/drive/MyDrive/flickr8k",
14     annotation_file="/content/drive/MyDrive/flickr8k/captions.txt",
15     transform=transform,
16     num_workers=2,
17 )
18
19 torch.backends.cudnn.benchmark = True
20 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
21 load_model = False
22 save_model = False
23 train_CNN = False
24
25 # Hyperparameters
26 embed_size = 256
27 hidden_size = 256
28 vocab_size = len(dataset.vocab)
29 num_layers = 1
30 learning_rate = 3e-4
31 num_epochs = 100
32
33 # for tensorboard
34 writer = SummaryWriter("runs/flickr")
35 step = 0
```

```

36
37 # initialize model, loss etc
38 model = CNNtoRNN(embed_size, hidden_size, vocab_size, num_layers).to(device)
39 criterion = nn.CrossEntropyLoss(ignore_index=dataset.vocab.stoi[ "<PAD>"])
40 optimizer = optim.Adam(model.parameters(), lr=learning_rate)
41
42 # Only finetune the CNN
43 for name, param in model.encoderCNN.inception.named_parameters():
44     if "fc.weight" in name or "fc.bias" in name:
45         param.requires_grad = True
46     else:
47         param.requires_grad = train_CNN
48
49 if load_model:
50     step = load_checkpoint(torch.load("my_checkpoint.pth.tar"), model, optimizer)
51
52 model.train()

```

Out[57]:

```

CNNtoRNN(
  (encoderCNN): EncoderCNN(
    (inception): Inception3(
      (Conv2d_1a_3x3): BasicConv2d(
        (conv): Conv2d(3, 32, kernel_size=(3, 3), stride=(2, 2), bias
=False)
        (bn): BatchNorm2d(32, eps=0.001, momentum=0.1, affine=True, t
rack_running_stats=True)
      )
      (Conv2d_2a_3x3): BasicConv2d(
        (conv): Conv2d(32, 32, kernel_size=(3, 3), stride=(1, 1), bia
s=False)
        (bn): BatchNorm2d(32, eps=0.001, momentum=0.1, affine=True, t
rack_running_stats=True)
      )
      (Conv2d_2b_3x3): BasicConv2d(
        (conv): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), pad

```

TRAINING: Let see what happens!

In []:

```
1 for epoch in range(num_epochs):
2     # Uncomment the line below to see a couple of test cases
3     # print_examples(model, device, dataset)
4
5     if save_model:
6         checkpoint = {
7             "state_dict": model.state_dict(),
8             "optimizer": optimizer.state_dict(),
9             "step": step,
10        }
11        save_checkpoint(checkpoint)
12
13    for idx, (imgs, captions) in tqdm(
14        enumerate(train_loader), total=len(train_loader), leave=False
15    ):
16        imgs = imgs.to(device)
17        captions = captions.to(device)
18
19        outputs = model(imgs, captions[:-1])
20        loss = criterion(
21            outputs.reshape(-1, outputs.shape[2]), captions.reshape(-1)
22        )
23
24        writer.add_scalar("Training loss", loss.item(), global_step=step)
25        step += 1
26
27        optimizer.zero_grad()
28        loss.backward(loss)
29        optimizer.step()
```

45% | ██████████ | 564/1265 [01:38<02:00, 5.81it/s]

Assignment: (15 points)

1. Summary the difference of LSTM decoder in the training and test phrases? (2 Points)
2. Coding1: Prepare test samples in order to see the performance during training. (3 Points)
3. Coding2: Write the inference code.(10 Points)

Input any image, using our trained model to generate sentence for the image.

Question 1: 2 Pts

Question 2: 3 Pts

In []:

```
1 # Question 3: 10Pts Write you code here
```

In []:

```
1
```