

## Video Caption Generation

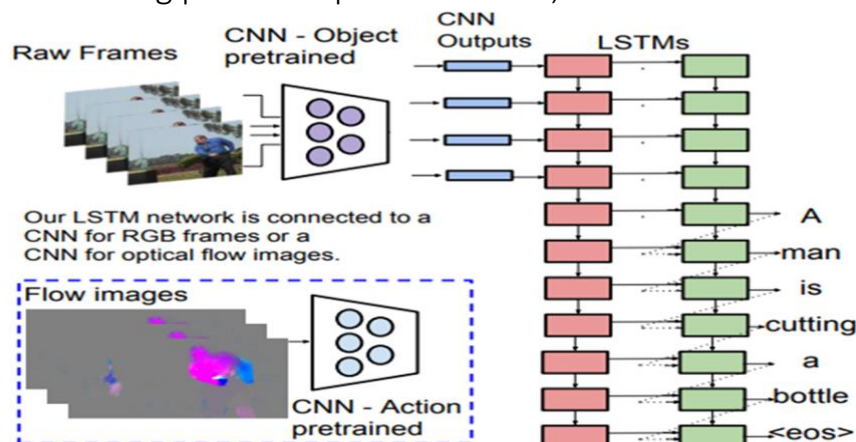
The basic idea of Video Caption Generation is to input a video and get a stream of captions for the actions happening in the video. This is achieved by following Deep learning techniques: Sequence to sequence model and training of the dataset provided.

### Requirements:

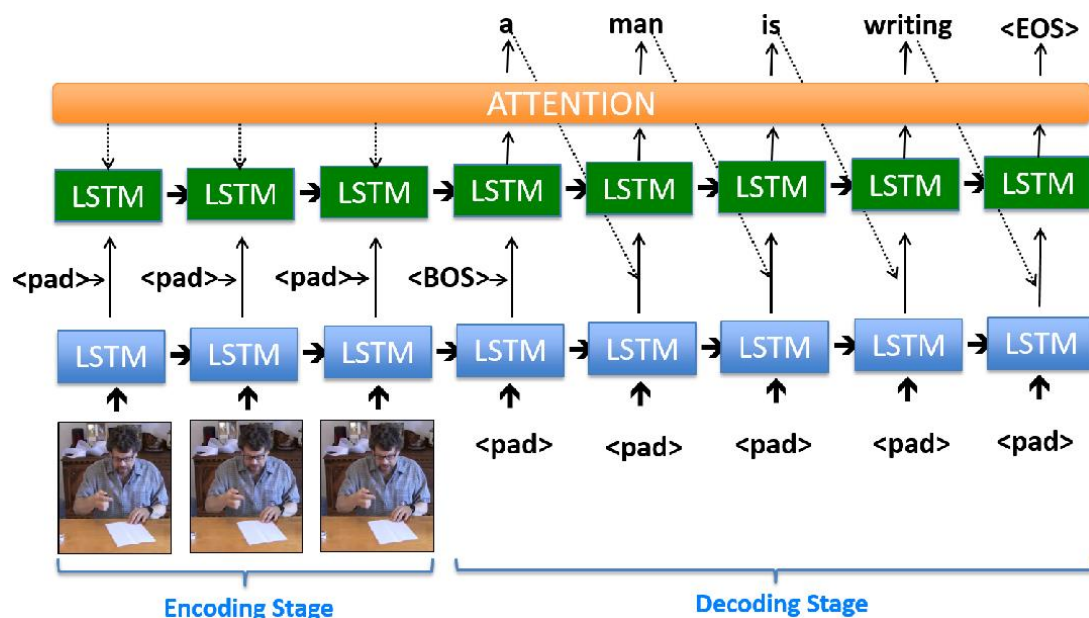
- Python 3
- Tensorflow 1.15
- MSVD Dataset (1450 videos for training, 100 for testing)
- Sequence to sequence model

### Processing of Video to Captions:

- The model uses 2 layers of LSTM (RNNs), in the first layer of the RNN, the Video is encoded, and an output is generated with the help of the decoder.
- In the decoding process, tokens are created to segment the captions based on the beginning and ending verse and perform processing over the video to produce the actual words based on the pictorial representation using a library.
- As part of sequence to sequence model, use of S2VT is to map the video frames to the words.
- The following picture depicts the same,



- While mapping the words and creating a caption, it is made sure that there is no anomaly in generation with the help of beam search which helps lean towards one path while generating.
- To reduce the exposure issues, it is important to give groundtruth as feed while processing at different input odds.
- We can reach the baseline of the model by computing the BLEU value which is based on the length and occurrences of the captions. BLEU score will unlock the efficiency at which the captions are generated.
- Attention mechanism can help enable sequence to sequence model to function effectively. Without the attention mechanism, BLEU score is 0.9854, with the attention mechanism, the BLEU score is 0.7854.



## Experiment and Results:

- The experimental parameters are the following for one of the scenarios:  
 rnn\_size = 1089  
 num\_of\_layers=2 (LSTM)  
 dim\_video\_feat = 2963  
 embed\_size = 5698  
 lr\_rate = 0.001 (learning rate)  
 batch\_size = 100  
 max\_gradnorm = 6 (gradient norm)  
 max\_encoder\_steps = 125  
 max\_decoder\_steps = 26  
 sample\_size = 1090

dim\_video\_frame = 100

Method	BLEU Score
Without Attention	0.9855
Attention with scale	0.7584
Attention with gradient norm	0.7100

- One of the results were following:

<BOS> a man

<BOS> a man is

<BOS> a man is

<BOS> a man is feeding

<BOS> a man is feeding a

<BOS> a man is feeding a dog

<BOS> a man is feeding a dog.

### Second data testing:

rnn\_size = 1024

num\_of\_layers=2 (LSTM)

dim\_video\_feat = 1859

embed\_size = 2564

lr\_rate = 0.0001 (learning rate)

batch\_size = 85

max\_gradnorm = 4 (gradient norm)

max\_encoder\_steps = 98

max\_decoder\_steps = 20

sample\_size = 985

dim\_video\_frame = 85

Method	BLEU Score
Without Attention	0.7026
Attention with scale	0.6127
Attention with gradient norm	0.6917

### Third data testing:

rnn\_size = 985  
 num\_of\_layers=2 (LSTM)  
 dim\_video\_feat = 1045  
 embed\_size = 2012  
 lr\_rate = 0.0001 (learning rate)  
 batch\_size = 50  
 max\_gradnorm = 5 (gradient norm)  
 max\_encoder\_steps = 70  
 max\_decoder\_steps = 18  
 sample\_size = 1852  
 dim\_video\_frame = 62

Method	BLEU Score
Without Attention	0.6924
Attention with scale	0.5988
Attention with gradient norm	0.7154